

PRINCIPLES OF MICROECONOMICS



Douglas Curtis and Ian Irvine
Trent University & Concordia University

Trent University & Concordia University
Principles of Microeconomics (Curtis and
Irvine)

Douglas Curtis and Ian Irvine

This text is disseminated via the Open Education Resource (OER) LibreTexts Project (<https://LibreTexts.org>) and like the hundreds of other texts available within this powerful platform, it is freely available for reading, printing and "consuming." Most, but not all, pages in the library have licenses that may allow individuals to make changes, save, and print this book. Carefully consult the applicable license(s) before pursuing such effects.

Instructors can adopt existing LibreTexts texts or Remix them to quickly build course-specific resources to meet the needs of their students. Unlike traditional textbooks, LibreTexts' web based origins allow powerful integration of advanced features and new technologies to support learning.



The LibreTexts mission is to unite students, faculty and scholars in a cooperative effort to develop an easy-to-use online platform for the construction, customization, and dissemination of OER content to reduce the burdens of unreasonable textbook costs to our students and society. The LibreTexts project is a multi-institutional collaborative venture to develop the next generation of open-access texts to improve postsecondary education at all levels of higher learning by developing an Open Access Resource environment. The project currently consists of 14 independently operating and interconnected libraries that are constantly being optimized by students, faculty, and outside experts to supplant conventional paper-based books. These free textbook alternatives are organized within a central environment that is both vertically (from advance to basic level) and horizontally (across different fields) integrated.

The LibreTexts libraries are Powered by [NICE CXOne](#) and are supported by the Department of Education Open Textbook Pilot Project, the UC Davis Office of the Provost, the UC Davis Library, the California State University Affordable Learning Solutions Program, and Merlot. This material is based upon work supported by the National Science Foundation under Grant No. 1246120, 1525057, and 1413739.

Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation nor the US Department of Education.

Have questions or comments? For information about adoptions or adaptations contact info@LibreTexts.org. More information on our activities can be found via Facebook (<https://facebook.com/Libretexts>), Twitter (<https://twitter.com/libretexts>), or our blog (<http://Blog.Libretexts.org>).

This text was compiled on 04/15/2025

TABLE OF CONTENTS

Licensing

Unit 1: The Building Blocks

- 1: Introduction to key ideas
 - 1.1: What's it all about?
 - 1.2: Understanding through the use of models
 - 1.3: Opportunity cost and the market
 - 1.4: A model of exchange and specialization
 - 1.5: Economy-wide production possibilities
 - 1.6: New Page
 - 1.7: Conclusion
 - 1.8: Key Terms
 - 1.9: Exercises for Chapter 1
- 2: Theories, data and beliefs
 - 2.1: Data analysis
 - 2.2: Data, theory and economic models
 - 2.3: Ethics, efficiency and beliefs
 - 2.4: Key Terms
 - 2.5: Exercises for Chapter 2
- 3: The classical marketplace – demand and supply
 - 3.1: The marketplace - trading
 - 3.2: The market's building blocks
 - 3.3: Demand and supply curves
 - 3.4: Non-price influences on demand
 - 3.5: Non-price influences on supply
 - 3.6: Simultaneous supply and demand impacts
 - 3.7: Market interventions - governments and interest groups
 - 3.8: Individual and market functions
 - 3.9: Useful techniques - demand and supply equations
 - 3.10: Conclusion
 - 3.11: Key Terms
 - 3.12: Exercises for Chapter 3

Unit 2: Responsiveness and the Value of Markets

- 4: Measures of response- Elasticities
 - 4.1: Price responsiveness of demand
 - 4.2: Price elasticities and public policy
 - 4.3: The time horizon and inflation
 - 4.4: Cross-price elasticities - cable or satellite
 - 4.5: The income elasticity of demand
 - 4.6: Elasticity of supply
 - 4.7: Elasticities and tax incidence
 - 4.8: Technical tricks with elasticities
 - 4.9: Key Terms
 - 4.10: Exercises for Chapter 4
- 5: Welfare economics and externalities

- 5.1: Equity and efficiency
- 5.2: Consumer and producer surplus
- 5.3: Efficient market outcomes
- 5.4: Taxation, surplus and efficiency
- 5.5: Market failures - externalities
- 5.6: Other market failures
- 5.7: Environmental policy and climate change
- 5.8: Conclusion
- 5.9: Key Terms
- 5.10: Exercises for Chapter 5

Unit 3: Decision Making by Consumer and Producers

- 6: Individual choice
 - 6.1: Rationality
 - 6.2: Choice with measurable utility
 - 6.3: Choice with ordinal utility
 - 6.4: Applications of indifference analysis
 - 6.5: Key Terms
 - 6.6: Exercises for Chapter 6
- 7: Firms, investors and capital markets
 - 7.1: Business organization
 - 7.2: Profit
 - 7.3: Risk and the investor
 - 7.4: Risk pooling and diversification
 - 7.5: Conclusion
 - 7.6: Key Terms
 - 7.7: Exercises for Chapter 7
- 8: Production and cost
 - 8.1: Efficient production
 - 8.2: The time frame
 - 8.3: Production in the short run
 - 8.4: Costs in the short run
 - 8.5: Fixed costs and sunk costs
 - 8.6: Long-run production and costs
 - 8.7: Technological change- globalization and localization
 - 8.8: Clusters, learning by doing, scope economics
 - 8.9: Conclusion
 - 8.10: Key Terms
 - 8.11: Exercises for Chapter 8

Unit 4: Market Structures

- 9: Perfect competition
 - 9.1: The perfect competition paradigm
 - 9.2: Market characteristics
 - 9.3: The firm's supply decision
 - 9.4: Dynamics- Entry and exit
 - 9.5: Long-run industry supply
 - 9.6: Globalization and technological change
 - 9.7: Efficient resource allocation
 - 9.8: Key Terms

- 9.9: Exercises for Chapter 9
- 10: Monopoly
 - 10.1: Monopolies
 - 10.2: Profit maximizing behaviour
 - 10.3: Long-run choices
 - 10.4: Output inefficiency
 - 10.5: Price discrimination
 - 10.6: Cartels- Acting like a monopolist
 - 10.7: Invention, innovation and rent seeking
 - 10.8: Conclusion
 - 10.9: Key Terms
 - 10.10: Exercises for Chapter 10
- 11: Imperfect competition
 - 11.1: The principle ideas
 - 11.2: Imperfect competitors
 - 11.3: Imperfect competitors- measures of structure and market power
 - 11.4: Imperfect competition- monopolistic competition
 - 11.5: Imperfect competition- economies of scope and platforms
 - 11.6: Strategic behaviour- Oligopoly and games
 - 11.7: Strategic behaviour- Duopoly and Cournot games
 - 11.8: Strategic behaviour- Entry, exit and potential competition
 - 11.9: Matching markets- design
 - 11.10: Conclusion
 - 11.11: Key Terms
 - 11.12: Exercises for Chapter 11

Unit 5: The Factors of Production

- 12: Labour and capital
 - 12.1: Labour - a derived demand
 - 12.2: The supply of labour
 - 12.3: Labour market equilibrium and mobility
 - 12.4: Capital - concepts
 - 12.5: The capital market
 - 12.6: Land
 - 12.7: Key Terms
 - 12.8: Exercises for Chapter 12
- 13: Human capital and the income distribution
 - 13.1: Human capital
 - 13.2: Productivity and education
 - 13.3: On-the-job training
 - 13.4: Education as signalling
 - 13.5: Education returns and quality
 - 13.6: Discrimination
 - 13.7: The income distribution
 - 13.8: Wealth and capitalism
 - 13.9: Key Terms
 - 13.10: Exercises for Chapter 13

Unit 6: Government and Trade

- 14: Government
 - 14.1: Market Failure
 - 14.2: Fiscal federalism- Taxing and spending
 - 14.3: Federal-provincial fiscal relations
 - 14.4: Government-to-individual transfers
 - 14.5: Regulation and competition policy
 - 14.6: Key Terms
 - 14.7: Exercises for Chapter 14
- 15: International trade
 - 15.1: Trade in our daily lives
 - 15.2: Canada in the world economy
 - 15.3: The gains from trade- Comparative advantage
 - 15.4: Returns to scale and dynamic gains from trade
 - 15.5: Trade barriers- Tariffs, subsidies and quotas
 - 15.6: The politics of protection
 - 15.7: Institutions governing trade
 - 15.8: Key Terms
 - 15.9: Exercises for Chapter 15

Index

Solutions To Exercises

- 1.1: Chapter 1 Solutions
- 1.2: Chapter 2 Solutions
- 1.3: Chapter 3 Solutions
- 1.4: Chapter 4 Solutions
- 1.5: Chapter 5 Solutions
- 1.6: Chapter 6 Solutions
- 1.7: Chapter 7 Solutions
- 1.8: Chapter 8 Solutions
- 1.9: Chapter 9 Solutions
- 1.10: Chapter 10 Solutions
- 1.11: Chapter 11 Solutions
- 1.12: Chapter 12 Solutions
- 1.13: Chapter 13 Solutions
- 1.14: Chapter 14 Solutions
- 1.15: Chapter 15 Solutions

Detailed Licensing

Licensing

A detailed breakdown of this resource's licensing can be found in [Back Matter/Detailed Licensing](#).

SECTION OVERVIEW

Unit 1: The Building Blocks

Economics is a social science; it analyzes human interactions in a scientific manner. We begin by defining the central aspects of this social science – trading, the marketplace, opportunity cost and resources. We explore how producers and consumers interact in society. Trade is central to improving the living standards of individuals. This material forms the subject matter of Chapter 1.

Methods of analysis are central to any science. Consequently we explore how data can be displayed and analyzed in order to better understand the economy around us in Chapter 2. Understanding the world is facilitated by the development of theories and models and then testing such theories with the use of data-driven models.

Trade is critical to individual well-being, whether domestically or internationally. To understand this trading process we analyze the behaviour of suppliers and buyers in the marketplace. Markets are formed by suppliers and demanders coming together for the purpose of trading. Thus, demand and supply are examined in Chapter 3 in tabular, graphical and mathematical form.

1: Introduction to key ideas

- 1.1: What's it all about?
- 1.2: Understanding through the use of models
- 1.3: Opportunity cost and the market
- 1.4: A model of exchange and specialization
- 1.5: Economy-wide production possibilities
- 1.6: New Page
- 1.7: Conclusion
- 1.8: Key Terms
- 1.9: Exercises for Chapter 1

2: Theories, data and beliefs

- 2.1: Data analysis
- 2.2: Data, theory and economic models
- 2.3: Ethics, efficiency and beliefs
- 2.4: Key Terms
- 2.5: Exercises for Chapter 2

3: The classical marketplace – demand and supply

- 3.1: The marketplace - trading
- 3.2: The market's building blocks
- 3.3: Demand and supply curves
- 3.4: Non-price influences on demand
- 3.5: Non-price influences on supply
- 3.6: Simultaneous supply and demand impacts
- 3.7: Market interventions - governments and interest groups
- 3.8: Individual and market functions
- 3.9: Useful techniques - demand and supply equations
- 3.10: Conclusion
- 3.11: Key Terms
- 3.12: Exercises for Chapter 3

This page titled [Unit 1: The Building Blocks](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine](#) (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1: Introduction to key ideas

Chapter 1: Introduction to key ideas

In this chapter we will explore:

1.1	What it's all about
1.2	Understanding through the use of models
1.3	Opportunity cost and the market
1.4	A model of exchange and specialization
1.5	Production possibilities for the economy
1.6	Aggregate output, growth and cycles

1.1 What's it all about?

The big issues

Economics is the study of human behaviour. Since it uses scientific methods it is called a social science. We study human behaviour to better understand and improve our world. During his acceptance speech, a recent Nobel Laureate in Economics suggested:

Economics, at its best, is a set of ideas and methods for the improvement of society. It is not, as so often seems the case today, a set of ideological rules for asserting why we cannot face the challenges of stagnation, job loss and widening inequality.

Christopher Sims, Nobel Laureate in Economics 2011

This is an elegant definition of economics and serves as a timely caution about the perils of ideology. Economics evolves continuously as current observations and experience provide new evidence about economic behaviour and relationships. Inference and policy recommendations based on earlier theories, observations and institutional structures require constant analysis and updating if they are to furnish valuable responses to changing conditions and problems.

Much of today's developed world still faces severe challenges as a result of the financial crisis that began in 2008. Unemployment rates among young people are at historically high levels in several economies, many government balance sheets are in disarray, and inequality is on the rise. In addition to the challenges posed by this severe economic cycle, the world simultaneously faces structural upheaval: Overpopulation, climate change, political instability and globalization challenge us to understand and modify our behaviour.

These challenges do not imply that our world is deteriorating. Literacy rates have been rising dramatically in the developing world for decades; child mortality has plummeted; family size is a fraction of what it was 50 years ago; prosperity is on the rise in much of Asia; life expectancy is increasing universally and deaths through wars are in a state of long-term decline.

These developments, good and bad, have a universal character and affect billions of individuals. They involve an understanding of economies as large organisms with interactive components.

Aggregate output in a national economy

A national economy is a complete multi-sector system, made up of household, business, financial, government and international sectors. Each of these sectors is an aggregate or sum of a many smaller economic units with very similar characteristics. The government sector, for example aggregates the taxing and spending activities of local, provincial and national governments. Similarly, the household sector is an aggregate of the income, spending and saving of all households but not the specifics of each individual household. Economic activity within any one of these sectors reflects, in part, the conditions and choices made in that sector. But it is also affects and is affected by conditions and actions in the other sectors. These *interactions* and *feedbacks* within the system mean that the workings of the macro-economy are more complex than the operation of the sum of its parts.

Macroeconomics: the study of the economy as a system in which interactions and feedbacks among sectors determine national output, employment and prices.

For example, consider a simple economy with just household, business and financial sectors. The household sector earns income by providing labour to the other sectors. Households make choices about spending or saving this income. Businesses make decisions about the sizes of their establishments, their labour forces, and their outputs of goods and services. The financial sector provides banking services: bank deposits, loans, and the payments system used by all three sectors.

Suppose households decide to spend more on goods and services and save less. That decision by itself does not change household sector income, but it does increase business sector sales and revenues. It also reduces the flow of household savings into bank deposits in the financial sector. As a result the business sector has an incentive to increase employment and output and perhaps to borrow from the financial sector to finance that expansion. Increased employment in the business sector increases incomes in the household sector and further increases household expenditure and savings. These inter-sector linkages and feedbacks produce a response in aggregate economy greater than the initial change.

Expanding this simple example to include more sectors increases its complexity but does not change the basics. A change in behaviour within a sector, or disturbance from outside that sector, changes aggregate levels of output, employment and prices. A complete multisector macroeconomic theory and model is required to understand the effects, on the aggregate economy, of changes in either internal or external economic conditions. It is also essential for the design of policies to manage the macroeconomic conditions.

Mitigating the effects of a large random shock from outside the economy, like the COVID-19 pandemic, disrupts all sectors of the economy. Flows of income, expenditure, revenue, and output among sectors are reduced sharply both by the pandemic and by government and financial sectors policy responses.

Application Box 1.1 COVID-19 and the Economy

The COVID-19 pandemic attacked Canada in early 2020. It revealed the complexities and interdependencies that drive the macro economy. Control and elimination of the disease depends on stopping person to person transmission. This is why the government mandated personal and social distancing for individuals plus self-isolation and quarantine in some cases. In addition businesses, mainly in the service sector, that relied on face to face interactions with customers or live audiences were forced to close.

As a result, businesses lost sales revenues and cut output. They reduced employment to cut labour costs, but overhead costs remained. Households lost employment and employment incomes. They reduced their discretionary spending, but their overhead costs continued. As a result the economy faced a unique, simultaneous collapse in overall private supply and demand and the risk of a deep recession. The government and financial sectors intervened with fiscal and monetary policy support. Government introduced a wide range of new income supports for the household sector, and loan and subsidy programs to support businesses, funded by large increases in the government's budget deficit. The central bank lowered interest rates and increased the monetary base to support the government's borrowing requirements, and the credit demands on private banks and other financial institutions. The banking system lost the normal growth in customer deposits, but worked to accommodate the needs of their business and household clients.

This unprecedented support from government fiscal policy and central bank monetary policy will offset part of the loss in national output and income. But it will not reverse it. Recovery will begin with the reopening of business and the growth of employment at some time in the uncertain future. The size of the estimated effect of COVID-19 on the Canadian economy is stark. In its Monetary Policy Report, April 2020, The Bank of Canada estimates that real GDP in Canada will be 1% to 7.5% lower in 2020Q1 and lower by 15-30% in 2020Q2 than in 2019Q4. The Monetary Policy Report is available on the Bank of Canada's website at www.bankofcanada.ca.

Individual behaviours

Economic actions, at the level of the person or organization, form the subject matter of microeconomics. Formally, microeconomics is the study of individual behaviour in the context of scarcity. Not all individual behaviours are motivated by self-interest; many are motivated by a concern for the well being of society-at-large. Philanthropic societies are goal-oriented and seek to attain their objectives in an efficient manner.

Microeconomics: the study of individual behaviour in the context of scarcity.

Individual economic decisions need not be world-changing events, or motivated by a search for profit. Microeconomics is also about how we choose to spend our time and money. There are quite a few options to choose from: Sleep, work, study, food, shelter,

transportation, entertainment, recreation and so forth. Because both time and income are limited we cannot do all things all the time. Many choices are routine or are driven by necessity. You have to eat and you need a place to live. If you have a job you have committed some of your time to work, or if you are a student some of your time is committed to lectures and study. There is more flexibility in other choices. Critically, microeconomics seeks to understand and explain how we make choices and how those choices affect our behaviour in the workplace, the marketplace, and society more generally.

A critical element in making choices is that there exists a *scarcity* of time, or income or productive resources. Decisions are invariably subject to limits or constraints, and it is these constraints that make decisions both challenging and scientific.

Microeconomics also concerns business choices. How does a business use its funds and management skill to produce goods and services? The individual business operator or firm has to decide what to produce, how to produce it, how to sell it and in many cases, how to price it. To make and sell pizza, for example, the pizza parlour needs, in addition to a source of pizza ingredients, a store location (land), a pizza oven (capital), a cook and a sales person (labour). Payments for the use of these inputs generate income to those supplying them. If revenue from the sale of pizzas is greater than the costs of production, the business earns a profit for the owner. A business fails if it cannot cover its costs.

In these micro-level behaviours the decision makers have a common goal: To do as well as they can, *given the constraints imposed by the operating environment*. The individual wants to mix work and leisure in a way that makes her as happy or contented as possible. The entrepreneur aims at making a profit. These actors, or agents as we sometimes call them, are *maximizing*. Such maximizing behaviour is a central theme in this book and in economics at large.

Markets and government

Markets play a key role in coordinating the choices of individuals with the decisions of business. In modern market economies goods and services are supplied by both business and government. Hence we call them mixed economies. Some products or services are available through the marketplace to those who wish to buy them and have the necessary income—as in cases like coffee and wireless services. Other services are provided to all people through government programs like law enforcement and health care.

Mixed economy: goods and services are supplied both by private suppliers and government.

Markets offer the choice of a wide range of goods and services at various prices. Individuals can use their incomes to decide the pattern of expenditures and the bundle of goods and services they prefer. Businesses sell goods and services in the expectation that the market price will cover costs and yield a profit.

The market also allows for specialization and separation between production and use. Rather than each individual growing her own food, for example, she can sell her time or labour to employers in return for income. That income can then support her desired purchases. If businesses can produce food more cheaply than individuals the individual obviously gains from using the market – by both having the food to consume, and additional income with which to buy other goods and services. Economics seeks to explain how markets and specialization might yield such gains for individuals and society.

We will represent individuals and firms by envisaging that they have explicit objectives – to maximize their happiness or profit. However, this does not imply that individuals and firms are concerned only with such objectives. On the contrary, much of microeconomics and macroeconomics focuses upon the role of government: How it manages the economy through fiscal and monetary policy, how it redistributes through the tax-transfer system, how it supplies information to buyers and sets safety standards for products.

Since governments perform all of these society-enhancing functions, in large measure governments reflect the social ethos of voters. So, while these voters may be maximizing at the individual level in their everyday lives, and our models of human behaviour in microeconomics certainly emphasize this optimization, economics does not see individuals and corporations as being devoid of civic virtue or compassion, nor does it assume that only market-based activity is important. Governments play a central role in modern economies, to the point where they account for more than one third of all economic activity in the modern mixed economy.

Governments supply goods and services in many spheres, for example, health and education. The provision of public education is motivated both by a concern for equality and a realization that an educated labour force increases the productivity of an economy. Likewise, the provision of law and order, through our legal system broadly defined, represents more than a commitment to a just society at the individual level; without a legal system that enforces contracts and respects property rights, the private sector of the

economy would diminish dramatically as a result of corruption, uncertainty and insecurity. It is the lack of such a secure environment in many of the world's economies that inhibits their growth and prosperity.

Let us consider now the methods of economics, methods that are common to science-based disciplines.

1.2 Understanding through the use of models

Most students have seen an image of Ptolemy's concept of our Universe. Planet Earth forms the centre, with the other planets and our sun revolving around it. The ancients' anthropocentric view of the universe necessarily placed their planet at the centre. Despite being false, this view of our world worked reasonably well – in the sense that the ancients could predict celestial motions, lunar patterns and the seasons quite accurately.

More than one Greek astronomer believed that it was more natural for smaller objects such as the earth to revolve around larger objects such as the sun, and they knew that the sun had to be larger as a result of having studied eclipses of the moon and sun. Nonetheless, the Ptolemaic description of the universe persisted until Copernicus wrote his treatise "On the Revolutions of the Celestial Spheres" in the early sixteenth century. And it was another hundred years before the Church accepted that our corner of the universe is heliocentric. During this time evidence accumulated as a result of the work of Brahe, Kepler and Galileo. The time had come for the Ptolemaic *model* of the universe to be supplanted with a better *model*.

All disciplines progress and develop and explain themselves using models of reality. A model is a formalization of theory that facilitates scientific inquiry. Any history or philosophy of science book will describe the essential features of a model. First, it is a stripped down, or reduced, version of the phenomenon that is under study. It incorporates the key elements while disregarding what are considered to be secondary elements. Second, it should accord with reality. Third, it should be able to make meaningful predictions. Ptolemy's model of the known universe met these criteria: It was not excessively complicated (for example distant stars were considered as secondary elements in the universe and were excluded); it corresponded to the known reality of the day, and made pretty good predictions. Evidently not all models are correct and this was the case here.

Model: a formalization of theory that facilitates scientific inquiry.

In short, models are frameworks we use to organize how we think about a problem. Economists sometimes interchange the terms theories and models, though they are conceptually distinct. A theory is a logical view of how things work, and is frequently formulated on the basis of observation. A model is a formalization of the essential elements of a theory, and has the characteristics we described above. As an example of an economic model, suppose we theorize that a household's expenditure depends on its key characteristics: A corresponding model might specify that wealth, income, and household size determine its expenditures, while it might ignore other, less important, traits such as the household's neighbourhood or its religious beliefs. The model reduces and simplifies the theory to manageable dimensions. From such a reduced picture of reality we develop an analysis of how an economy and its components work.

Theory: a logical view of how things work, and is frequently formulated on the basis of observation.

An economist uses a model as a tourist uses a map. Any city map misses out some detail—traffic lights and speed bumps, for example. But with careful study you can get a good idea of the best route to take. Economists are not alone in this approach; astronomers, meteorologists, physicists, and genetic scientists operate similarly. Meteorologists disregard weather conditions in South Africa when predicting tomorrow's conditions in Winnipeg. Genetic scientists concentrate on the interactions of limited subsets of genes that they believe are the most important for their purpose. Even with huge computers, all of these scientists build *models* that concentrate on the essentials.

1.3 Opportunity cost and the market

Individuals face choices at every turn: In deciding to go to the hockey game tonight, you may have to forgo a concert; or you will have to forgo some leisure time this week in order to earn additional income for the hockey game ticket. Indeed, there is no such thing as a free lunch, a free hockey game or a free concert. In economics we say that these limits or constraints reflect opportunity cost. The opportunity cost of a choice is what must be sacrificed when a choice is made. That cost may be financial; it may be measured in time, or simply the alternative foregone.

Opportunity cost: what must be sacrificed when a choice is made.

Opportunity costs play a determining role in markets. It is precisely because individuals and organizations have different opportunity costs that they enter into exchange agreements. If you are a skilled plumber and an unskilled gardener, while your neighbour is a skilled gardener and an unskilled plumber, then you and your neighbour not only have different capabilities, you also

have different opportunity costs, and *you could gain by trading your skills*. Here's why. Fixing a leaking pipe has a low opportunity cost for you in terms of time: You can do it quickly. But pruning your apple trees will be costly because you must first learn how to avoid killing them and this may require many hours. Your neighbour has exactly the same problem, with the tasks in reverse positions. In a sensible world you would fix your own pipes *and* your neighbour's pipes, and she would ensure the health of the apple trees in both backyards.

If you reflect upon this 'sensible' solution—one that involves each of you achieving your objectives while minimizing the time input—you will quickly realize that it resembles the solution provided by the marketplace. You may not have a gardener as a neighbour, so you buy the services of a gardener in the marketplace. Likewise, your immediate neighbour may not need a leaking pipe repaired, but many others in your neighbourhood do, so you sell your service to them. You each specialize in the performance of specific tasks as a result of having different opportunity costs or different efficiencies. Let us now develop a model of exchange to illustrate the advantages of specialization and trade, and hence the markets that facilitate these activities. This model is developed with the help of some two-dimensional graphics.

1.4 A model of exchange and specialization

Production and specialization

We have two producers and two goods: Amanda and Zoe produce vegetables (V) and or fish (F). Their production capabilities are defined in Table 1.1 and in Figure 1.1, where the quantity of V appears on the vertical axis and the quantity of F on the horizontal axis. Zoe and Amanda each have 36-hour weeks and they devote that time to producing the two goods. But their efficiencies differ: Amanda requires two hours to produce a unit of V and three hours for a unit of F . As a consequence, if she devotes all of her time to V she can produce 18 units, or if she devotes all of her time to F she can produce 12 units. Or, she could share her time between the two. This environment can also be illustrated and analyzed graphically, as in Figure 1.1.

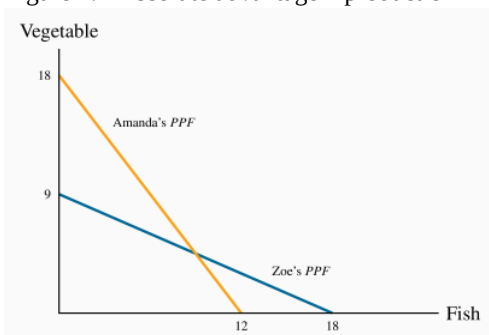
Table 1.1 Production possibilities in a two-person economy

	Hours/ fish	Hours/ vegetable	Fish specialization	Vegetable specialization
Amanda	3	2	12	18
Zoe	2	4	18	9

Each producer has a time allocation of 36 hours. By allocating total time to one activity, Amanda can produce 12 F or 18 V , Zoe can produce 18 F or 9 V . By splitting their time each person can also produce a combination of the two.

Two-dimensional graphics are a means of portraying the operation of a model, as defined above. We will use these graphical representations throughout the text. In this case, Amanda's production capability is represented by the line that meets the vertical axis at 18 and the horizontal axis at 12. The vertical point indicates that she can produce 18 units of V if she produces zero units of F – keep in mind that where V has a value of 18, Amanda has no time left for fish production. Likewise, if she devotes all of her time to fish she can produce 12 units, since each unit requires 3 of her 36 hours. The point $F=12$ is thus another possibility for her. In addition to these two possibilities, which we can term 'specialization', she could allocate her time to producing some of each good. For example, by dividing her 36 hours equally she could produce 6 units of F and 9 units of V . A little computation will quickly convince us that different allocations of her time will lead to combinations of the two goods that lie along a straight line joining the specialization points.

Figure 1.1 Absolute advantage – production



Amanda's *PPF* indicates that she can produce either 18V (and zero *F*), or 12F (and zero *V*), or some combination. Zoe's *PPF* indicates she can produce either 9V (and zero *F*), or 18F (and zero *V*), or some combination. Amanda is more efficient in producing *V* and Zoe is more efficient at producing *F*.

We will call this straight line Amanda's production possibility frontier (*PPF*): It is the combination of goods she can produce while using all of her resources – time. She could not produce combinations of goods represented by points beyond this line (to the top right). She could indeed produce combinations below it (lower left) – for example, a combination of 4 units of *V* and 4 units of *F*; but such points would not require all of her time. The (4,4) combination would require just 20 hours. In sum, points beyond this line are not feasible, and points within it do not require all of her time resources.

Production possibility frontier (*PPF*): the combination of goods that can be produced using all of the resources available.

Having developed Amanda's *PPF*, it is straightforward to develop a corresponding set of possibilities for Zoe. If she requires 4 hours to produce a unit of *V* and 2 hours to produce a unit of *F*, then her 36 hours will enable her to specialize in 9 units of *V* or 18 units of *F*; or she could produce a combination represented by the straight line that joins these two specialty extremes.

Consider now the opportunity costs for each person. Suppose Amanda is currently producing 18 *V* and zero *F*, and considers producing some *F* and less *V*. For each unit of *F* she wishes to produce, it is evident from her *PPF* that she must sacrifice 1.5 units of *V*. This is because *F* requires 50% more hours than *V*. Her trade-off is 1.5:1.0. The additional time requirement is also expressed in the intercept values: She could give up 18 units of *V* and produce 12 units of *F* instead; this again is a ratio of 1.5:1.0. This ratio defines her opportunity cost: The cost of an additional unit of *F* is that 1.5 units of *V* must be 'sacrificed'.

Applying the same reasoning to Zoe's *PPF*, her opportunity cost is 0.5:1; she must sacrifice one half of a unit of *V* to free up enough time to produce one unit of *F*.

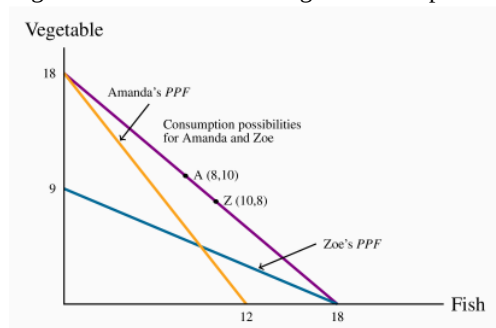
So we have established two things about Amanda and Zoe's production possibilities. First, if Amanda specializes in *V* she can produce more than Zoe, just as Zoe can produce more than Amanda if Zoe specializes in *F*. Second, their opportunity costs are different: Amanda must sacrifice more *V* than Zoe in producing one more unit of *F*. The different opportunity costs translate into potential gains for each individual.

The gains from exchange

We shall illustrate the gains that arise from specialization and exchange graphically. Note first that if these individuals are self-sufficient, in the sense that they consume their own production, each individual's consumption combination will lie on their own *PPF*. For example, Amanda could allocate half of her time to each good, and produce (and consume) 6F and 9V. Such a point necessarily lies on her *PPF*. Likewise for Zoe. So, *in the absence of exchange*, each individual's *PPF* is also her consumption possibility frontier (*CPF*). In Figure 1.1 the *PPF* for each individual is thus also her *CPF*.

Consumption possibility frontier (*CPF*): the combination of goods that can be consumed as a result of a given production choice.

Figure 1.2 Absolute advantage – consumption



With specialization and trade at a rate of 1:1 they consume along the line joining the specialization points. If Amanda trades 8V to Zoe in return for 8F, Amanda moves to the point A(8,10) and Zoe to Z(10,8). Each can consume more after specialization than before specialization.

Upon realizing that they are not equally efficient in producing the two goods, they decide to specialize completely in producing just the single good where they are most efficient. Amanda specializes in *V* and Zoe in *F*. Next they must agree to a rate at which to exchange *V* for *F*. Since Amanda's opportunity cost is 1.5:1 and Zoe's is 0.5:1, suppose they agree to exchange *V* for *F* at an intermediate rate of 1:1. There are many trading, or exchange, rates possible; our purpose is to illustrate that gains are possible for *both* individuals at some exchange rate. The choice of this rate also makes the graphic as simple as possible. At this exchange rate,

18V must exchange for 18F. In Figure 1.2, this means that each individual is now able to consume along the line joining the coordinates (0,18) and (18,0).¹ This is because Amanda produces 18V and she can trade at a rate of 1:1, while Zoe produces 18F and trades at the same rate of 1:1.

The fundamental result illustrated in Figure 1.2 is that, as a result of specialization and trade, each individual can consume combinations of goods that lie on a line beyond her initial consumption possibilities. Their consumption well-being has thus improved. For example, suppose Amanda trades away 8V to Zoe and obtains 8F in return. The points 'A' and 'Z' with coordinates (8,10) and (10,8) respectively define their final consumption. Pre-specialization, if Amanda wished to consume 8F she would have been constrained to consume 6V rather than the 10V now possible. Zoe benefits correspondingly.²

The foregoing example illustrates that trade is not a zero-sum game; it has a positive net value because both parties to the trade can gain. A zero-sum gain is where the gains to one party exactly offset the losses to another. This is an extraordinarily important principle in trade negotiations, whether international or domestic.

A zero-sum game is an interaction where the gain to one party equals the loss to another party.

Market design

In the preceding example we have shown that specialization provides scope for gains that can accrue to those participating in the exchange. But this tells us little about how a market for these products comes into being: how does the exchange take place, and how is information transmitted? The answer is that while some markets have evolved historically to their current state, many markets are designed by an institution or a firm. Fruit and vegetable markets have been with us for thousands of years - since we ceased being purely a hunter-gatherer society. They exist in every community in the world economy. In contrast, the Dutch tulip auction was designed in the early 1600s and exists in basically the same form to this day: the auctioneer begins with a high price, lowers it at known time intervals (measured in seconds or minutes) until some buyer signals that she is willing to purchase the lot on offer. Supermarkets in contrast offer goods at a fixed price. Government contracts are normally signed after a tendering process, in which interested suppliers submit bids. Amazon Inc. is currently experimenting with cashierless 'bricks and mortar' stores that monitor all transactions electronically. Craig's List and E-Bay have their own sets of rules.

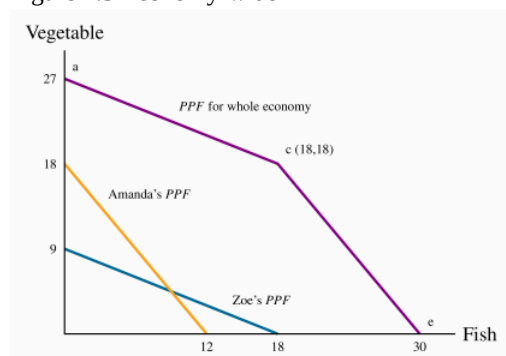
In each of these cases markets are designed, frequently with a specific objective on the part of the supplier or the mediating institution: Amazon wants to increase its share of all goods trades; governments wish to limit costs. Markets do not all grow spontaneously and the structure of a market will influence how the gains from trade are distributed.

1.5 Economy-wide production possibilities

The *PPFs* in Figures 1.1 and 1.2 define the amounts of the goods that each *individual* can produce while using all of their fixed productive capacity—time in this instance. The national, or economy-wide, *PPF* for this two-person economy reflects these individual possibilities combined. Such a frontier can be constructed using the individual frontiers as the component blocks.

First let us define this economy-wide frontier precisely. The economy-wide *PPF* is the set of goods and services combinations that can be produced in the economy when all available productive resources are in use. Figure 1.3 contains both of the individual frontiers plus the aggregate of these, represented by the kinked line *ace*. The point on the V axis, *a*=27, represents the total amount of V that could be produced if both individuals devoted all of their time to it. The point *e*=30 on the horizontal axis is the corresponding total for fish.

Figure 1.3 Economy-wide PPF



From *a*, to produce Fish it is more efficient to use Zoe because her opportunity cost is less (segment *ac*). When Zoe is completely specialized, Amanda produces (*ce*). With complete specialization this economy can produce 27V or 30F.

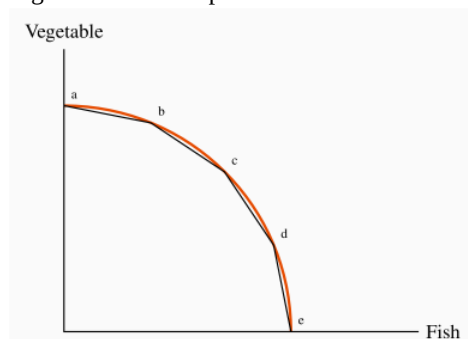
Economy-wide PPF: the set of goods and services combinations that can be produced in the economy when all available productive resources are in use.

To understand the point c , imagine initially that all resources are devoted to V . From such a point, a , consider a reduction in V and an increase in F . The most efficient way of increasing F production at the point a is to use the individual whose opportunity cost is lower. Zoe can produce one unit of F by sacrificing just 0.5 units of V , whereas Amanda must sacrifice 1.5 units of V to produce 1 unit of F . Hence, at this stage Amanda should stick to V and Zoe should devote some time to fish. In fact as long as we want to produce more fish Zoe should be the one to do it, until she has exhausted her time resource. This occurs after she has produced 18 F and has ceased producing V . At this point the economy will be producing 18 V and 18 F – the point c .

From this combination, if the economy wishes to produce more fish Amanda must become involved. Since her opportunity cost is 1.5 units of V for each unit of F , the next segment of the economy-wide PPF must see a reduction of 1.5 units of V for each additional unit of F . This is reflected in the segment ce . When both producers allocate all of their time to F the economy can produce 30 units. Hence the economy's PPF is the two-segment line ace . Since this has an outward kink, we call it concave (rather than convex).

As a final step consider what this PPF would resemble if the economy were composed of many persons with differing efficiencies. A little imagination suggests (correctly) that it will have a segment for each individual and continue to have its outward concave form. Hence, a four-person economy in which each person had a different opportunity cost could be represented by the segmented line $abcde$, in Figure 1.4. Furthermore, we could represent the PPF of an economy with a very large number of such individuals by a somewhat smooth PPF that accompanies the 4-person PPF. The logic for its shape continues to be the same: As we produce less V and more F we progressively bring into play resources, or individuals, whose opportunity cost, in terms of reduced V is higher.

Figure 1.4 A multi-person PPF



The PPF for the whole economy, $abcde$, is obtained by allocating productive resources most efficiently. With many individuals we can think of the PPF as the *concave envelope* of the individual capabilities.

The outputs V and F in our economic model require just one input – time, but if other productive resources were required the result would be still a concave PPF. Furthermore, we generally interpret the PPF to define the output possibilities *when the economy is running at its normal capacity*. In this example, we consider a work week of 36 hours to be the 'norm'. Yet it is still possible that the economy's producers might work some additional time in exceptional circumstances, and this would increase total production possibilities. This event would be represented by an outward movement of the PPF.

1.6 Aggregate output, growth and business cycles

The PPF can be used to illustrate several aspects of macroeconomics: In particular, the level of an economy's output, the growth of national and per capita output over time, and short-run business-cycle fluctuations in national output and employment.

Aggregate output

An economy's capacity to produce goods and services depends on its endowment of resources and the productivity of those resources. The two-person, two-product examples in the previous section reflect this.

The productivity of labour, defined as output per worker or per hour, depends on:

- The skill, knowledge and experience of the labour force;
- The capital stock: Buildings, machinery, equipment, and software the labour force has to work with; and
- The current state of technology.

The **productivity of labour** is the output of goods and services per worker.

An economy's **capital stock** is the buildings, machinery, equipment and software used in producing goods and services.

The economy's output, which we define by Y , can be defined as the output per worker times the number of workers; hence, we can write:

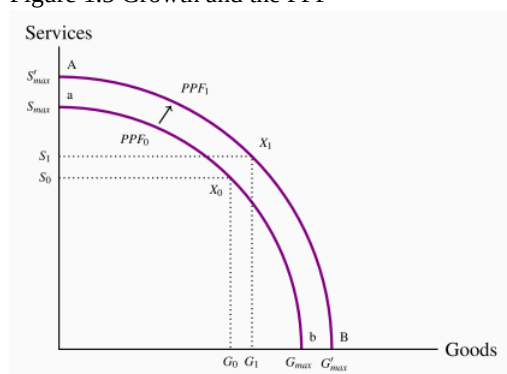
$$Y = (\text{number of workers employed}) \times (\text{output per worker}).$$

When the employment of labour corresponds to 'full employment' in the sense that everyone willing to work at current wage rates and normal hours of work is working, the economy's actual output is also its capacity output Y_c . We also term this capacity output as full employment output:

$$\text{Full employment output } Y_c = (\text{number of workers at full employment}) \times (\text{output per worker}).$$

Suppose the economy is operating with full employment of resources producing outputs of two types: Goods and services. In Figure 1.5, PPF_0 shows the different combinations of goods and services the economy can produce in a particular year using all its labour, capital and the best technology available at the time.

Figure 1.5 Growth and the PPF



Economic growth is illustrated by an outward shift in the PPF from PPF_0 to PPF_1 . PPF_1 shows the economy can produce more in both sectors than with PPF_0 .

An aggregate economy produces a large variety of outputs in two broad categories. Goods are the products of the agriculture, forestry, mining, manufacturing and construction industries. Services are provided by the wholesale and retail trade, transportation, hospitality, finance, health care, education, legal and other service sectors. As in the two-product examples used earlier, the shape of the PPF illustrates the opportunity cost of increasing the output of either product type. We are not concerned with who supplies the products for the moment: It may be the private sector or the government.

Point X_0 on PPF_0 shows one possible structure of capacity output. This combination may reflect the pattern of demand and hence expenditures in this economy. Output structures and therefore the shapes of $PPFs$ differ among economies with different income levels. High-income economies spend more on services than goods and produce higher ratios of services to goods. Middle income countries produce lower ratios of services to goods, and low income countries much lower ratios of services to goods. For example, in 2017, the structure of national output in Canada was 70 percent services and 30 percent goods, while in Mexico the structure was 48 percent services and 52 percent goods.

Different countries also have different $PPFs$ and different output structures, depending on their resource endowments, labour forces, capital stocks, technology and expenditure patterns.

Economic growth

Three things contribute to growth in the economy. The labour supply grows as the population expands; the stock of capital grows as spending by business (and government) on buildings, machinery, information technology and so forth increases; and labour-force productivity grows as a result of experience, the development of scientific knowledge combined with product and process innovations, and advances in the technology of production. Combined, these developments expand capacity output over time. In Figure 1.5 economic growth shifts the PPF out from PPF_0 to PPF_1 .

This basic description covers the key sources of growth in total output. Economies differ in their rates of overall economic growth as a result of different rates of growth in labour force, in capital stock, and improvements in technology. But improvements in standards of living require more than growth in total output. Increases in output *per worker* and *per person* are necessary. Sustained

increases in living standards require sustained growth in labour productivity, which in turn is based on advances in the technology along with the amount of capital each worker has to work with. Furthermore, if the growth in output is to benefit society at large, workers across the board need to see an increase in their earnings. As we shall explore in Chapter 13, several developed countries have seen the fruits of growth concentrated in the hands of the highest income earners.

Recessions and booms

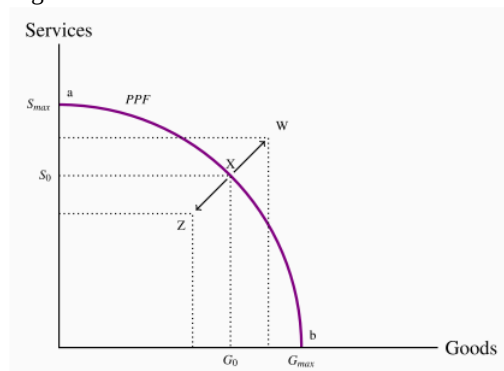
A prime objective of economic policy is to ensure that the economy operates on or near the *PPF* – it should use its resources to capacity and have minimal unemployment. However, economic conditions are seldom tranquil for long periods of time. Unpredictable changes in business expectations of future profits, in consumer confidence, in financial markets, in commodity and energy prices, in trade agreements and disputes, in economic conditions in major trading partners, in government policy and many other events disrupt patterns of expenditure and output. Some of these changes disturb the level of total expenditure and thus the demand for total output. Others disturb the conditions of production and thus the economy's production capacity. Whatever the exact cause, the economy may be pushed off its current *PPF*. If expenditures on goods and services decline, the economy may experience a recession. Output would fall short of capacity output and unemployment would rise. Alternatively, times of rapidly growing expenditure and output may result in an economic boom: Output and employment expand beyond capacity levels.

An **economic recession** occurs when output falls below the economy's capacity output.

A **boom** is a period of high growth that raises output above normal capacity output.

Recent history provides examples. Following the financial crisis of 2008-09 that hit the US and many other developed economies, many economies were pushed into recessions. Expenditure on new residential construction collapsed for lack of income and secure financing, as did business investment, consumption spending and exports. Lower expenditures reduced producers' revenues, forcing cuts in output and employment and reducing household incomes. Lower incomes led to further cutbacks in spending. In Canada in 2009 aggregate output declined by 2.9 percent, employment declined by 1.6 percent and the unemployment rate rose from 6.1 percent in 2008 to 8.3 percent by 2010. The world's economies have been slow to recover, and even by 2019 the output in several developed economies was no higher than it was in 2008. Canada's recession was not nearly as severe as the recessions in economies such as Spain, Italy and Greece; but output between 2009 and 2019 has been below the potential of the Canadian economy. In the third quarter of 2019 the national output was about 0.7 percent below potential output and the unemployment rate was 5.5 percent.

Figure 1.6 Booms and recessions



Economic recessions leave the economy below its normal capacity; the economy might be driven to a point such as Z. Economic expansions, or booms, may drive capacity above its normal level, to a point such as W.

An economy in a recession is operating inside its *PPF*. The fall in output from X to Z in Figure 1.6 illustrates the effect of a recession. Expenditures on goods and services have declined. Output is less than capacity output, unemployment is up and some plant capacity is idle. Labour income and business profits are lower. More people would like to work and business would like to produce and sell more output, but it takes time for interdependent product, labour and financial markets in the economy to adjust and increase employment and output. Monetary and fiscal policy may be productive in specific circumstances, to stimulate demand, increase output and employment and move the economy back to capacity output and full employment. The development and implementation of such policies form the core of macroeconomics.

Alternatively, an unexpected increase in demand for exports would increase output and employment. Higher employment and output would increase incomes and expenditure, and in the process spread the effects of higher output sales to other sectors of the

economy. The economy would move outside its *PPF*, for example to *W* in Figure 1.6, by using its resources more intensively than normal. Unemployment would fall and overtime work would increase. Extra production shifts would run plant and equipment for longer hours and work days than were planned when it was designed and installed. Output at this level may not be sustainable, because shortages of labour and materials along with excessive rates of equipment wear and tear would push costs and prices up. Again, we will examine how the economy reacts to such a state in our macroeconomic analysis.

Output and employment in the Canadian economy over the past twenty years fluctuated about growth trend in the way Figure 1.6 illustrates. For several years prior to 2008 the Canadian economy operated slightly above its capacity; but once the recession arrived monetary and fiscal policy were used to fight it – to bring the economy back from a point such as *Z* towards a point such as *X* on the *PPF*.

Macroeconomic models and policy

The *PPF* diagrams illustrate the main dimensions of macroeconomics: Capacity output, growth in capacity output and business cycle fluctuations in actual output relative to capacity. But these diagrams do not offer explanations and analysis of macroeconomic activity. We need a macroeconomic *model* to understand and evaluate the causes and consequences of business cycle fluctuations. As we shall see, these models are based on explanations of expenditure decisions by households and business, financial market conditions, production costs and producer pricing decisions at different levels of output. Models also capture the objectives of fiscal and monetary policies and provide a framework for policy evaluation. A full macroeconomic model integrates different sector behaviours and the feedbacks across sectors that can moderate or amplify the effects of changes in one sector on national output and employment.

Conclusion

We have covered a lot of ground in this introductory chapter. It is intended to open up the vista of economics to the new student in the discipline. Economics is powerful and challenging, and the ideas we have developed here will serve as conceptual foundations for our exploration of the subject.

Key Terms

Macroeconomics studies the economy as system in which linkages and feedbacks among sectors determine national output, employment and prices.

Microeconomics is the study of individual behaviour in the context of scarcity.

Mixed economy: goods and services are supplied both by private suppliers and government.

Model is a formalization of theory that facilitates scientific inquiry.

Theory is a logical view of how things work, and is frequently formulated on the basis of observation.

Opportunity cost of a choice is what must be sacrificed when a choice is made.

Production possibility frontier (PPF) defines the combination of goods that can be produced using all of the resources available.

Consumption possibility frontier (CPF): the combination of goods that can be consumed as a result of a given production choice.

A **zero-sum game** is an interaction where the gain to one party equals the loss to another party.

Economy-wide PPF is the set of goods combinations that can be produced in the economy when all available productive resources are in use.

Productivity of labour is the output of goods and services per worker.

Capital stock: the buildings, machinery, equipment and software used in producing goods and services.

Full employment output $Y_c = (\text{number of workers at full employment}) \times (\text{output per worker})$. **Recession:** when output falls below the economy's capacity output. **Boom:** a period of high growth that raises output above normal capacity output.

Exercises for Chapter 1

EXERCISE 1.1

An economy has 100 identical workers. Each one can produce four cakes or three shirts, regardless of the number of other individuals producing each good.

1. How many cakes can be produced in this economy when all the workers are cooking?

2. How many shirts can be produced in this economy when all the workers are sewing?
3. On a diagram with cakes on the vertical axis, and shirts on the horizontal axis, join these points with a straight line to form the *PPF*.
4. Label the inefficient and unattainable regions on the diagram.

EXERCISE 1.2

In the table below are listed a series of points that define an economy's production possibility frontier for goods *Y* and *X*.

Y	1000	900	800	700	600	500	400	300	200	100	0
X	0	1600	2500	3300	4000	4600	5100	5500	5750	5900	6000

1. Plot these pairs of points to scale, on graph paper, or with the help of a spreadsheet.
2. Given the shape of this *PPF* is the economy made up of individuals who are similar or different in their production capabilities?
3. What is the opportunity cost of producing 100 more *Y* at the combination ($X=5500, Y=300$).
4. Suppose next there is technological change so that at every output level of good *Y* the economy can produce 20 percent more *X*. Enter a new row in the table containing the new values, and plot the new *PPF*.

EXERCISE 1.3

Using the *PPF* that you have graphed using the data in Exercise 1.2, determine if the following combinations are attainable or not: ($X=3000, Y=720$), ($X=4800, Y=480$).

EXERCISE 1.4

You and your partner are highly efficient people. You can earn \$20 per hour in the workplace; your partner can earn \$30 per hour.

1. What is the opportunity cost of one hour of leisure for you?
2. What is the opportunity cost of one hour of leisure for your partner?
3. Now consider what a *PPF* would look like: You can produce/consume two things, leisure and income. Since income buys things you can think of the *PPF* as having these two 'products' – leisure and consumption goods/services. So, with leisure on the horizontal axis and income in dollars is on the vertical axis, plot your *PPF*. You can assume that you have 12 hours per day to allocate to either leisure or income. [Hint: the leisure axis will have an intercept of 12 hours. The income intercept will have a dollar value corresponding to where all hours are devoted to work.]
4. Draw the *PPF* for your partner.

EXERCISE 1.5

Louis and Carrie Anne are students who have set up a summer business in their neighbourhood. They cut lawns and clean cars. Louis is particularly efficient at cutting the grass – he requires one hour to cut a typical lawn, while Carrie Anne needs one and one half hours. In contrast, Carrie Anne can wash a car in a half hour, while Louis requires three quarters of an hour.

1. If they decide to specialize in the tasks, who should cut the grass and who should wash cars?
2. If they each work a twelve hour day, how many lawns can they cut and how many cars can they wash if they each specialize in performing the task where they are most efficient?
3. Illustrate the *PPF* for each individual where lawns are on the horizontal axis and car washes on the vertical axis, if each individual has twelve hours in a day.

EXERCISE 1.6

Continuing with the same data set, suppose Carrie Anne's productivity improves so that she can now cut grass as efficiently as Louis; that is, she can cut grass in one hour, and can still wash a car in one half of an hour.

1. In a new diagram draw the *PPF* for each individual.
2. In this case does specialization matter if they are to be as productive as possible as a team?
3. Draw the *PPF* for the whole economy, labelling the intercepts and the 'kink' point coordinates.

EXERCISE 1.7

Going back to the simple *PPF* plotted for Exercise 1.1 where each of 100 workers can produce either four cakes or three shirts, suppose a recession reduces demand for the outputs to 220 cakes and 129 shirts.

1. Plot this combination of outputs in the diagram that also shows the *PPF*.
 2. How many workers are needed to produce this output of cakes and shirts?
 3. What percentage of the 100 worker labour force is unemployed?
1. When two values, separated by a comma, appear in parentheses, the first value refers to the horizontal-axis variable, and the second to the vertical-axis variable.
2. In the situation we describe above one individual is absolutely more efficient in producing one of the goods and absolutely less efficient in the other. We will return to this model in Chapter 15 and illustrate that consumption gains of the type that arise here can also result if one of the individuals is absolutely more efficient in producing both goods, but that the degree of such advantage differs across goods.

This page titled [1: Introduction to key ideas](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine](#) (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.1: What's it all about?

The big issues

Economics is the study of human behaviour. Since it uses scientific methods it is called a social science. We study human behaviour to better understand and improve our world. During his acceptance speech, a recent Nobel Laureate in Economics suggested:

Economics, at its best, is a set of ideas and methods for the improvement of society. It is not, as so often seems the case today, a set of ideological rules for asserting why we cannot face the challenges of stagnation, job loss and widening inequality.

Christopher Sims, Nobel Laureate in Economics 2011

This is an elegant definition of economics and serves as a timely caution about the perils of ideology. Economics evolves continuously as current observations and experience provide new evidence about economic behaviour and relationships. Inference and policy recommendations based on earlier theories, observations and institutional structures require constant analysis and updating if they are to furnish valuable responses to changing conditions and problems.

Much of today's developed world still faces severe challenges as a result of the financial crisis that began in 2008. Unemployment rates among young people are at historically high levels in several economies, many government balance sheets are in disarray, and inequality is on the rise. In addition to the challenges posed by this severe economic cycle, the world simultaneously faces structural upheaval: Overpopulation, climate change, political instability and globalization challenge us to understand and modify our behaviour.

These challenges do not imply that our world is deteriorating. Literacy rates have been rising dramatically in the developing world for decades; child mortality has plummeted; family size is a fraction of what it was 50 years ago; prosperity is on the rise in much of Asia; life expectancy is increasing universally and deaths through wars are in a state of long-term decline.

These developments, good and bad, have a universal character and affect billions of individuals. They involve an understanding of economies as large organisms with interactive components.

Aggregate output in a national economy

A national economy is a complete multi-sector system, made up of household, business, financial, government and international sectors. Each of these sectors is an aggregate or sum of a many smaller economic units with very similar characteristics. The government sector, for example aggregates the taxing and spending activities of local, provincial and national governments. Similarly, the household sector is an aggregate of the income, spending and saving of all households but not the specifics of each individual household. Economic activity within any one of these sectors reflects, in part, the conditions and choices made in that sector. But it is also affects and is affected by conditions and actions in the other sectors. These *interactions* and *feedbacks* within the system mean that the workings of the macro-economy are more complex than the operation of the sum of its parts.

Macroeconomics: the study of the economy as a system in which interactions and feedbacks among sectors determine national output, employment and prices.

For example, consider a simple economy with just household, business and financial sectors. The household sector earns income by providing labour to the other sectors. Households make choices about spending or saving this income. Businesses make decisions about the sizes of their establishments, their labour forces, and their outputs of goods and services. The financial sector provides banking services: bank deposits, loans, and the payments system used by all three sectors.

Suppose households decide to spend more on goods and services and save less. That decision by itself does not change household sector income, but it does increase business sector sales and revenues. It also reduces the flow of household savings into bank deposits in the financial sector. As a result the business sector has an incentive to increase employment and output and perhaps to borrow from the financial sector to finance that expansion. Increased employment in the business sector increases incomes in the household sector and further increases household expenditure and savings. These inter-sector linkages and feedbacks produce a response in aggregate economy greater than the initial change.

Expanding this simple example to include more sectors increases its complexity but does not change the basics. A change in behaviour within a sector, or disturbance from outside that sector, changes aggregate levels of output, employment and prices. A complete multisector macroeconomic theory and model is required to understand the effects, on the aggregate economy, of changes in either internal or external economic conditions. It is also essential for the design of policies to manage the macroeconomic conditions.

Mitigating the effects of a large random shock from outside the economy, like the COVID-19 pandemic, disrupts all sectors of the economy. Flows of income, expenditure, revenue, and output among sectors are reduced sharply both by the pandemic and by government and financial sectors policy responses.

Application Box 1.1 COVID-19 and the Economy

The COVID-19 pandemic attacked Canada in early 2020. It revealed the complexities and interdependencies that drive the macro economy. Control and elimination of the disease depends on stopping person to person transmission. This is why the government mandated personal and social distancing for individuals plus self-isolation and quarantine in some cases. In addition businesses, mainly in the service sector, that relied on face to face interactions with customers or live audiences were forced to close.

As a result, businesses lost sales revenues and cut output. They reduced employment to cut labour costs, but overhead costs remained. Households lost employment and employment incomes. They reduced their discretionary spending, but their overhead costs continued. As a result the economy faced a unique, simultaneous collapse in overall private supply and demand and the risk of a deep recession. The government and financial sectors intervened with fiscal and monetary policy support. Government introduced a wide range of new income supports for the household sector, and loan and subsidy programs to support businesses, funded by large increases in the government's budget deficit. The central bank lowered interest rates and increased the monetary base to support the government's borrowing requirements, and the credit demands on private banks and other financial institutions. The banking system lost the normal growth in customer deposits, but worked to accommodate the needs of their business and household clients.

This unprecedented support from government fiscal policy and central bank monetary policy will offset part of the loss in national output and income. But it will not reverse it. Recovery will begin with the reopening of business and the growth of employment at some time in the uncertain future. The size of the estimated effect of COVID-19 on the Canadian economy is stark. In its Monetary Policy Report, April 2020, The Bank of Canada estimates that real GDP in Canada will be 1% to 7.5% lower in 2020Q1 and lower by 15-30% in 2020Q2 than in 2019Q4. The Monetary Policy Report is available on the Bank of Canada's website at www.bankofcanada.ca.

Individual behaviours

Economic actions, at the level of the person or organization, form the subject matter of microeconomics. Formally, microeconomics is the study of individual behaviour in the context of scarcity. Not all individual behaviours are motivated by self-interest; many are motivated by a concern for the well being of society-at-large. Philanthropic societies are goal-oriented and seek to attain their objectives in an efficient manner.

Microeconomics: the study of individual behaviour in the context of scarcity.

Individual economic decisions need not be world-changing events, or motivated by a search for profit. Microeconomics is also about how we choose to spend our time and money. There are quite a few options to choose from: Sleep, work, study, food, shelter, transportation, entertainment, recreation and so forth. Because both time and income are limited we cannot do all things all the time. Many choices are routine or are driven by necessity. You have to eat and you need a place to live. If you have a job you have committed some of your time to work, or if you are a student some of your time is committed to lectures and study. There is more flexibility in other choices. Critically, microeconomics seeks to understand and explain how we make choices and how those choices affect our behaviour in the workplace, the marketplace, and society more generally.

A critical element in making choices is that there exists a *scarcity* of time, or income or productive resources. Decisions are invariably subject to limits or constraints, and it is these constraints that make decisions both challenging and scientific.

Microeconomics also concerns business choices. How does a business use its funds and management skill to produce goods and services? The individual business operator or firm has to decide what to produce, how to produce it, how to sell it and in many cases, how to price it. To make and sell pizza, for example, the pizza parlour needs, in addition to a source of pizza ingredients, a store location (land), a pizza oven (capital), a cook and a sales person (labour). Payments for the use of these inputs generate

income to those supplying them. If revenue from the sale of pizzas is greater than the costs of production, the business earns a profit for the owner. A business fails if it cannot cover its costs.

In these micro-level behaviours the decision makers have a common goal: To do as well as they can, *given the constraints imposed by the operating environment*. The individual wants to mix work and leisure in a way that makes her as happy or contented as possible. The entrepreneur aims at making a profit. These actors, or agents as we sometimes call them, are *maximizing*. Such maximizing behaviour is a central theme in this book and in economics at large.

Markets and government

Markets play a key role in coordinating the choices of individuals with the decisions of business. In modern market economies goods and services are supplied by both business and government. Hence we call them mixed economies. Some products or services are available through the marketplace to those who wish to buy them and have the necessary income—as in cases like coffee and wireless services. Other services are provided to all people through government programs like law enforcement and health care.

Mixed economy: goods and services are supplied both by private suppliers and government.

Markets offer the choice of a wide range of goods and services at various prices. Individuals can use their incomes to decide the pattern of expenditures and the bundle of goods and services they prefer. Businesses sell goods and services in the expectation that the market price will cover costs and yield a profit.

The market also allows for specialization and separation between production and use. Rather than each individual growing her own food, for example, she can sell her time or labour to employers in return for income. That income can then support her desired purchases. If businesses can produce food more cheaply than individuals the individual obviously gains from using the market – by both having the food to consume, and additional income with which to buy other goods and services. Economics seeks to explain how markets and specialization might yield such gains for individuals and society.

We will represent individuals and firms by envisaging that they have explicit objectives – to maximize their happiness or profit. However, this does not imply that individuals and firms are concerned only with such objectives. On the contrary, much of microeconomics and macroeconomics focuses upon the role of government: How it manages the economy through fiscal and monetary policy, how it redistributes through the tax-transfer system, how it supplies information to buyers and sets safety standards for products.

Since governments perform all of these society-enhancing functions, in large measure governments reflect the social ethos of voters. So, while these voters may be maximizing at the individual level in their everyday lives, and our models of human behaviour in microeconomics certainly emphasize this optimization, economics does not see individuals and corporations as being devoid of civic virtue or compassion, nor does it assume that only market-based activity is important. Governments play a central role in modern economies, to the point where they account for more than one third of all economic activity in the modern mixed economy.

Governments supply goods and services in many spheres, for example, health and education. The provision of public education is motivated both by a concern for equality and a realization that an educated labour force increases the productivity of an economy. Likewise, the provision of law and order, through our legal system broadly defined, represents more than a commitment to a just society at the individual level; without a legal system that enforces contracts and respects property rights, the private sector of the economy would diminish dramatically as a result of corruption, uncertainty and insecurity. It is the lack of such a secure environment in many of the world's economies that inhibits their growth and prosperity.

Let us consider now the methods of economics, methods that are common to science-based disciplines.

This page titled [1.1: What's it all about?](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.2: Understanding through the use of models

Most students have seen an image of Ptolemy's concept of our Universe. Planet Earth forms the centre, with the other planets and our sun revolving around it. The ancients' anthropocentric view of the universe necessarily placed their planet at the centre. Despite being false, this view of our world worked reasonably well – in the sense that the ancients could predict celestial motions, lunar patterns and the seasons quite accurately.

More than one Greek astronomer believed that it was more natural for smaller objects such as the earth to revolve around larger objects such as the sun, and they knew that the sun had to be larger as a result of having studied eclipses of the moon and sun. Nonetheless, the Ptolemaic description of the universe persisted until Copernicus wrote his treatise "On the Revolutions of the Celestial Spheres" in the early sixteenth century. And it was another hundred years before the Church accepted that our corner of the universe is heliocentric. During this time evidence accumulated as a result of the work of Brahe, Kepler and Galileo. The time had come for the Ptolemaic *model* of the universe to be supplanted with a better *model*.

All disciplines progress and develop and explain themselves using models of reality. A model is a formalization of theory that facilitates scientific inquiry. Any history or philosophy of science book will describe the essential features of a model. First, it is a stripped down, or reduced, version of the phenomenon that is under study. It incorporates the key elements while disregarding what are considered to be secondary elements. Second, it should accord with reality. Third, it should be able to make meaningful predictions. Ptolemy's model of the known universe met these criteria: It was not excessively complicated (for example distant stars were considered as secondary elements in the universe and were excluded); it corresponded to the known reality of the day, and made pretty good predictions. Evidently not all models are correct and this was the case here.

Model: a formalization of theory that facilitates scientific inquiry.

In short, models are frameworks we use to organize how we think about a problem. Economists sometimes interchange the terms theories and models, though they are conceptually distinct. A theory is a logical view of how things work, and is frequently formulated on the basis of observation. A model is a formalization of the essential elements of a theory, and has the characteristics we described above. As an example of an economic model, suppose we theorize that a household's expenditure depends on its key characteristics: A corresponding model might specify that wealth, income, and household size determine its expenditures, while it might ignore other, less important, traits such as the household's neighbourhood or its religious beliefs. The model reduces and simplifies the theory to manageable dimensions. From such a reduced picture of reality we develop an analysis of how an economy and its components work.

Theory: a logical view of how things work, and is frequently formulated on the basis of observation.

An economist uses a model as a tourist uses a map. Any city map misses out some detail—traffic lights and speed bumps, for example. But with careful study you can get a good idea of the best route to take. Economists are not alone in this approach; astronomers, meteorologists, physicists, and genetic scientists operate similarly. Meteorologists disregard weather conditions in South Africa when predicting tomorrow's conditions in Winnipeg. Genetic scientists concentrate on the interactions of limited subsets of genes that they believe are the most important for their purpose. Even with huge computers, all of these scientists build *models* that concentrate on the essentials.

This page titled [1.2: Understanding through the use of models](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.3: Opportunity cost and the market

Individuals face choices at every turn: In deciding to go to the hockey game tonight, you may have to forgo a concert; or you will have to forgo some leisure time this week in order to earn additional income for the hockey game ticket. Indeed, there is no such thing as a free lunch, a free hockey game or a free concert. In economics we say that these limits or constraints reflect opportunity cost. The opportunity cost of a choice is what must be sacrificed when a choice is made. That cost may be financial; it may be measured in time, or simply the alternative foregone.

Opportunity cost: what must be sacrificed when a choice is made.

Opportunity costs play a determining role in markets. It is precisely because individuals and organizations have different opportunity costs that they enter into exchange agreements. If you are a skilled plumber and an unskilled gardener, while your neighbour is a skilled gardener and an unskilled plumber, then you and your neighbour not only have different capabilities, you also have different opportunity costs, and *you could gain by trading your skills*. Here's why. Fixing a leaking pipe has a low opportunity cost for you in terms of time: You can do it quickly. But pruning your apple trees will be costly because you must first learn how to avoid killing them and this may require many hours. Your neighbour has exactly the same problem, with the tasks in reverse positions. In a sensible world you would fix your own pipes *and* your neighbour's pipes, and she would ensure the health of the apple trees in both backyards.

If you reflect upon this 'sensible' solution—one that involves each of you achieving your objectives while minimizing the time input—you will quickly realize that it resembles the solution provided by the marketplace. You may not have a gardener as a neighbour, so you buy the services of a gardener in the marketplace. Likewise, your immediate neighbour may not need a leaking pipe repaired, but many others in your neighbourhood do, so you sell your service to them. You each specialize in the performance of specific tasks as a result of having different opportunity costs or different efficiencies. Let us now develop a model of exchange to illustrate the advantages of specialization and trade, and hence the markets that facilitate these activities. This model is developed with the help of some two-dimensional graphics.

This page titled [1.3: Opportunity cost and the market](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.4: A model of exchange and specialization

Production and specialization

We have two producers and two goods: Amanda and Zoe produce vegetables (V) and or fish (F). Their production capabilities are defined in Table 1.1 and in Figure 1.1, where the quantity of V appears on the vertical axis and the quantity of F on the horizontal axis. Zoe and Amanda each have 36-hour weeks and they devote that time to producing the two goods. But their efficiencies differ: Amanda requires two hours to produce a unit of V and three hours for a unit of F . As a consequence, if she devotes all of her time to V she can produce 18 units, or if she devotes all of her time to F she can produce 12 units. Or, she could share her time between the two. This environment can also be illustrated and analyzed graphically, as in Figure 1.1.

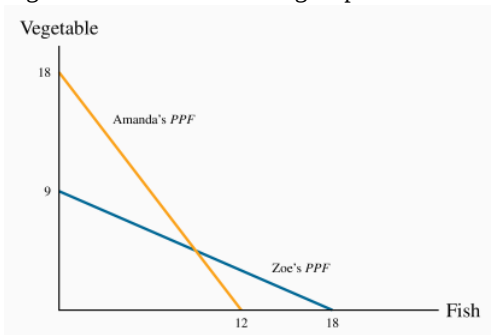
Table 1.1 Production possibilities in a two-person economy

	Hours/ fish	Hours/ vegetable	Fish specialization	Vegetable specialization
Amanda	3	2	12	18
Zoe	2	4	18	9

Each producer has a time allocation of 36 hours. By allocating total time to one activity, Amanda can produce 12 F or 18 V , Zoe can produce 18 F or 9 V . By splitting their time each person can also produce a combination of the two.

Two-dimensional graphics are a means of portraying the operation of a model, as defined above. We will use these graphical representations throughout the text. In this case, Amanda's production capability is represented by the line that meets the vertical axis at 18 and the horizontal axis at 12. The vertical point indicates that she can produce 18 units of V if she produces zero units of F – keep in mind that where V has a value of 18, Amanda has no time left for fish production. Likewise, if she devotes all of her time to fish she can produce 12 units, since each unit requires 3 of her 36 hours. The point $F=12$ is thus another possibility for her. In addition to these two possibilities, which we can term 'specialization', she could allocate her time to producing some of each good. For example, by dividing her 36 hours equally she could produce 6 units of F and 9 units of V . A little computation will quickly convince us that different allocations of her time will lead to combinations of the two goods that lie along a straight line joining the specialization points.

Figure 1.1 Absolute advantage – production



Amanda's *PPF* indicates that she can produce either 18 V (and zero F), or 12 F (and zero V), or some combination. Zoe's *PPF* indicates she can produce either 9 V (and zero F), or 18 F (and zero V), or some combination. Amanda is more efficient in producing V and Zoe is more efficient at producing F .

We will call this straight line Amanda's production possibility frontier (*PPF*): It is the combination of goods she can produce while using all of her resources – time. She could not produce combinations of goods represented by points beyond this line (to the top right). She could indeed produce combinations below it (lower left) – for example, a combination of 4 units of V and 4 units of F ; but such points would not require all of her time. The (4,4) combination would require just 20 hours. In sum, points beyond this line are not feasible, and points within it do not require all of her time resources.

Production possibility frontier (*PPF*): the combination of goods that can be produced using all of the resources available.

Having developed Amanda's *PPF*, it is straightforward to develop a corresponding set of possibilities for Zoe. If she requires 4 hours to produce a unit of V and 2 hours to produce a unit of F , then her 36 hours will enable her to specialize in 9 units of V or 18

units of F ; or she could produce a combination represented by the straight line that joins these two specialty extremes.

Consider now the opportunity costs for each person. Suppose Amanda is currently producing 18 V and zero F , and considers producing some F and less V . For each unit of F she wishes to produce, it is evident from her PPF that she must sacrifice 1.5 units of V . This is because F requires 50% more hours than V . Her trade-off is 1.5:1.0. The additional time requirement is also expressed in the intercept values: She could give up 18 units of V and produce 12 units of F instead; this again is a ratio of 1.5:1.0. This ratio defines her opportunity cost: The cost of an additional unit of F is that 1.5 units of V must be 'sacrificed'.

Applying the same reasoning to Zoe's PPF , her opportunity cost is 0.5:1; she must sacrifice one half of a unit of V to free up enough time to produce one unit of F .

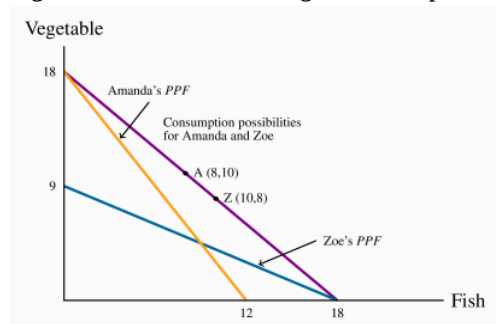
So we have established two things about Amanda and Zoe's production possibilities. First, if Amanda specializes in V she can produce more than Zoe, just as Zoe can produce more than Amanda if Zoe specializes in F . Second, their opportunity costs are different: Amanda must sacrifice more V than Zoe in producing one more unit of F . The different opportunity costs translate into potential gains for each individual.

The gains from exchange

We shall illustrate the gains that arise from specialization and exchange graphically. Note first that if these individuals are self-sufficient, in the sense that they consume their own production, each individual's consumption combination will lie on their own PPF . For example, Amanda could allocate half of her time to each good, and produce (and consume) 6 F and 9 V . Such a point necessarily lies on her PPF . Likewise for Zoe. So, *in the absence of exchange*, each individual's PPF is also her consumption possibility frontier (CPF). In Figure 1.1 the PPF for each individual is thus also her CPF .

Consumption possibility frontier (CPF): the combination of goods that can be consumed as a result of a given production choice.

Figure 1.2 Absolute advantage – consumption



With specialization and trade at a rate of 1:1 they consume along the line joining the specialization points. If Amanda trades 8 V to Zoe in return for 8 F , Amanda moves to the point A(8,10) and Zoe to Z(10,8). Each can consume more after specialization than before specialization.

Upon realizing that they are not equally efficient in producing the two goods, they decide to specialize completely in producing just the single good where they are most efficient. Amanda specializes in V and Zoe in F . Next they must agree to a rate at which to exchange V for F . Since Amanda's opportunity cost is 1.5:1 and Zoe's is 0.5:1, suppose they agree to exchange V for F at an intermediate rate of 1:1. There are many trading, or exchange, rates possible; our purpose is to illustrate that gains are possible for *both* individuals at some exchange rate. The choice of this rate also makes the graphic as simple as possible. At this exchange rate, 18 V must exchange for 18 F . In Figure 1.2, this means that each individual is now able to consume along the line joining the coordinates (0,18) and (18,0).¹ This is because Amanda produces 18 V and she can trade at a rate of 1:1, while Zoe produces 18 F and trades at the same rate of 1:1.

The fundamental result illustrated in Figure 1.2 is that, as a result of specialization and trade, each individual can consume combinations of goods that lie on a line beyond her initial consumption possibilities. Their consumption well-being has thus improved. For example, suppose Amanda trades away 8 V to Zoe and obtains 8 F in return. The points 'A' and 'Z' with coordinates (8,10) and (10,8) respectively define their final consumption. Pre-specialization, if Amanda wished to consume 8 F she would have been constrained to consume 6 V rather than the 10 V now possible. Zoe benefits correspondingly.²

The foregoing example illustrates that trade is not a zero-sum game; it has a positive net value because both parties to the trade can gain. A zero-sum gain is where the gains to one party exactly offset the losses to another. This is an extraordinarily important principle in trade negotiations, whether international or domestic.

A **zero-sum game** is an interaction where the gain to one party equals the loss to another party.

Market design

In the preceding example we have shown that specialization provides scope for gains that can accrue to those participating in the exchange. But this tells us little about how a market for these products comes into being: how does the exchange take place, and how is information transmitted? The answer is that while some markets have evolved historically to their current state, many markets are designed by an institution or a firm. Fruit and vegetable markets have been with us for thousands of years - since we ceased being purely a hunter-gatherer society. They exist in every community in the world economy. In contrast, the Dutch tulip auction was designed in the early 1600s and exists in basically the same form to this day: the auctioneer begins with a high price, lowers it at known time intervals (measured in seconds or minutes) until some buyer signals that she is willing to purchase the lot on offer. Supermarkets in contrast offer goods at a fixed price. Government contracts are normally signed after a tendering process, in which interested suppliers submit bids. Amazon Inc. is currently experimenting with cashierless 'bricks and mortar' stores that monitor all transactions electronically. Craig's List and E-Bay have their own sets of rules.

In each of these cases markets are designed, frequently with a specific objective on the part of the supplier or the mediating institution: Amazon wants to increase its share of all goods trades; governments wish to limit costs. Markets do not all grow spontaneously and the structure of a market will influence how the gains from trade are distributed.

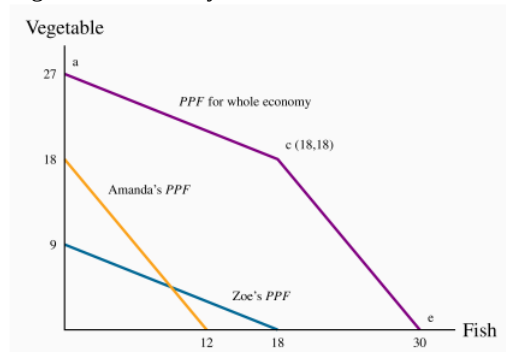
This page titled [1.4: A model of exchange and specialization](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.5: Economy-wide production possibilities

The *PPFs* in Figures 1.1 and 1.2 define the amounts of the goods that each *individual* can produce while using all of their fixed productive capacity—time in this instance. The national, or economy-wide, *PPF* for this two-person economy reflects these individual possibilities combined. Such a frontier can be constructed using the individual frontiers as the component blocks.

First let us define this economy-wide frontier precisely. The economy-wide *PPF* is the set of goods and services combinations that can be produced in the economy when all available productive resources are in use. Figure 1.3 contains both of the individual frontiers plus the aggregate of these, represented by the kinked line *ace*. The point on the *V* axis, $a=27$, represents the total amount of *V* that could be produced if both individuals devoted all of their time to it. The point $e=30$ on the horizontal axis is the corresponding total for fish.

Figure 1.3 Economy-wide PPF



From *a*, to produce Fish it is more efficient to use Zoe because her opportunity cost is less (segment *ac*). When Zoe is completely specialized, Amanda produces (*ce*). With complete specialization this economy can produce 27*V* or 30*F*.

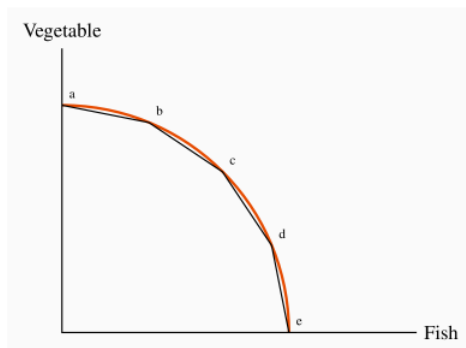
Economy-wide *PPF*: the set of goods and services combinations that can be produced in the economy when all available productive resources are in use.

To understand the point *c*, imagine initially that all resources are devoted to *V*. From such a point, *a*, consider a reduction in *V* and an increase in *F*. The most efficient way of increasing *F* production at the point *a* is to use the individual whose opportunity cost is lower. Zoe can produce one unit of *F* by sacrificing just 0.5 units of *V*, whereas Amanda must sacrifice 1.5 units of *V* to produce 1 unit of *F*. Hence, at this stage Amanda should stick to *V* and Zoe should devote some time to fish. In fact as long as we want to produce more fish Zoe should be the one to do it, until she has exhausted her time resource. This occurs after she has produced 18*F* and has ceased producing *V*. At this point the economy will be producing 18*V* and 18*F* – the point *c*.

From this combination, if the economy wishes to produce more fish Amanda must become involved. Since her opportunity cost is 1.5 units of *V* for each unit of *F*, the next segment of the economy-wide *PPF* must see a reduction of 1.5 units of *V* for each additional unit of *F*. This is reflected in the segment *ce*. When both producers allocate all of their time to *F* the economy can produce 30 units. Hence the economy's *PPF* is the two-segment line *ace*. Since this has an outward kink, we call it concave (rather than convex).

As a final step consider what this *PPF* would resemble if the economy were composed of many persons with differing efficiencies. A little imagination suggests (correctly) that it will have a segment for each individual and continue to have its outward concave form. Hence, a four-person economy in which each person had a different opportunity cost could be represented by the segmented line *abcde*, in Figure 1.4. Furthermore, we could represent the *PPF* of an economy with a very large number of such individuals by a somewhat smooth *PPF* that accompanies the 4-person *PPF*. The logic for its shape continues to be the same: As we produce less *V* and more *F* we progressively bring into play resources, or individuals, whose opportunity cost, in terms of reduced *V* is higher.

Figure 1.4 A multi-person PPF



The PPF for the whole economy, abcde, is obtained by allocating productive resources most efficiently. With many individuals we can think of the PPF as the *concave envelope* of the individual capabilities.

The outputs V and F in our economic model require just one input – time, but if other productive resources were required the result would be still a concave *PPF*. Furthermore, we generally interpret the *PPF* to define the output possibilities *when the economy is running at its normal capacity*. In this example, we consider a work week of 36 hours to be the 'norm'. Yet it is still possible that the economy's producers might work some additional time in exceptional circumstances, and this would increase total production possibilities. This event would be represented by an outward movement of the *PPF*.

This page titled [1.5: Economy-wide production possibilities](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.6: New Page

The *PPF* can be used to illustrate several aspects of macroeconomics: In particular, the level of an economy's output, the growth of national and per capita output over time, and short-run business-cycle fluctuations in national output and employment.

Aggregate output

An economy's capacity to produce goods and services depends on its endowment of resources and the productivity of those resources. The two-person, two-product examples in the previous section reflect this.

The productivity of labour, defined as output per worker or per hour, depends on:

- The skill, knowledge and experience of the labour force;
- The capital stock: Buildings, machinery, equipment, and software the labour force has to work with; and
- The current state of technology.

The **productivity of labour** is the output of goods and services per worker.

An economy's **capital stock** is the buildings, machinery, equipment and software used in producing goods and services.

The economy's output, which we define by Y , can be defined as the output per worker times the number of workers; hence, we can write:

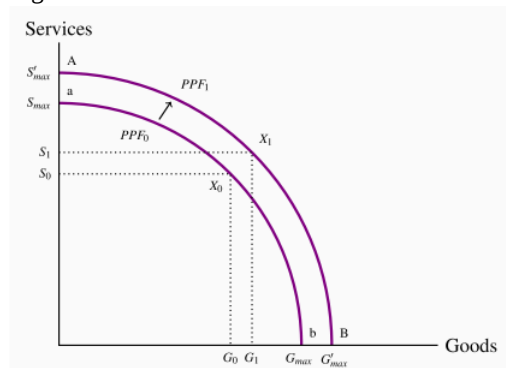
$$Y = (\text{number of workers employed}) \times (\text{output per worker}).$$

When the employment of labour corresponds to 'full employment' in the sense that everyone willing to work at current wage rates and normal hours of work is working, the economy's actual output is also its capacity output Y_c . We also term this capacity output as full employment output:

Full employment output $Y_c = (\text{number of workers at full employment}) \times (\text{output per worker})$.

Suppose the economy is operating with full employment of resources producing outputs of two types: Goods and services. In Figure 1.5, PPF_0 shows the different combinations of goods and services the economy can produce in a particular year using all its labour, capital and the best technology available at the time.

Figure 1.5 Growth and the PPF



Economic growth is illustrated by an outward shift in the *PPF* from PPF_0 to PPF_1 . PPF_1 shows the economy can produce more in both sectors than with PPF_0 .

An aggregate economy produces a large variety of outputs in two broad categories. Goods are the products of the agriculture, forestry, mining, manufacturing and construction industries. Services are provided by the wholesale and retail trade, transportation, hospitality, finance, health care, education, legal and other service sectors. As in the two-product examples used earlier, the shape of the *PPF* illustrates the opportunity cost of increasing the output of either product type. We are not concerned with who supplies the products for the moment: It may be the private sector or the government.

Point X_0 on PPF_0 shows one possible structure of capacity output. This combination may reflect the pattern of demand and hence expenditures in this economy. Output structures and therefore the shapes of *PPFs* differ among economies with different income levels. High-income economies spend more on services than goods and produce higher ratios of services to goods. Middle income countries produce lower ratios of services to goods, and low income countries much lower ratios of services to goods. For example,

in 2017, the structure of national output in Canada was 70 percent services and 30 percent goods, while in Mexico the structure was 48 percent services and 52 percent goods.

Different countries also have different *PPFs* and different output structures, depending on their resource endowments, labour forces, capital stocks, technology and expenditure patterns.

Economic growth

Three things contribute to growth in the economy. The labour supply grows as the population expands; the stock of capital grows as spending by business (and government) on buildings, machinery, information technology and so forth increases; and labour-force productivity grows as a result of experience, the development of scientific knowledge combined with product and process innovations, and advances in the technology of production. Combined, these developments expand capacity output over time. In Figure 1.5 economic growth shifts the *PPF* out from PPF_0 to PPF_1 .

This basic description covers the key sources of growth in total output. Economies differ in their rates of overall economic growth as a result of different rates of growth in labour force, in capital stock, and improvements in technology. But improvements in standards of living require more than growth in total output. Increases in output *per worker* and *per person* are necessary. Sustained increases in living standards require sustained growth in labour productivity, which in turn is based on advances in the technology along with the amount of capital each worker has to work with. Furthermore, if the growth in output is to benefit society at large, workers across the board need to see an increase in their earnings. As we shall explore in Chapter 13, several developed countries have seen the fruits of growth concentrated in the hands of the highest income earners.

Recessions and booms

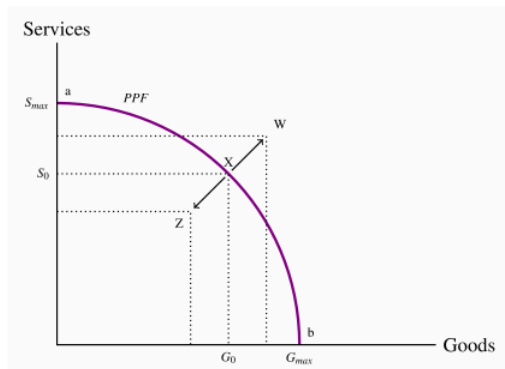
A prime objective of economic policy is to ensure that the economy operates on or near the *PPF* – it should use its resources to capacity and have minimal unemployment. However, economic conditions are seldom tranquil for long periods of time. Unpredictable changes in business expectations of future profits, in consumer confidence, in financial markets, in commodity and energy prices, in trade agreements and disputes, in economic conditions in major trading partners, in government policy and many other events disrupt patterns of expenditure and output. Some of these changes disturb the level of total expenditure and thus the demand for total output. Others disturb the conditions of production and thus the economy's production capacity. Whatever the exact cause, the economy may be pushed off its current *PPF*. If expenditures on goods and services decline, the economy may experience a recession. Output would fall short of capacity output and unemployment would rise. Alternatively, times of rapidly growing expenditure and output may result in an economic boom: Output and employment expand beyond capacity levels.

An **economic recession** occurs when output falls below the economy's capacity output.

A **boom** is a period of high growth that raises output above normal capacity output.

Recent history provides examples. Following the financial crisis of 2008-09 that hit the US and many other developed economies, many economies were pushed into recessions. Expenditure on new residential construction collapsed for lack of income and secure financing, as did business investment, consumption spending and exports. Lower expenditures reduced producers' revenues, forcing cuts in output and employment and reducing household incomes. Lower incomes led to further cutbacks in spending. In Canada in 2009 aggregate output declined by 2.9 percent, employment declined by 1.6 percent and the unemployment rate rose from 6.1 percent in 2008 to 8.3 percent by 2010. The world's economies have been slow to recover, and even by 2019 the output in several developed economies was no higher than it was in 2008. Canada's recession was not nearly as severe as the recessions in economies such as Spain, Italy and Greece; but output between 2009 and 2019 has been below the potential of the Canadian economy. In the third quarter of 2019 the national output was about 0.7 percent below potential output and the unemployment rate was 5.5 percent.

Figure 1.6 Booms and recessions



Economic recessions leave the economy below its normal capacity; the economy might be driven to a point such as Z. Economic expansions, or booms, may drive capacity above its normal level, to a point such as W.

An economy in a recession is operating inside its *PPF*. The fall in output from X to Z in Figure 1.6 illustrates the effect of a recession. Expenditures on goods and services have declined. Output is less than capacity output, unemployment is up and some plant capacity is idle. Labour income and business profits are lower. More people would like to work and business would like to produce and sell more output, but it takes time for interdependent product, labour and financial markets in the economy to adjust and increase employment and output. Monetary and fiscal policy may be productive in specific circumstances, to stimulate demand, increase output and employment and move the economy back to capacity output and full employment. The development and implementation of such policies form the core of macroeconomics.

Alternatively, an unexpected increase in demand for exports would increase output and employment. Higher employment and output would increase incomes and expenditure, and in the process spread the effects of higher output sales to other sectors of the economy. The economy would move outside its *PPF*, for example to W in Figure 1.6, by using its resources more intensively than normal. Unemployment would fall and overtime work would increase. Extra production shifts would run plant and equipment for longer hours and work days than were planned when it was designed and installed. Output at this level may not be sustainable, because shortages of labour and materials along with excessive rates of equipment wear and tear would push costs and prices up. Again, we will examine how the economy reacts to such a state in our macroeconomic analysis.

Output and employment in the Canadian economy over the past twenty years fluctuated about growth trend in the way Figure 1.6 illustrates. For several years prior to 2008 the Canadian economy operated slightly above its capacity; but once the recession arrived monetary and fiscal policy were used to fight it – to bring the economy back from a point such as Z towards a point such as X on the *PPF*.

Macroeconomic models and policy

The *PPF* diagrams illustrate the main dimensions of macroeconomics: Capacity output, growth in capacity output and business cycle fluctuations in actual output relative to capacity. But these diagrams do not offer explanations and analysis of macroeconomic activity. We need a macroeconomic *model* to understand and evaluate the causes and consequences of business cycle fluctuations. As we shall see, these models are based on explanations of expenditure decisions by households and business, financial market conditions, production costs and producer pricing decisions at different levels of output. Models also capture the objectives of fiscal and monetary policies and provide a framework for policy evaluation. A full macroeconomic model integrates different sector behaviours and the feedbacks across sectors that can moderate or amplify the effects of changes in one sector on national output and employment.

This page titled 1.6: New Page is shared under a [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.7: Conclusion

We have covered a lot of ground in this introductory chapter. It is intended to open up the vista of economics to the new student in the discipline. Economics is powerful and challenging, and the ideas we have developed here will serve as conceptual foundations for our exploration of the subject.

This page titled [1.7: Conclusion](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.8: Key Terms

Macroeconomics studies the economy as system in which linkages and feedbacks among sectors determine national output, employment and prices.

Microeconomics is the study of individual behaviour in the context of scarcity.

Mixed economy: goods and services are supplied both by private suppliers and government.

Model is a formalization of theory that facilitates scientific inquiry.

Theory is a logical view of how things work, and is frequently formulated on the basis of observation.

Opportunity cost of a choice is what must be sacrificed when a choice is made.

Production possibility frontier (PPF) defines the combination of goods that can be produced using all of the resources available.

Consumption possibility frontier (CPF): the combination of goods that can be consumed as a result of a given production choice.

A **zero-sum game** is an interaction where the gain to one party equals the loss to another party.

Economy-wide PPF is the set of goods combinations that can be produced in the economy when all available productive resources are in use.

Productivity of labour is the output of goods and services per worker.

Capital stock: the buildings, machinery, equipment and software used in producing goods and services.

Full employment output $Y_c = (\text{number of workers at full employment}) \times (\text{output per worker})$. **Recession:** when output falls below the economy's capacity output. **Boom:** a period of high growth that raises output above normal capacity output.

This page titled [1.8: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

1.9: Exercises for Chapter 1

EXERCISE 1.1

An economy has 100 identical workers. Each one can produce four cakes or three shirts, regardless of the number of other individuals producing each good.

- How many cakes can be produced in this economy when all the workers are cooking?
- How many shirts can be produced in this economy when all the workers are sewing?
- On a diagram with cakes on the vertical axis, and shirts on the horizontal axis, join these points with a straight line to form the *PPF*.
- Label the inefficient and unattainable regions on the diagram.

EXERCISE 1.2

In the table below are listed a series of points that define an economy's production possibility frontier for goods Y and X.

Y	1000	900	800	700	600	500	400	300	200	100	0
X	0	1600	2500	3300	4000	4600	5100	5500	5750	5900	6000

- Plot these pairs of points to scale, on graph paper, or with the help of a spreadsheet.
- Given the shape of this *PPF* is the economy made up of individuals who are similar or different in their production capabilities?
- What is the opportunity cost of producing 100 more Y at the combination ($X=5500, Y=300$).
- Suppose next there is technological change so that at every output level of good Y the economy can produce 20 percent more X. Enter a new row in the table containing the new values, and plot the new *PPF*.

EXERCISE 1.3

Using the *PPF* that you have graphed using the data in Exercise 1.2, determine if the following combinations are attainable or not: ($X=3000, Y=720$), ($X=4800, Y=480$).

EXERCISE 1.4

You and your partner are highly efficient people. You can earn \$20 per hour in the workplace; your partner can earn \$30 per hour.

- What is the opportunity cost of one hour of leisure for you?
- What is the opportunity cost of one hour of leisure for your partner?
- Now consider what a *PPF* would look like: You can produce/consume two things, leisure and income. Since income buys things you can think of the *PPF* as having these two 'products' – leisure and consumption goods/services. So, with leisure on the horizontal axis and income in dollars is on the vertical axis, plot your *PPF*. You can assume that you have 12 hours per day to allocate to either leisure or income. [Hint: the leisure axis will have an intercept of 12 hours. The income intercept will have a dollar value corresponding to where all hours are devoted to work.]
- Draw the *PPF* for your partner.

EXERCISE 1.5

Louis and Carrie Anne are students who have set up a summer business in their neighbourhood. They cut lawns and clean cars. Louis is particularly efficient at cutting the grass – he requires one hour to cut a typical lawn, while Carrie Anne needs one and one half hours. In contrast, Carrie Anne can wash a car in a half hour, while Louis requires three quarters of an hour.

- If they decide to specialize in the tasks, who should cut the grass and who should wash cars?
- If they each work a twelve hour day, how many lawns can they cut and how many cars can they wash if they each specialize in performing the task where they are most efficient?
- Illustrate the *PPF* for each individual where lawns are on the horizontal axis and car washes on the vertical axis, if each individual has twelve hours in a day.

EXERCISE 1.6

Continuing with the same data set, suppose Carrie Anne's productivity improves so that she can now cut grass as efficiently as Louis; that is, she can cut grass in one hour, and can still wash a car in one half of an hour.

- a. In a new diagram draw the *PPF* for each individual.
- b. In this case does specialization matter if they are to be as productive as possible as a team?
- c. Draw the *PPF* for the whole economy, labelling the intercepts and the 'kink' point coordinates.

EXERCISE 1.7

Going back to the simple *PPF* plotted for Exercise 1.1 where each of 100 workers can produce either four cakes or three shirts, suppose a recession reduces demand for the outputs to 220 cakes and 129 shirts.

- a. Plot this combination of outputs in the diagram that also shows the *PPF*.
 - b. How many workers are needed to produce this output of cakes and shirts?
 - c. What percentage of the 100 worker labour force is unemployed?
1. When two values, separated by a comma, appear in parentheses, the first value refers to the horizontal-axis variable, and the second to the vertical-axis variable.
 2. In the situation we describe above one individual is absolutely more efficient in producing one of the goods and absolutely less efficient in the other. We will return to this model in Chapter 15 and illustrate that consumption gains of the type that arise here can also result if one of the individuals is absolutely more efficient in producing both goods, but that the degree of such advantage differs across goods.

This page titled [1.9: Exercises for Chapter 1](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

2: Theories, data and beliefs

Chapter 2: Theories, data and beliefs

In this chapter we will explore:

2.1	Data analysis
2.2	Data, theory and economic models
2.3	Ethics, efficiency and beliefs

Economists, like other scientists and social scientists, observe and analyze behaviour and events. Economists are concerned primarily with the economic causes and consequences of what they observe. They want to understand an extensive range of human experience, including: money, government finances, industrial production, household consumption, inequality in income distribution, war, monopoly power, health, professional and amateur sports, pollution, marriage, the arts, and much more.

Economists approach these issues using theories and models. To present, explain, illustrate and evaluate their theories and models they have developed a set of techniques or tools. These involve verbal descriptions and explanations, diagrams, algebraic equations, data tables and charts and statistical tests of economic relationships.

This chapter covers some of these basic techniques of analysis.

2.1 Data analysis

The analysis of behaviour necessarily involves data. Data may serve to validate or contradict a theory. Data analysis, even without being motivated by economic theory, frequently displays patterns of behaviour that merit examination. The terms *variables* and *data* are related. Variables are measures that can take on different magnitudes. The interest rate on a student loan, for example, is a variable with a certain value at a point in time but perhaps a different value at an earlier or later date. Economic theories and models explain the causal relationships between variables. In contrast, Data are the recorded values of variables. Sets of data provide specific values for the variables we want to study and analyze. Knowing that gross domestic product (a variable) declined in 2009 is just a partial description of events. If the data indicate that it decreased by exactly 3%, we know a great deal more – the decline was large.

Variables: measures that can take on different values.

Data: recorded values of variables.

Sets of data help us to test our models or theories, but first we need to pay attention to the economic logic involved in observations and modelling. For example, if sunspots or baggy pants were found to be correlated with economic expansion, would we consider these events a coincidence or a key to understanding economic growth? The observation is based on facts or data, but it need not have any economic content. The economist's task is to distinguish between coincidence and economic causation. Merely because variables are associated or correlated does not mean that one causes the other.

While the more frequent wearing of loose clothing in the past may have been associated with economic growth because they both occurred at the same time (correlation), one could not argue on a logical basis that this behaviour causes good economic times. Therefore, the past association of these variables should be considered as no more than a coincidence. Once specified on the basis of economic logic, a model must be tested to determine its usefulness in explaining observed economic events.

Table 2.1 House prices and price indexes

Year	House prices in dollars (P_H)	Percentage change in P_H	Percentage change in consumer prices	Real percentage change in P_H	Index for price of housing	5-year mortgage rate
2001	350,000				100	7.75
2002	360,000				102.9	6.85

2003	395,000	35,000/360,000=9.7%	3%	6.7%	112.9	6.6
2004	434,000				124.0	5.8
2005	477,000				136.3	6.1
2006	580,000				165.7	6.3
2007	630,000				180.0	6.65
2008	710,000				202.9	7.3
2009	605,000	-105,000/710,000=-14.8%	1.6%	-16.4%	172.9	5.8
2010	740,000				211.4	5.4
2011	800,000				228.6	5.2

Note: Data on changes in consumer prices come from Statistics Canada, CANSIM series V41692930; data on house prices are for N. Vancouver from *Royal Le Page*; data on mortgage rates from <http://www.ratehub.ca>. Index for house prices obtained by scaling each entry in column 2 by 100/350,000. The real percentage change in the price of housing is: The percentage change in the price of housing minus the percentage change in consumer prices.

Data types

Data come in several forms. One form is time-series, which reflects a set of measurements made in sequence at different points in time. The first column in Table 2.1 reports the values for house prices in North Vancouver for the first quarter of each year, between 2001 and 2011. Evidently this is a time series. Annual data report one observation per year. We could, alternatively, have presented the data in monthly, weekly, or even daily form. The frequency we use depends on the purpose: If we are interested in the longer-term trend in house prices, then the annual form suffices. In contrast, financial economists, who study the behaviour of stock prices, might not be content with daily or even hourly prices; they may need prices minute-by-minute. Such data are called high-frequency data, whereas annual data are low-frequency data.

Table 2.2 Unemployment rates, Canada and Provinces, monthly 2012, seasonally adjusted

	Jan	Feb	Mar	Apr	May	Jun
CANADA	7.6	7.4	7.2	7.3	7.3	7.2
NFLD	13.5	12.9	13.0	12.3	12.0	13.0
PEI	12.2	10.5	11.3	11.0	11.3	11.3
NS	8.4	8.2	8.3	9.0	9.2	9.6
NB	9.5	10.1	12.2	9.8	9.4	9.5
QUE	8.4	8.4	7.9	8.0	7.8	7.7
ONT	8.1	7.6	7.4	7.8	7.8	7.8
MAN	5.4	5.6	5.3	5.3	5.1	5.2
SASK	5.0	5.0	4.8	4.9	4.5	4.9
ALTA	4.9	5.0	5.3	4.9	4.5	4.6
BC	6.9	6.9	7.0	6.2	7.4	6.6

Source: Statistics Canada CANSIM Table 282-0087.

Time-series: a set of measurements made sequentially at different points in time.

High (low) frequency data: series with short (long) intervals between observations.

In contrast to time-series data, cross-section data record the values of different variables at a point in time. Table 2.2 contains a cross-section of unemployment rates for Canada and Canadian provinces economies. For January 2012 we have a snapshot of the

provincial economies at that point in time, likewise for the months until June. This table therefore contains repeated cross-sections. When the unit of observation is the same over time such repeated cross sections are called longitudinal data. For example, a health survey that followed and interviewed the same individuals over time would yield longitudinal data. If the individuals differ each time the survey is conducted, the data are repeated cross sections. Longitudinal data therefore follow the same units of observation through time.

Cross-section data: values for different variables recorded at a point in time.

Repeated cross-section data: cross-section data recorded at regular or irregular intervals.

Longitudinal data: follow the same units of observation through time.

Graphing the data

Data can be presented in graphical as well as tabular form. Figure 2.1 plots the house price data from the second column of Table 2.1. Each asterisk in the figure represents a price value and a corresponding time period. The horizontal axis reflects time, the vertical axis price in dollars. The graphical presentation of data simply provides a visual rather than numeric perspective. It is immediately evident that house prices increased consistently during this 11-year period, with a single downward 'correction' in 2009. We have plotted the data a second time in Figure 2.2 to illustrate the need to read graphs carefully. The greater *apparent* slope in Figure 2.1 might easily be interpreted to mean that prices increased more steeply than suggested in Figure 2.2. But a careful reading of the axes reveals that this is not so; using different scales when plotting data or constructing diagrams can mislead the unaware viewer.

Figure 2.1 House prices in dollars 1999-2012

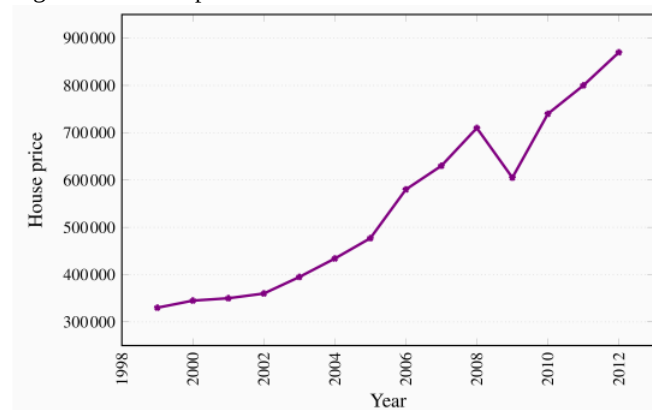
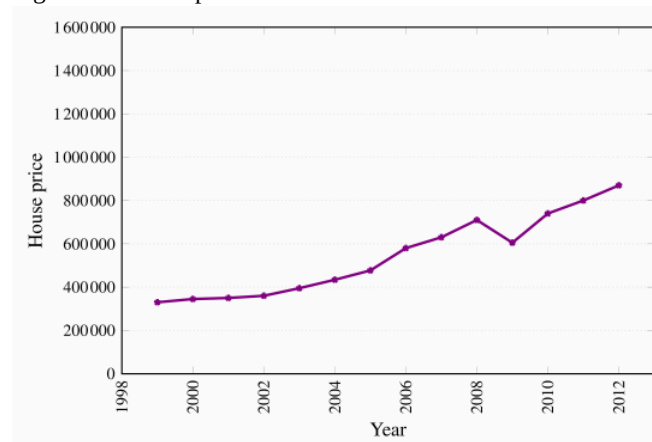


Figure 2.2 House prices in dollars 1999-2012



Percentage changes

The use of percentages makes the analysis of data particularly simple. Suppose we wanted to compare the prices of New York luxury condominiums with the prices of homes in rural Mississippi. In the latter case, a change in average prices of \$10,000 might be considered enormous, whereas a change of one million dollars in New York might be pretty normal – because the average price

in New York is so much higher than in Mississippi. To make comparisons between the two markets, we can use the concept of a percentage change. This is defined as the change in the value of the variable, relative to its initial value, multiplied by 100.

Percentage change = $[(\text{change in values})/(\text{original value})] \times 100$.

The third column of Table 2.1 contains the values of the percentage change in house prices for two pairs of years. Between 2002 and 2003 the price change was \$35,000. Relative to the price in the first of these two years this change was the fraction $35,000/395,000=0.097$. If we multiply this fraction by 100 we obtain a percentage price change of 9.7%. Evidently we could calculate the percentage price changes for all pairs of years. A second price change is calculated for the 2008-2009 pair of years. Here price declined and the result is thus a negative percentage change.

Consumer prices

Most variables in economics are averages of the components that go into them. When variables are denominated in dollar terms it is important to be able to interpret them correctly. While the house price series above indicates a strong pattern of price increases, it is vital to know if the price of housing increased more or less rapidly than other prices in the economy. If all prices in the economy were increasing in line with house prices there would be no special information in the house price series. However, if house prices increased more rapidly than prices in general, then the data indicate that something special took place in the housing market during the decade in question. To determine an answer to this we need to know the degree to which the general price level changed each year.

Statistics Canada regularly surveys the price of virtually every product produced in the economy. One such survey records the prices of goods and services purchased by consumers. Statistics Canada then computes an average price level for all of these goods combined for each time period the survey is carried out (monthly). Once Statistics Canada has computed the average consumer price, it can compute the change in the price level from one period to the next. In Table 2.1 two such values are entered in the following data column: Consumer prices increased by 3% between 2002 and 2003, and by 1.6% between 2008 and 2009. These percentage changes in the general price level represent inflation if prices increase, and deflation if prices decline.

In this market it is clear that housing price changes were substantially larger than the changes in consumer prices for these two pairs of years. The next column provides information on the difference between the house price changes and changes in the general consumer price level, in percentage terms. This is (approximately) the change in the relative price of housing, or what economists call the real price of housing. It is obtained by subtracting the rate of change in the general price index from the rate of change in the variable of interest.

Consumer price index: the average price level for consumer goods and services.

Inflation (deflation) rate: the annual percentage increase (decrease) in the level of consumer prices.

Real price: the actual price adjusted by the general (consumer) price level in the economy.

Index numbers

Statistics Canada and other statistical agencies frequently present data in index number form. An index number provides an easy way to read the data. For example, suppose we wanted to compute the percentage change in the price of housing between 2001 and 2007. We could do this by entering the two data points in a spreadsheet or calculator and do the computation. But suppose the prices were entered in another form. In particular, by dividing each price value by the first year value and multiplying the result by 100 we obtain a series of prices that are all relative to the initial year – which we call the base year. The resulting series in column 6 of Table 2.1 is an index of house price values. Each entry is the corresponding value in column 2, divided by the first entry in column 2. For example, the value 124.0 in row 4 is obtained as $(434,000/350,000) \times 100 = 124.0$. The key characteristics of indexes are that they are *not dependent upon the units of measurement of the data in question*, and they are interpretable easily with reference to a given base value. To illustrate, suppose we wish to know how prices behaved between 2001 and 2007. The index number column immediately tells us that prices increased by 80%, because relative to 2001, the 2007 value is 80% higher.

Index number: value for a variable, or an average of a set of variables, expressed relative to a given base value.

Furthermore, index numbers enable us to make *comparisons with the price patterns for other goods* much more easily. If we had constructed a price index for automobiles, which also had a base value of 100 in 2001, we could make immediate comparisons without having to compare one set of numbers defined in thousands of dollars with another defined in hundreds of thousands of dollars. In short, index numbers simplify the interpretation of data.

2.2 Data, theory and economic models

Let us now investigate the interplay between economic theories on the one hand and data on the other. We will develop two examples. The first will be based upon the data on house prices, the second upon a new data set.

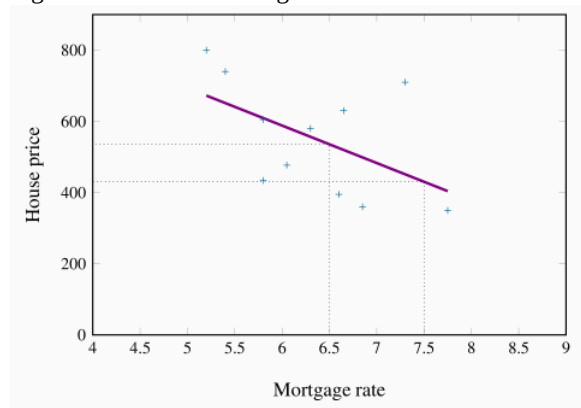
House prices – theory

Remember from Chapter 1 that a theory is a logical argument regarding economic relationships. A theory of house prices would propose that the price of housing depends upon a number of elements in the economy. In particular, if borrowing costs are low then buyers are able to afford the interest costs on larger borrowings. This in turn might mean they are willing to pay higher prices. Conversely, if borrowing rates are higher. Consequently, the borrowing rate, or mortgage rate, is a variable for an economic model of house prices. A second variable might be available space for development: If space in a given metropolitan area is tight then the land value will reflect this, and consequently the higher land price should be reflected in higher house prices. A third variable would be the business climate: If there is a high volume of high-value business transacted in a given area then buildings will be more in demand, and that in turn should be reflected in higher prices. For example, both business and residential properties are more highly priced in San Francisco and New York than in Moncton, New Brunswick. A fourth variable might be environmental attractiveness: Vancouver may be more enticing than other towns in Canada. A fifth variable might be the climate.

House prices – evidence

These and other variables could form the basis of a *theory* of house prices. A *model* of house prices, as explained in Chapter 1, focuses upon what we would consider to be the most important subset of these variables. In the limit, we could have an extremely simple model that specified a dependence between the price of housing and the mortgage rate alone. To test such a simple model we need data on house prices and mortgage rates. The final column of Table 2.1 contains data on the 5-year fixed-rate mortgage for the period in question. Since our simple model proposes that prices depend (primarily) upon mortgage rates, in Figure 2.3 we plot the house price series on the vertical axis, and the mortgage rate on the horizontal axis, for each year from 2001 to 2011. As before, each point (shown as a '+') represents a pair of price and mortgage rate values.

Figure 2.3 Price of housing



The resulting plot (called a scatter diagram) suggests that there is a negative relationship between these two variables. That is, higher prices are correlated with lower mortgage rates. Such a correlation is consistent with our theory of house prices, and so we might conclude that changes in mortgage rates *cause* changes in house prices. Or at least the data suggest that we should not reject the idea that such causation is in the data.

House prices – inference

To summarize the relationship between these variables, the pattern suggests that a straight line through the scatter plot would provide a reasonably good description of the relationship between these variables. Obviously it is important to define the most appropriate line – one that 'fits' the data well.¹ The line we have drawn through the data points is informative, because it relates the two variables in a *quantitative manner*. It is called a regression line. It predicts that, on average, if the mortgage rate increases, the price of housing will respond in the downward direction. This particular line states that a one point change in the mortgage rate will move prices in the opposing direction by \$105,000. This is easily verified by considering the dollar value corresponding to say a mortgage value of 6.5, and then the value corresponding to a mortgage value of 7.5. Projecting vertically to the regression line from each of these points on the horizontal axis, and from there across to the vertical axis will produce a change in price of \$105,000.

Note that the line is not at all a 'perfect' fit. For example, the mortgage rate declined between 2008 and 2009, but the price declined also – contrary to our theory. The model is not a perfect predictor; it states that *on average* a change in the magnitude of the x-axis variable leads to a change of a specific amount in the magnitude of the y-axis variable.

In this instance the slope of the line is given by $-105,000/1$, which is the vertical distance divided by the corresponding horizontal distance. Since the line is straight, this slope is unchanging.

Regression line: representation of the average relationship between two variables in a scatter diagram.

Road fatalities – theory, evidence and inference

Table 2.3 contains data on annual road fatalities per 100,000 drivers for various age groups. In the background, we have a *theory*, proposing that driver fatalities depend upon the age of the driver, the quality of roads and signage, speed limits, the age of the automobile stock and perhaps some other variables. Our model focuses upon a subset of these variables, and in order to present the example in graphical terms we specify fatalities as being dependent upon a single variable – age of driver.

Table 2.3 Non-linearity: Driver fatality rates Canada, 2009

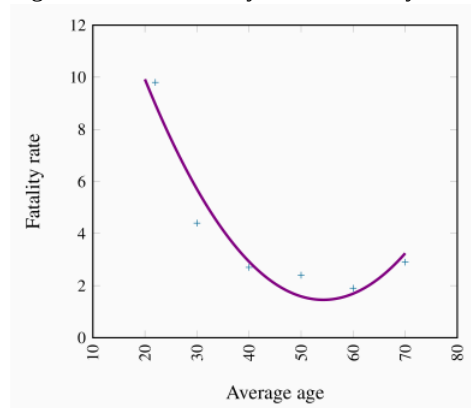
Age of driver	Fatality rate per 100,000 drivers
20-24	9.8
25-34	4.4
35-44	2.7
45-54	2.4
55-64	1.9
65+	2.9

Source: Transport Canada, Canadian motor vehicle traffic collision statistics, 2009.

The scatter diagram is presented in Figure 2.4. Two aspects of this plot stand out. First, there is an exceedingly steep decline in the fatality rate when we go from the youngest age group to the next two age groups. The decline in fatalities between the youngest and second youngest groups is 5.4 points, and between the second and third age groups is 1.7 points. The decline between the third and fourth groups is minimal - just 0.3 points. Hence, behaviour is not the same throughout the age distribution. Second, we notice that fatalities increase for the oldest age group, perhaps indicating that the oldest drivers are not as good as middle-aged drivers.

These two features suggest that the relationship between fatalities and age differs across the age spectrum. Accordingly, a straight line would not be an accurate way of representing the behaviours in these data. A straight line through the plot implies that a given change in age should have a similar impact on fatalities, no matter the age group. Accordingly we have an example of a *non-linear relationship*. Such a non-linear relationship might be represented by the curve going through the plot. Clearly the slope of this line varies as we move from one age category to another.

Figure 2.4 Non-linearity: Driver fatality rates Canada, 2009



Fatality rates vary non-linearly with age: At first they decline, then increase again, relative to the youngest age group.

2.3 Ethics, efficiency and beliefs

Positive economics studies objective or scientific explanations of how the economy functions. Its aim is to understand and generate predictions about how the economy may respond to changes and policy initiatives. In this effort economists strive to act as detached scientists, regardless of political sympathies or ethical code. Personal judgments and preferences are (ideally) kept apart. In this particular sense, economics is similar to the natural sciences such as physics or biology. To date in this chapter we have been exploring economics primarily from a positive standpoint.

In contrast, normative economics offers recommendations based partly on value judgments. While economists of different political persuasions can agree that raising the income tax rate would lead to some reduction in the number of hours worked, they may yet differ in their views on the advisability of such a rise. One economist may believe that the additional revenue that may come in to government coffers is not worth the disincentives to work; another may think that, if such monies can be redistributed to benefit the needy, or provide valuable infrastructure, the negative impact on the workers paying the income tax is worth it.

Positive economics studies objective or scientific explanations of how the economy functions.

Normative economics offers recommendations that incorporate value judgments.

Scientific research can frequently resolve differences that arise in positive economics—not so in normative economics. For example, if we claim that "the elderly have high medical bills, and the government should cover all of the bills", we are making both a positive and a normative statement. The first part is positive, and its truth is easily established. The latter part is normative, and individuals of different beliefs may reasonably differ. Some people may believe that the money would be better spent on the environment and have the aged cover at least part of their own medical costs. Positive economics does not attempt to show that one of these views is correct and the other false. The views are based on value judgments, and are motivated by a concern for equity. Equity is a vital guiding principle in the formation of policy and is frequently, though not always, seen as being in competition with the drive for economic growth. Equity is driven primarily by normative considerations. Few economists would disagree with the assertion that a government should implement policies that improve the lot of the poor—but to what degree?

Economic equity is concerned with the distribution of well-being among members of the economy.

Application Box 2.1 Wealth Tax

US Senator Elizabeth Warren, in seeking her (Democratic) Party's nomination as candidate for the Presidency in 2020, proposed that individuals with high wealth should pay a wealth tax. Her proposal was to levy a tax of 2% on individual wealth holdings above \$50 million and a 6% tax on wealth above one billion dollars. This is clearly a *normative* approach to the issue of wealth concentration; it represented her ethical solution to what she perceived as socially unjust inequality.

In contrast, others in her party (Professor Larry Summers of Harvard for example) argued that the impact of such a tax would be to incentivize wealthy individuals to reclassify their wealth, or to give it to family members, or to offshore it, in order to avoid such a tax. If individuals behaved in this way the tax take would be far less than envisaged by Senator Warren. Such an analysis by Professor Summers is *positive* in nature; it attempts to define what might happen in response to the normative policy of Senator Warren. If he took the further step of saying that wealth should not be taxed, then he would be venturing into normative territory. Henry Aaron of the Brookings Institution argued that a more progressive inheritance tax than currently exists would be easier to implement, and would be more effective in both generating tax revenue and equalizing wealth holdings. Since he also advocated implementing such a proposal, as an alternative to Senator Warren's proposals, he was being both *normative and positive*.

Most economists hold normative views, sometimes very strongly. They frequently see themselves, not just as cold hearted scientists, but as champions for their (normative) cause in addition. Conservative economists see a smaller role for government than left-leaning economists.

Many economists see a conflict between equity and the efficiency considerations that we developed in Chapter 1. For example, high taxes may provide disincentives to work in the marketplace and therefore reduce the efficiency of the economy: Plumbers and gardeners may decide to do their own gardening and their own plumbing because, by staying out of the marketplace where monetary transactions are taxed, they can avoid the taxes. And avoiding the taxes may turn out to be as valuable as the efficiency gains they forgo.

In other areas the equity-efficiency trade-off is not so obvious: If taxes (that may have disincentive effects) are used to educate individuals who otherwise would not develop the skills that follow education, then economic growth may be higher as a result of the intervention.

Revisiting the definition of economics – core beliefs

This is an appropriate point at which to return to the definition of economics in Chapter 1 that we borrowed from Nobel Laureate Christopher Sims: Economics is a set of ideas and methods for the betterment of society.

If economics is concerned about the betterment of society, clearly there are ethical as well as efficiency considerations at play. And given the philosophical differences among scientists (including economists), can we define an approach to economics that is shared by the economics profession at large? Most economists would answer that the profession shares a set of beliefs, and that differences refer to the extent to which one consideration may collide with another.

- First of all we believe that *markets are critical* because they facilitate exchange and therefore encourage efficiency. Specialization and trade creates benefits for the trading parties. For example, Canada has not the appropriate climate for growing coffee beans, and Colombia has not the terrain for wheat. If Canada had to be self-sufficient, we might have to grow coffee beans in green-houses—a costly proposition. But with trade we can specialize, and then exchange some of our wheat for Colombian coffee. Similar benefits arise for the Colombians.

A frequent complaint against trade is that its modern-day form (globalization) does not benefit the poor. For example, workers in the Philippines may earn only a few dollars per day manufacturing clothing for Western markets. From this perspective, most of the gains from trade go to the Western consumers and capitalists, come at the expense of jobs to western workers, and provide Asian workers with meagre rewards.

- A corollary of the centrality of markets is *that incentives matter*. If the price of business class seats on your favourite airline is reduced, you may consider upgrading. Economists believe that the price mechanism influences behaviour, and therefore favour the use of price incentives in the marketplace and public policy more generally. Environmental economists, for example, advocate the use of pollution permits that can be traded at a price between users, or carbon taxes on the emission of greenhouse gases. We will develop such ideas in *Principles of Microeconomics* Chapter 5 more fully.
- In saying that economists believe in incentives, we are not proposing that human beings are purely mercenary. People have many motivations: Self-interest, a sense of public duty, kindness, etc. Acting out of a sense of self-interest does not imply that people are morally empty or have no altruistic sense.
- Economists believe universally in the *importance of the rule of law*, no matter where they sit on the political spectrum. Legal institutions that govern contracts are critical to the functioning of an economy. If goods and services are to be supplied in a market economy, the suppliers must be guaranteed that they will be remunerated. And this requires a developed legal structure with penalties imposed on individuals or groups who violate contracts. Markets alone will not function efficiently.

Modern development economics sees the implementation of the rule of law as perhaps the central challenge facing poorer economies. There is a strong correlation between economic growth and national wealth on the one hand, and an effective judicial and policing system on the other. The consequence on the world stage is that numerous 'economic' development projects now focus upon training jurists, police officers and bureaucrats in the rule of law!

- Finally, economists believe in the centrality of government. Governments can solve a number of problems that arise in market economies that cannot be addressed by the private market place. For example, governments can best address the potential abuses of monopoly power. Monopoly power, as we shall see in *Microeconomics* Chapter 10, not only has equity impacts it may also reduce economic efficiency. Governments are also best positioned to deal with environmental or other types of externalities – the impact of economic activity on sectors of the economy that are not directly involved in the activity under consideration.

In summary, governments have a variety of roles to play in the economy. These roles involve making the economy more equitable and more efficient by using their many powers.

Key Terms

Variables: measures that can take on different sizes.

Data: recorded values of variables.

Time series data: a set of measurements made sequentially at different points in time.

High (low) frequency data series have short (long) intervals between observations.

Cross-section data: values for different variables recorded at a point in time.

Repeated cross-section data: cross-section data recorded at regular or irregular intervals.

Longitudinal data follow the same units of observation through time.

Percentage change= (change in values)/original value \times 100.

Consumer price index: the average price level for consumer goods and services.

Inflation (deflation) rate: the annual percentage increase (decrease) in the level of consumer prices.

Real price: the actual price adjusted by the general (consumer) price level in the economy.

Index number: value for a variable, or an average of a set of variables, expressed relative to a given base value.

Regression line: representation of the average relationship between two variables in a scatter diagram.

Positive economics studies objective or scientific explanations of how the economy functions.

Normative economics offers recommendations that incorporate value judgments.

Economic equity is concerned with the distribution of well-being among members of the economy.

Exercises for Chapter 2

EXERCISE 2.1

An examination of a country's recent international trade flows yields the data in the table below.

Year	National Income (\$b)	Imports (\$b)
2011	1,500	550
2012	1,575	573
2013	1,701	610
2014	1,531	560
2015	1,638	591

1. Based on an examination of these data do you think the national income and imports are not related, positively related, or negatively related?
2. Plot each pair of observations in a two-dimensional line diagram to illustrate your view of the import/income relationship. Measure income on the horizontal axis and imports on the vertical axis. This can be done using graph paper or a spreadsheet-cum-graphics software.

EXERCISE 2.2

The average price of a medium coffee at *Wakeup Coffee Shop* in each of the past ten years is given in the table below.

2005	2006	2007	2008	2009	2010	2011	2012	2013	2014
\$1.05	\$1.10	\$1.14	\$1.20	\$1.25	\$1.25	\$1.33	\$1.35	\$1.45	\$1.49

1. Construct an annual 'coffee price index' for this time period using 2005 as the base year. [Hint: follow the procedure detailed in the chapter – divide each yearly price by the base year price.]
2. Based on your price index, what was the percentage change in the price of a medium coffee from 2005 to 2012?
3. Based on your index, what was the average annual percentage change in the price of coffee from 2005 to 2010?
4. Assuming the inflation rate in this economy was 2% every year, what was the real change in the price of coffee between 2007 and 2008; and between 2009 and 2010?

EXERCISE 2.3

The following table shows hypothetical consumption spending by households and income of households in billions of dollars.

Year	Income	Consumption
2006	476	434
2007	482	447
2008	495	454
2009	505	471
2010	525	489
2011	539	509
2012	550	530
2013	567	548

1. Plot the scatter diagram with consumption on the vertical axis and income on the horizontal axis.
2. Fit a line through these points.
3. Does the line indicate that these two variables are related to each other?
4. How would you describe the *causal relationship* between income and consumption?

EXERCISE 2.4

Using the data from Exercise 2.3, compute the percentage change in consumption and the percentage change in income for each pair of adjoining years between 2006 and 2013.

EXERCISE 2.5

You are told that the relationship between two variables, X and Y , has the form $Y=10+2X$. By trying different values for X you can obtain the corresponding predicted value for Y (e.g., if $X=3$, then $Y = 10 + 2 \times 3 = 16$). For values of X between 0 and 12, compute the matching value of Y and plot the scatter diagram.

EXERCISE 2.6

For the data below, plot a scatter diagram with variable Y on the vertical axis and variable X on the horizontal axis.

Y	40	33	29	56	81	19	20
X	5	7	9	3	1	11	10

1. Is the relationship between the variables positive or negative?
 2. Do you think that a linear or non-linear line better describes the relationship?
1. This task is the job of econometricians, who practice econometrics. Econometrics is the science of examining and quantifying relationships between economic variables. It attempts to determine the separate influences of each variable, in an environment where many things move simultaneously. Computer algorithms that do this are plentiful. Computers can also work in many dimensions in order to capture the influences of *several variables simultaneously* if the model requires that.

This page titled [2: Theories, data and beliefs](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

2.1: Data analysis

The analysis of behaviour necessarily involves data. Data may serve to validate or contradict a theory. Data analysis, even without being motivated by economic theory, frequently displays patterns of behaviour that merit examination. The terms *variables* and *data* are related. Variables are measures that can take on different magnitudes. The interest rate on a student loan, for example, is a variable with a certain value at a point in time but perhaps a different value at an earlier or later date. Economic theories and models explain the causal relationships between variables. In contrast, Data are the recorded values of variables. Sets of data provide specific values for the variables we want to study and analyze. Knowing that gross domestic product (a variable) declined in 2009 is just a partial description of events. If the data indicate that it decreased by exactly 3%, we know a great deal more – the decline was large.

Variables: measures that can take on different values.

Data: recorded values of variables.

Sets of data help us to test our models or theories, but first we need to pay attention to the economic logic involved in observations and modelling. For example, if sunspots or baggy pants were found to be correlated with economic expansion, would we consider these events a coincidence or a key to understanding economic growth? The observation is based on facts or data, but it need not have any economic content. The economist's task is to distinguish between coincidence and economic causation. Merely because variables are associated or correlated does not mean that one causes the other.

While the more frequent wearing of loose clothing in the past may have been associated with economic growth because they both occurred at the same time (correlation), one could not argue on a logical basis that this behaviour causes good economic times. Therefore, the past association of these variables should be considered as no more than a coincidence. Once specified on the basis of economic logic, a model must be tested to determine its usefulness in explaining observed economic events.

Table 2.1 House prices and price indexes

Year	House prices in dollars (P_H)	Percentage change in P_H	Percentage change in consumer prices	Real percentage change in P_H	Index for price of housing	5-year mortgage rate
2001	350,000				100	7.75
2002	360,000				102.9	6.85
2003	395,000	35,000/360,000=9.7 %	3%	6.7%	112.9	6.6
2004	434,000				124.0	5.8
2005	477,000				136.3	6.1
2006	580,000				165.7	6.3
2007	630,000				180.0	6.65
2008	710,000				202.9	7.3
2009	605,000	-105,000/710,000=-14.8%	1.6%	-16.4%	172.9	5.8
2010	740,000				211.4	5.4
2011	800,000				228.6	5.2

Note: Data on changes in consumer prices come from Statistics Canada, CANSIM series V41692930; data on house prices are for N. Vancouver from *Royal Le Page*; data on mortgage rates from <http://www.ratehub.ca>. Index for house prices obtained by scaling each entry in column 2 by 100/350,000. The real percentage change in the price of housing is: The percentage change in the price of housing minus the percentage change in consumer prices.

Data types

Data come in several forms. One form is time-series, which reflects a set of measurements made in sequence at different points in time. The first column in Table 2.1 reports the values for house prices in North Vancouver for the first quarter of each year, between 2001 and 2011. Evidently this is a time series. Annual data report one observation per year. We could, alternatively, have presented the data in monthly, weekly, or even daily form. The frequency we use depends on the purpose: If we are interested in the longer-term trend in house prices, then the annual form suffices. In contrast, financial economists, who study the behaviour of stock prices, might not be content with daily or even hourly prices; they may need prices minute-by-minute. Such data are called high-frequency data, whereas annual data are low-frequency data.

Table 2.2 Unemployment rates, Canada and Provinces, monthly 2012, seasonally adjusted

	Jan	Feb	Mar	Apr	May	Jun
CANADA	7.6	7.4	7.2	7.3	7.3	7.2
NFLD	13.5	12.9	13.0	12.3	12.0	13.0
PEI	12.2	10.5	11.3	11.0	11.3	11.3
NS	8.4	8.2	8.3	9.0	9.2	9.6
NB	9.5	10.1	12.2	9.8	9.4	9.5
QUE	8.4	8.4	7.9	8.0	7.8	7.7
ONT	8.1	7.6	7.4	7.8	7.8	7.8
MAN	5.4	5.6	5.3	5.3	5.1	5.2
SASK	5.0	5.0	4.8	4.9	4.5	4.9
ALTA	4.9	5.0	5.3	4.9	4.5	4.6
BC	6.9	6.9	7.0	6.2	7.4	6.6

Source: Statistics Canada CANSIM Table 282-0087.

Time-series: a set of measurements made sequentially at different points in time.

High (low) frequency data: series with short (long) intervals between observations.

In contrast to time-series data, cross-section data record the values of different variables at a point in time. Table 2.2 contains a cross-section of unemployment rates for Canada and Canadian provinces economies. For January 2012 we have a snapshot of the provincial economies at that point in time, likewise for the months until June. This table therefore contains repeated cross-sections.

When the unit of observation is the same over time such repeated cross sections are called longitudinal data. For example, a health survey that followed and interviewed the same individuals over time would yield longitudinal data. If the individuals differ each time the survey is conducted, the data are repeated cross sections. Longitudinal data therefore follow the same units of observation through time.

Cross-section data: values for different variables recorded at a point in time.

Repeated cross-section data: cross-section data recorded at regular or irregular intervals.

Longitudinal data: follow the same units of observation through time.

Graphing the data

Data can be presented in graphical as well as tabular form. Figure 2.1 plots the house price data from the second column of Table 2.1. Each asterisk in the figure represents a price value and a corresponding time period. The horizontal axis reflects time, the vertical axis price in dollars. The graphical presentation of data simply provides a visual rather than numeric perspective. It is immediately evident that house prices increased consistently during this 11-year period, with a single downward 'correction' in 2009. We have plotted the data a second time in Figure 2.2 to illustrate the need to read graphs carefully. The greater *apparent* slope in Figure 2.1 might easily be interpreted to mean that prices increased more steeply than suggested in Figure 2.2. But a

careful reading of the axes reveals that this is not so; using different scales when plotting data or constructing diagrams can mislead the unaware viewer.

Figure 2.1 House prices in dollars 1999-2012

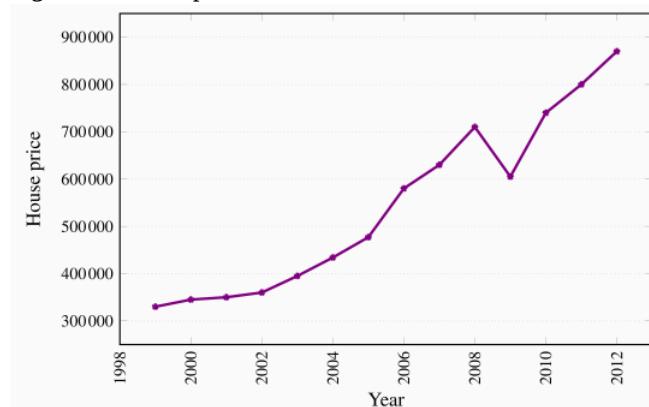
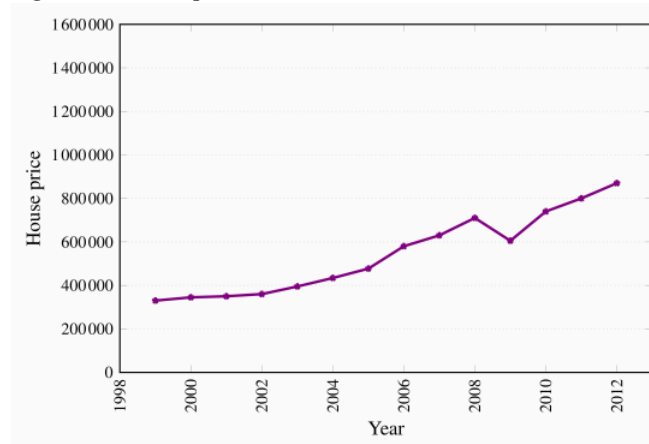


Figure 2.2 House prices in dollars 1999-2012



Percentage changes

The use of percentages makes the analysis of data particularly simple. Suppose we wanted to compare the prices of New York luxury condominiums with the prices of homes in rural Mississippi. In the latter case, a change in average prices of \$10,000 might be considered enormous, whereas a change of one million dollars in New York might be pretty normal – because the average price in New York is so much higher than in Mississippi. To make comparisons between the two markets, we can use the concept of a percentage change. This is defined as the change in the value of the variable, relative to its initial value, multiplied by 100.

Percentage change = $\frac{(\text{change in values})}{(\text{original value})} \times 100$.

The third column of Table 2.1 contains the values of the percentage change in house prices for two pairs of years. Between 2002 and 2003 the price change was \$35,000. Relative to the price in the first of these two years this change was the fraction $35,000/395,000=0.097$. If we multiply this fraction by 100 we obtain a percentage price change of 9.7%. Evidently we could calculate the percentage price changes for all pairs of years. A second price change is calculated for the 2008-2009 pair of years. Here price declined and the result is thus a negative percentage change.

Consumer prices

Most variables in economics are averages of the components that go into them. When variables are denominated in dollar terms it is important to be able to interpret them correctly. While the house price series above indicates a strong pattern of price increases, it is vital to know if the price of housing increased more or less rapidly than other prices in the economy. If all prices in the economy were increasing in line with house prices there would be no special information in the house price series. However, if house prices increased more rapidly than prices in general, then the data indicate that something special took place in the housing market during

the decade in question. To determine an answer to this we need to know the degree to which the general price level changed each year.

Statistics Canada regularly surveys the price of virtually every product produced in the economy. One such survey records the prices of goods and services purchased by consumers. Statistics Canada then computes an average price level for all of these goods combined for each time period the survey is carried out (monthly). Once Statistics Canada has computed the average consumer price, it can compute the change in the price level from one period to the next. In Table 2.1 two such values are entered in the following data column: Consumer prices increased by 3% between 2002 and 2003, and by 1.6% between 2008 and 2009. These percentage changes in the general price level represent inflation if prices increase, and deflation if prices decline.

In this market it is clear that housing price changes were substantially larger than the changes in consumer prices for these two pairs of years. The next column provides information on the difference between the house price changes and changes in the general consumer price level, in percentage terms. This is (approximately) the change in the relative price of housing, or what economists call the real price of housing. It is obtained by subtracting the rate of change in the general price index from the rate of change in the variable of interest.

Consumer price index: the average price level for consumer goods and services.

Inflation (deflation) rate: the annual percentage increase (decrease) in the level of consumer prices.

Real price: the actual price adjusted by the general (consumer) price level in the economy.

Index numbers

Statistics Canada and other statistical agencies frequently present data in index number form. An index number provides an easy way to read the data. For example, suppose we wanted to compute the percentage change in the price of housing between 2001 and 2007. We could do this by entering the two data points in a spreadsheet or calculator and do the computation. But suppose the prices were entered in another form. In particular, by dividing each price value by the first year value and multiplying the result by 100 we obtain a series of prices that are all relative to the initial year – which we call the base year. The resulting series in column 6 of Table 2.1 is an index of house price values. Each entry is the corresponding value in column 2, divided by the first entry in column 2. For example, the value 124.0 in row 4 is obtained as $(434,000/350,000) * 100 = 124.0$. The key characteristics of indexes are that they are *not dependent upon the units of measurement of the data in question*, and they are interpretable easily with reference to a given base value. To illustrate, suppose we wish to know how prices behaved between 2001 and 2007. The index number column immediately tells us that prices increased by 80%, because relative to 2001, the 2007 value is 80% higher.

Index number: value for a variable, or an average of a set of variables, expressed relative to a given base value.

Furthermore, index numbers enable us to make *comparisons with the price patterns for other goods* much more easily. If we had constructed a price index for automobiles, which also had a base value of 100 in 2001, we could make immediate comparisons without having to compare one set of numbers defined in thousands of dollars with another defined in hundreds of thousands of dollars. In short, index numbers simplify the interpretation of data.

This page titled [2.1: Data analysis](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

2.2: Data, theory and economic models

Let us now investigate the interplay between economic theories on the one hand and data on the other. We will develop two examples. The first will be based upon the data on house prices, the second upon a new data set.

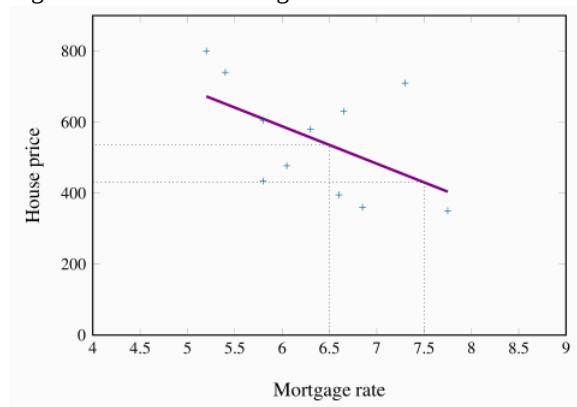
House prices – theory

Remember from Chapter 1 that a theory is a logical argument regarding economic relationships. A theory of house prices would propose that the price of housing depends upon a number of elements in the economy. In particular, if borrowing costs are low then buyers are able to afford the interest costs on larger borrowings. This in turn might mean they are willing to pay higher prices. Conversely, if borrowing rates are higher. Consequently, the borrowing rate, or mortgage rate, is a variable for an economic model of house prices. A second variable might be available space for development: If space in a given metropolitan area is tight then the land value will reflect this, and consequently the higher land price should be reflected in higher house prices. A third variable would be the business climate: If there is a high volume of high-value business transacted in a given area then buildings will be more in demand, and that in turn should be reflected in higher prices. For example, both business and residential properties are more highly priced in San Francisco and New York than in Moncton, New Brunswick. A fourth variable might be environmental attractiveness: Vancouver may be more enticing than other towns in Canada. A fifth variable might be the climate.

House prices – evidence

These and other variables could form the basis of a *theory* of house prices. A *model* of house prices, as explained in Chapter 1, focuses upon what we would consider to be the most important subset of these variables. In the limit, we could have an extremely simple model that specified a dependence between the price of housing and the mortgage rate alone. To test such a simple model we need data on house prices and mortgage rates. The final column of Table 2.1 contains data on the 5-year fixed-rate mortgage for the period in question. Since our simple model proposes that prices depend (primarily) upon mortgage rates, in Figure 2.3 we plot the house price series on the vertical axis, and the mortgage rate on the horizontal axis, for each year from 2001 to 2011. As before, each point (shown as a '+') represents a pair of price and mortgage rate values.

Figure 2.3 Price of housing



The resulting plot (called a scatter diagram) suggests that there is a negative relationship between these two variables. That is, higher prices are correlated with lower mortgage rates. Such a correlation is consistent with our theory of house prices, and so we might conclude that changes in mortgage rates *cause* changes in house prices. Or at least the data suggest that we should not reject the idea that such causation is in the data.

House prices – inference

To summarize the relationship between these variables, the pattern suggests that a straight line through the scatter plot would provide a reasonably good description of the relationship between these variables. Obviously it is important to define the most appropriate line – one that 'fits' the data well.¹ The line we have drawn through the data points is informative, because it relates the two variables in a *quantitative manner*. It is called a regression line. It predicts that, on average, if the mortgage rate increases, the price of housing will respond in the downward direction. This particular line states that a one point change in the mortgage rate will move prices in the opposing direction by \$105,000. This is easily verified by considering the dollar value corresponding to say a

mortgage value of 6.5, and then the value corresponding to a mortgage value of 7.5. Projecting vertically to the regression line from each of these points on the horizontal axis, and from there across to the vertical axis will produce a change in price of \$105,000.

Note that the line is not at all a 'perfect' fit. For example, the mortgage rate declined between 2008 and 2009, but the price declined also – contrary to our theory. The model is not a perfect predictor; it states that *on average* a change in the magnitude of the x-axis variable leads to a change of a specific amount in the magnitude of the y-axis variable.

In this instance the slope of the line is given by $-105,000/1$, which is the vertical distance divided by the corresponding horizontal distance. Since the line is straight, this slope is unchanging.

Regression line: representation of the average relationship between two variables in a scatter diagram.

Road fatalities – theory, evidence and inference

Table 2.3 contains data on annual road fatalities per 100,000 drivers for various age groups. In the background, we have a *theory*, proposing that driver fatalities depend upon the age of the driver, the quality of roads and signage, speed limits, the age of the automobile stock and perhaps some other variables. Our model focuses upon a subset of these variables, and in order to present the example in graphical terms we specify fatalities as being dependent upon a single variable – age of driver.

Table 2.3 Non-linearity: Driver fatality rates Canada, 2009

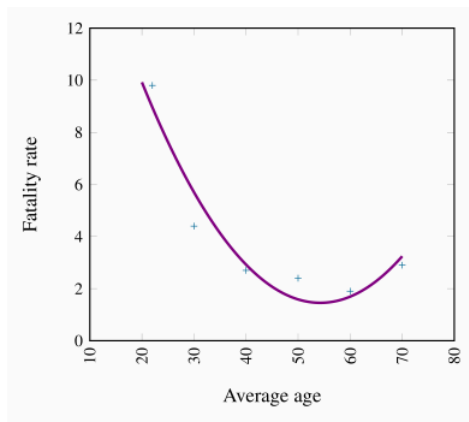
Age of driver	Fatality rate per 100,000 drivers
20-24	9.8
25-34	4.4
35-44	2.7
45-54	2.4
55-64	1.9
65+	2.9

Source: Transport Canada, Canadian motor vehicle traffic collision statistics, 2009.

The scatter diagram is presented in Figure 2.4. Two aspects of this plot stand out. First, there is an exceedingly steep decline in the fatality rate when we go from the youngest age group to the next two age groups. The decline in fatalities between the youngest and second youngest groups is 5.4 points, and between the second and third age groups is 1.7 points. The decline between the third and fourth groups is minimal - just 0.3 points. Hence, behaviour is not the same throughout the age distribution. Second, we notice that fatalities increase for the oldest age group, perhaps indicating that the oldest drivers are not as good as middle-aged drivers.

These two features suggest that the relationship between fatalities and age differs across the age spectrum. Accordingly, a straight line would not be an accurate way of representing the behaviours in these data. A straight line through the plot implies that a given change in age should have a similar impact on fatalities, no matter the age group. Accordingly we have an example of a *non-linear relationship*. Such a non-linear relationship might be represented by the curve going through the plot. Clearly the slope of this line varies as we move from one age category to another.

Figure 2.4 Non-linearity: Driver fatality rates Canada, 2009



Fatality rates vary non-linearly with age: At first they decline, then increase again, relative to the youngest age group.

This page titled [2.2: Data, theory and economic models](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

2.3: Ethics, efficiency and beliefs

Positive economics studies objective or scientific explanations of how the economy functions. Its aim is to understand and generate predictions about how the economy may respond to changes and policy initiatives. In this effort economists strive to act as detached scientists, regardless of political sympathies or ethical code. Personal judgments and preferences are (ideally) kept apart. In this particular sense, economics is similar to the natural sciences such as physics or biology. To date in this chapter we have been exploring economics primarily from a positive standpoint.

In contrast, normative economics offers recommendations based partly on value judgments. While economists of different political persuasions can agree that raising the income tax rate would lead to some reduction in the number of hours worked, they may yet differ in their views on the advisability of such a rise. One economist may believe that the additional revenue that may come in to government coffers is not worth the disincentives to work; another may think that, if such monies can be redistributed to benefit the needy, or provide valuable infrastructure, the negative impact on the workers paying the income tax is worth it.

Positive economics studies objective or scientific explanations of how the economy functions.

Normative economics offers recommendations that incorporate value judgments.

Scientific research can frequently resolve differences that arise in positive economics—not so in normative economics. For example, if we claim that "the elderly have high medical bills, and the government should cover all of the bills", we are making both a positive and a normative statement. The first part is positive, and its truth is easily established. The latter part is normative, and individuals of different beliefs may reasonably differ. Some people may believe that the money would be better spent on the environment and have the aged cover at least part of their own medical costs. Positive economics does not attempt to show that one of these views is correct and the other false. The views are based on value judgments, and are motivated by a concern for equity. Equity is a vital guiding principle in the formation of policy and is frequently, though not always, seen as being in competition with the drive for economic growth. Equity is driven primarily by normative considerations. Few economists would disagree with the assertion that a government should implement policies that improve the lot of the poor—but to what degree?

Economic equity is concerned with the distribution of well-being among members of the economy.

Application Box 2.1 Wealth Tax

US Senator Elizabeth Warren, in seeking her (Democratic) Party's nomination as candidate for the Presidency in 2020, proposed that individuals with high wealth should pay a wealth tax. Her proposal was to levy a tax of 2% on individual wealth holdings above \$50 million and a 6% tax on wealth above one billion dollars. This is clearly a *normative* approach to the issue of wealth concentration; it represented her ethical solution to what she perceived as socially unjust inequality.

In contrast, others in her party (Professor Larry Summers of Harvard for example) argued that the impact of such a tax would be to incentivize wealthy individuals to reclassify their wealth, or to give it to family members, or to offshore it, in order to avoid such a tax. If individuals behaved in this way the tax take would be far less than envisaged by Senator Warren. Such an analysis by Professor Summers is *positive* in nature; it attempts to define what might happen in response to the normative policy of Senator Warren. If he took the further step of saying that wealth should not be taxed, then he would be venturing into normative territory. Henry Aaron of the Brookings Institution argued that a more progressive inheritance tax than currently exists would be easier to implement, and would be more effective in both generating tax revenue and equalizing wealth holdings. Since he also advocated implementing such a proposal, as an alternative to Senator Warren's proposals, he was being both *normative and positive*.

Most economists hold normative views, sometimes very strongly. They frequently see themselves, not just as cold hearted scientists, but as champions for their (normative) cause in addition. Conservative economists see a smaller role for government than left-leaning economists.

Many economists see a conflict between equity and the efficiency considerations that we developed in Chapter 1. For example, high taxes may provide disincentives to work in the marketplace and therefore reduce the efficiency of the economy: Plumbers and gardeners may decide to do their own gardening and their own plumbing because, by staying out of the marketplace where monetary transactions are taxed, they can avoid the taxes. And avoiding the taxes may turn out to be as valuable as the efficiency gains they forgo.

In other areas the equity-efficiency trade-off is not so obvious: If taxes (that may have disincentive effects) are used to educate individuals who otherwise would not develop the skills that follow education, then economic growth may be higher as a result of the intervention.

Revisiting the definition of economics – core beliefs

This is an appropriate point at which to return to the definition of economics in Chapter 1 that we borrowed from Nobel Laureate Christopher Sims: Economics is a set of ideas and methods for the betterment of society.

If economics is concerned about the betterment of society, clearly there are ethical as well as efficiency considerations at play. And given the philosophical differences among scientists (including economists), can we define an approach to economics that is shared by the economics profession at large? Most economists would answer that the profession shares a set of beliefs, and that differences refer to the extent to which one consideration may collide with another.

- First of all we believe that *markets are critical* because they facilitate exchange and therefore encourage efficiency. Specialization and trade creates benefits for the trading parties. For example, Canada has not the appropriate climate for growing coffee beans, and Colombia has not the terrain for wheat. If Canada had to be self-sufficient, we might have to grow coffee beans in green-houses—a costly proposition. But with trade we can specialize, and then exchange some of our wheat for Colombian coffee. Similar benefits arise for the Colombians.

A frequent complaint against trade is that its modern-day form (globalization) does not benefit the poor. For example, workers in the Philippines may earn only a few dollars per day manufacturing clothing for Western markets. From this perspective, most of the gains from trade go to the Western consumers and capitalists, come at the expense of jobs to western workers, and provide Asian workers with meagre rewards.

- A corollary of the centrality of markets is *that incentives matter*. If the price of business class seats on your favourite airline is reduced, you may consider upgrading. Economists believe that the price mechanism influences behaviour, and therefore favour the use of price incentives in the marketplace and public policy more generally. Environmental economists, for example, advocate the use of pollution permits that can be traded at a price between users, or carbon taxes on the emission of greenhouse gases. We will develop such ideas in *Principles of Microeconomics* Chapter 5 more fully.
- In saying that economists believe in incentives, we are not proposing that human beings are purely mercenary. People have many motivations: Self-interest, a sense of public duty, kindness, etc. Acting out of a sense of self-interest does not imply that people are morally empty or have no altruistic sense.
- Economists believe universally in the *importance of the rule of law*, no matter where they sit on the political spectrum. Legal institutions that govern contracts are critical to the functioning of an economy. If goods and services are to be supplied in a market economy, the suppliers must be guaranteed that they will be remunerated. And this requires a developed legal structure with penalties imposed on individuals or groups who violate contracts. Markets alone will not function efficiently.

Modern development economics sees the implementation of the rule of law as perhaps the central challenge facing poorer economies. There is a strong correlation between economic growth and national wealth on the one hand, and an effective judicial and policing system on the other. The consequence on the world stage is that numerous 'economic' development projects now focus upon training jurists, police officers and bureaucrats in the rule of law!

- Finally, economists believe in the centrality of government. Governments can solve a number of problems that arise in market economies that cannot be addressed by the private market place. For example, governments can best address the potential abuses of monopoly power. Monopoly power, as we shall see in *Microeconomics* Chapter 10, not only has equity impacts it may also reduce economic efficiency. Governments are also best positioned to deal with environmental or other types of externalities – the impact of economic activity on sectors of the economy that are not directly involved in the activity under consideration.

In summary, governments have a variety of roles to play in the economy. These roles involve making the economy more equitable and more efficient by using their many powers.

This page titled [2.3: Ethics, efficiency and beliefs](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

2.4: Key Terms

Variables: measures that can take on different sizes.

Data: recorded values of variables.

Time series data: a set of measurements made sequentially at different points in time.

High (low) frequency data series have short (long) intervals between observations.

Cross-section data: values for different variables recorded at a point in time.

Repeated cross-section data: cross-section data recorded at regular or irregular intervals.

Longitudinal data follow the same units of observation through time.

Percentage change= (change in values)/original value \times 100.

Consumer price index: the average price level for consumer goods and services.

Inflation (deflation) rate: the annual percentage increase (decrease) in the level of consumer prices.

Real price: the actual price adjusted by the general (consumer) price level in the economy.

Index number: value for a variable, or an average of a set of variables, expressed relative to a given base value.

Regression line: representation of the average relationship between two variables in a scatter diagram.

Positive economics studies objective or scientific explanations of how the economy functions.

Normative economics offers recommendations that incorporate value judgments.

Economic equity is concerned with the distribution of well-being among members of the economy.

This page titled [2.4: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

2.5: Exercises for Chapter 2

EXERCISE 2.1

An examination of a country's recent international trade flows yields the data in the table below.

Year	National Income (\$b)	Imports (\$b)
2011	1,500	550
2012	1,575	573
2013	1,701	610
2014	1,531	560
2015	1,638	591

- Based on an examination of these data do you think the national income and imports are not related, positively related, or negatively related?
- Plot each pair of observations in a two-dimensional line diagram to illustrate your view of the import/income relationship. Measure income on the horizontal axis and imports on the vertical axis. This can be done using graph paper or a spreadsheet-cum-graphics software.

EXERCISE 2.2

The average price of a medium coffee at *Wakeup Coffee Shop* in each of the past ten years is given in the table below.

2005	2006	2007	2008	2009	2010	2011	2012	2013	2014
\$1.05	\$1.10	\$1.14	\$1.20	\$1.25	\$1.25	\$1.33	\$1.35	\$1.45	\$1.49

- Construct an annual 'coffee price index' for this time period using 2005 as the base year. [*Hint*: follow the procedure detailed in the chapter – divide each yearly price by the base year price.]
- Based on your price index, what was the percentage change in the price of a medium coffee from 2005 to 2012?
- Based on your index, what was the average annual percentage change in the price of coffee from 2005 to 2010?
- Assuming the inflation rate in this economy was 2% every year, what was the real change in the price of coffee between 2007 and 2008; and between 2009 and 2010?

EXERCISE 2.3

The following table shows hypothetical consumption spending by households and income of households in billions of dollars.

Year	Income	Consumption
2006	476	434
2007	482	447
2008	495	454
2009	505	471
2010	525	489
2011	539	509
2012	550	530
2013	567	548

- Plot the scatter diagram with consumption on the vertical axis and income on the horizontal axis.

- b. Fit a line through these points.
- c. Does the line indicate that these two variables are related to each other?
- d. How would you describe the *causal relationship* between income and consumption?

EXERCISE 2.4

Using the data from Exercise 2.3, compute the percentage change in consumption and the percentage change in income for each pair of adjoining years between 2006 and 2013.

EXERCISE 2.5

You are told that the relationship between two variables, X and Y , has the form $Y=10+2X$. By trying different values for X you can obtain the corresponding predicted value for Y (e.g., if $X=3$, then $Y = 10 + 2 \times 3 = 16$). For values of X between 0 and 12, compute the matching value of Y and plot the scatter diagram.

EXERCISE 2.6

For the data below, plot a scatter diagram with variable Y on the vertical axis and variable X on the horizontal axis.

Y	40	33	29	56	81	19	20
X	5	7	9	3	1	11	10

- a. Is the relationship between the variables positive or negative?
 - b. Do you think that a linear or non-linear line better describes the relationship?
1. This task is the job of econometricians, who practice econometrics. Econometrics is the science of examining and quantifying relationships between economic variables. It attempts to determine the separate influences of each variable, in an environment where many things move simultaneously. Computer algorithms that do this are plentiful. Computers can also work in many dimensions in order to capture the influences of *several variables simultaneously* if the model requires that.

This page titled [2.5: Exercises for Chapter 2](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3: The classical marketplace – demand and supply

Chapter 3: The classical marketplace – demand and supply

In this chapter we will explore:

3.1	The marketplace – trading
3.2	The market's building blocks
3.3	Demand curves and supply curves
3.4	Non-price determinants of demand
3.5	Non-price determinants of supply
3.6	Simultaneous demand and supply movements
3.7	Free and managed markets – interventions
3.8	From individuals to markets
3.9	Useful techniques – demand and supply equations

3.1 The marketplace – trading

The marketplace in today's economy has evolved from earlier times. It no longer has a unique form – one where buyers and sellers physically come together for the purpose of exchange. Indeed, supermarkets usually require individuals to be physically present to make their purchases. But when purchasing an airline ticket, individuals simply go online and interact with perhaps a number of different airlines (suppliers) simultaneously. Or again, individuals may simply give an instruction to their stock broker, who will execute a purchase on their behalf – the broker performs the role of a middleman, who may additionally give advice to the purchaser. Or a marketing agency may decide to subcontract work to a translator or graphic artist who resides in Mumbai. The advent of the coronavirus has shifted grocery purchases from in-store presence to home delivery for many buyers. In pure auctions (where a single work of art or a single residence is offered for sale) buyers compete one against the other for the single item supplied. Accommodations in private homes are supplied to potential visitors (buyers) through *Airbnb*. Taxi rides are mediated through *Lyft* or *Uber*. These institutions are all different types of markets; they serve the purpose of facilitating exchange and trade.

Not all goods and services in the modern economy are obtained through the marketplace. Schooling and health care are allocated in Canada primarily by government decree. In some instances the market plays a supporting role: Universities and colleges may levy fees, and most individuals must pay, at least in part, for their pharmaceuticals. In contrast, broadcasting services may carry a price of zero; BuzzFeed or other news and social media come free of payment. Furthermore, some markets have no price, yet they find a way of facilitating an exchange. For example, graduating medical students need to be matched with hospitals for their residencies. Matching mechanisms are a form of market in that they bring together suppliers and demanders. We explore their operation in Chapter 11.

The importance of the marketplace springs from its role as an allocating mechanism. Elevated prices effectively send a signal to suppliers that the buyers in the market place a high value on the product being traded; conversely when prices are low. Accordingly, suppliers may decide to cease supplying markets where prices do not remunerate them sufficiently, and redirect their energies and the productive resources under their control to other markets – markets where the product being traded is more highly valued, and where the buyer is willing to pay more.

Whatever their form, the marketplace is central to the economy we live in. Not only does it facilitate trade, it also provides a means of earning a livelihood. Suppliers must hire resources – human and non-human – in order to bring their supplies to market and these resources must be paid a return: income is generated.

In this chapter we will examine the process of price formation – how the prices that we observe in the marketplace come to be what they are. We will illustrate that the price for a good is inevitably linked to the quantity of a good; price and quantity are different sides of the same coin and cannot generally be analyzed separately. To understand this process more fully, we need to *model* a typical market. The essentials are demand and supply.

3.2 The market's building blocks

In economics we use the terminology that describes trade in a particular manner. Non-economists frequently describe microeconomics by saying "it's all about supply and demand". While this is largely true we need to define exactly what we mean by these two central words. Demand is the quantity of a good or service that buyers wish to purchase at each conceivable price, with all other influences on demand remaining unchanged. It reflects a multitude of values, not a single value. It is not a single or unique quantity such as two cell phones, but rather a full description of the quantity of a good or service that buyers would purchase at various prices.

Demand is the quantity of a good or service that buyers wish to purchase at each possible price, with all other influences on demand remaining unchanged.

As a hypothetical example, the first column of Table 3.1 shows the price of natural gas per cubic foot. The second column shows the quantity that would be purchased in a given time period at each price. It is therefore a schedule of quantities demanded at various prices. For example, at a price \$6 per unit, buyers would like to purchase 4 units, whereas at the lower price of \$3 buyers would like to purchase 7 units. Note also that this is a homogeneous good. A cubic foot of natural gas is considered to be the same product no matter which supplier brings it to the market. In contrast, accommodations supplied through *Airbnb* are heterogeneous – they vary in size and quality.

Table 3.1 Demand and supply for natural gas

Price (\$)	Demand (thousands of cu feet)	Supply (thousands of cu feet)	Excess
10	0	18	Excess Supply
9	1	16	
8	2	14	
7	3	12	
6	4	10	
5	5	8	Equilibrium
4	6	6	
3	7	4	Excess Demand
2	8	2	
1	9	0	
0	10	0	

Supply is interpreted in a similar manner. It is not a single value; we say that supply is the quantity of a good or service that sellers are willing to sell at each possible price, with all other influences on supply remaining unchanged. Such a supply schedule is defined in the third column of the table. It is assumed that no supplier can make a profit (on account of their costs) unless the price is at least \$2 per unit, and therefore a zero quantity is supplied below that price. The higher price is more profitable, and therefore induces a greater quantity supplied, perhaps by attracting more suppliers. This is reflected in the data. For example, at a price of \$3 suppliers are willing to supply 4 units, whereas with a price of \$7 they are willing to supply 12 units. There is thus a positive relationship between price and quantity for the supplier – a higher price induces a greater quantity; whereas on the demand side of the market a higher price induces a lower quantity demanded – a negative relationship.

Supply is the quantity of a good or service that sellers are willing to sell at each possible price, with all other influences on supply remaining unchanged.

We can now identify a key difference in terminology – between the words demand and quantity demanded, and between supply and quantity supplied. While the words demand and supply refer to the complete schedules of demand and supply, the terms quantity demanded and quantity supplied each define a single value of demand or supply at a particular price.

Quantity demanded defines the amount purchased at a particular price.

Quantity supplied refers to the amount supplied at a particular price.

Thus while the non-economist may say that when some fans did not get tickets to the Stanley Cup it was a case of demand exceeding supply, as economists we say that the quantity demanded exceeded the quantity supplied *at the going price of tickets*. In this instance, had every ticket been offered at a sufficiently high price, the market could have generated an excess supply rather than an excess demand. A higher ticket price would reduce the *quantity demanded*; yet would not change *demand*, because demand refers to the whole schedule of possible quantities demanded at different prices.

Other things equal – ceteris paribus

The demand and supply schedules rest on the assumption that all other influences on supply and demand remain the same as we move up and down the possible price values. The expression *other things being equal*, or its Latin counterpart *ceteris paribus*, describes this constancy of other influences. For example, we assume on the demand side that the prices of other goods remain constant, and that tastes and incomes are unchanging. On the supply side we assume, for example, that there is no technological change in production methods. If any of these elements change then the market supply or demand schedules will reflect such changes. For example, if coal or oil prices increase (decline) then some buyers may switch to (away from) gas or solar power. This will be reflected in the data: At any given price more (or less) will be demanded. We will illustrate this in graphic form presently.

Market equilibrium

Let us now bring the demand and supply schedules together in an attempt to analyze what the marketplace will produce – will a single price emerge that will equate supply and demand? We will keep other things constant for the moment, and explore what materializes at different prices. At low prices, the data in Table 3.1 indicate that the quantity demanded exceeds the quantity supplied – for example, verify what happens when the price is \$3 per unit. The opposite occurs when the price is high – what would happen if the price were \$8? Evidently, there exists an intermediate price, where the quantity demanded equals the quantity supplied. At this point we say that the market is in equilibrium. The equilibrium price equates demand and supply – it clears the market.

The **equilibrium price** equilibrates the market. It is the price at which quantity demanded equals the quantity supplied.

In Table 3.1 the equilibrium price is \$4, and the equilibrium quantity is 6 thousand cubic feet of gas (we use the notation 'k' to denote thousands). At higher prices there is an excess supply—suppliers wish to sell more than buyers wish to buy. Conversely, at lower prices there is an excess demand. Only at the equilibrium price is the quantity supplied equal to the quantity demanded.

Excess supply exists when the quantity supplied exceeds the quantity demanded at the going price.

Excess demand exists when the quantity demanded exceeds the quantity supplied at the going price.

Does the market automatically reach equilibrium? To answer this question, suppose initially that the sellers choose a price of \$10. Here suppliers would like to supply 18k cubic feet, but there are no buyers—a situation of extreme excess supply. At the price of \$7 the excess supply is reduced to 9k, because both the quantity demanded is now higher at 3k units, and the quantity supplied is lower at 12k. But excess supply means that there are suppliers willing to supply at a lower price, and this willingness exerts continual downward pressure on any price above the price that equates demand and supply.

At prices below the equilibrium there is, conversely, an excess demand. In this situation, suppliers could force the price upward, knowing that buyers will continue to buy at a price at which the suppliers are willing to sell. Such upward pressure would continue until the excess demand is eliminated.

In general then, above the equilibrium price excess supply exerts downward pressure on price, and below the equilibrium excess demand exerts upward pressure on price. This process implies that the buyers and sellers have information on the various elements that make up the marketplace.

We will explore later in this chapter some specific circumstances in which trading could take place at prices above or below the equilibrium price. In such situations the quantity actually traded always corresponds to the short side of the market: this means that at high prices the quantity demanded is less than the quantity supplied, and it is the quantity demanded that is traded because buyers will not buy the amount suppliers would like to supply. Correspondingly, at low prices the quantity demanded exceeds quantity supplied, and it is the amount that suppliers are willing to sell that is traded. In sum, when trading takes place at prices other than the equilibrium price it is always the lesser of the quantity demanded or supplied that is traded. Hence we say that at non-equilibrium prices the short side dominates. We will return to this in a series of examples later in this chapter.

The **short side of the market** determines outcomes at prices other than the equilibrium.

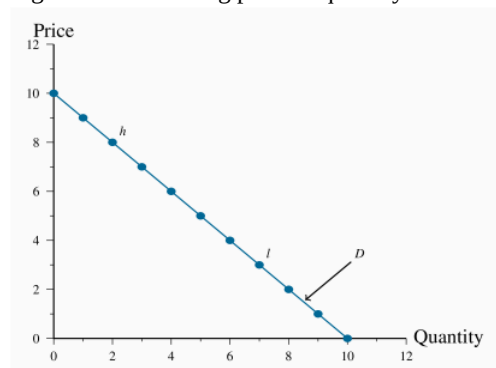
Supply and the nature of costs

Before progressing to a graphical analysis, we should add a word about costs. The supply schedules are based primarily on the cost of producing the product in question, and we frequently assume that all of the costs associated with supply are incorporated in the supply schedules. In Chapter 6 we will explore cases where costs additional to those incurred by producers may be relevant. For example, coal burning power plants emit pollutants into the atmosphere; but the individual supplier may not take account of these pollutants, which are costs to society at large, in deciding how much to supply at different prices. Stated another way, the private costs of production would not reflect the total, or full social costs of production. Conversely, if some individuals immunize themselves against a rampant virus, other individuals gain from that action because they become less likely to contract the virus - the social value thus exceeds the private value. For the moment the assumption is that no such additional costs are associated with the markets we analyze.

3.3 Demand and supply curves

The demand curve is a graphical expression of the relationship between price and quantity demanded, holding other things constant. Figure 3.1 measures price on the vertical axis and quantity on the horizontal axis. The curve D represents the data from the first two columns of Table 3.1. Each combination of price and quantity demanded lies on the curve. In this case the curve is *linear*—it is a straight line. The demand curve slopes downward (technically we say that its slope is negative), reflecting the fact that buyers wish to purchase more when the price is less.

Figure 3.1 Measuring price & quantity



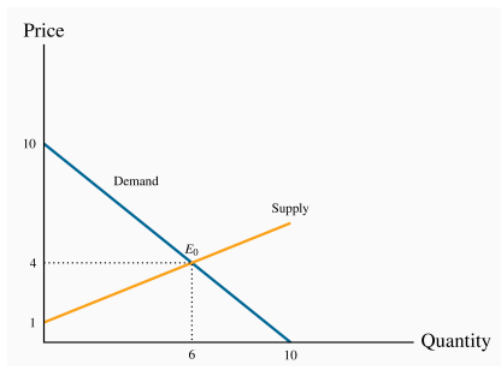
To derive this demand curve we take each price-quantity combination from the demand schedule in Table 3.1 and insert a point that corresponds to those combinations. For example, point h defines the combination $\{P = \$8, Q_d = 2\}$, the point l denotes the combination $\{P = \$3, Q_d = 7\}$. If we join all such points we obtain the demand curve in Figure 3.2. The same process yields the supply curve in Figure 3.2. In this example the supply and the demand curves are each linear. There is no reason why this linear property characterizes demand and supply curves in the real world; they are frequently found to have curvature. But straight lines are easier to work with, so we continue with them for the moment.

The **demand curve** is a graphical expression of the relationship between price and quantity demanded, with other influences remaining unchanged.

The supply curve is a graphical representation of the relationship between price and quantity supplied, holding other things constant. The supply curve S in Figure 3.2 is based on the data from columns 1 and 3 in Table 3.1. It has a positive slope indicating that suppliers wish to supply more at higher prices.

The **supply curve** is a graphical expression of the relationship between price and quantity supplied, with other influences remaining unchanged.

Figure 3.2 Supply, demand, equilibrium



The demand and supply curves intersect at point E_0 , corresponding to a price of \$4 which, as illustrated above, is the equilibrium price for this market. At any price below this the horizontal distance between the supply and demand curves represents excess demand, because demand exceeds supply. Conversely, at any price above \$4 there is an excess supply that is again measured by the horizontal distance between the two curves. Market forces tend to eliminate excess demand and excess supply as we explained above. In the final section of the chapter we illustrate how the supply and demand curves can be 'solved' for the equilibrium price and quantity.

3.4 Non-price influences on demand

We have emphasized several times the importance of the *ceteris paribus* assumption when exploring the impact of different prices on the quantity demanded: We assume all other influences on the purchase decision are unchanged (at least momentarily). These other influences fall into several broad categories: The prices of related goods; the incomes of buyers; buyer tastes; and expectations about the future. Before proceeding, note that we are dealing with *market* demand rather than demand by one *individual* (the precise relationship between the two is developed later in this chapter).

The prices of related goods – oil and gas, Kindle and paperbacks

We expect that the price of other forms of energy would impact the price of natural gas. For example, if hydro-electricity, oil or solar becomes less expensive we would expect some buyers to switch to these other products. Alternatively, if gas-burning furnaces experience a technological breakthrough that makes them more efficient and cheaper we would expect some users of other fuels to move to gas. Among these examples, oil and electricity are substitute fuels for gas; in contrast a more fuel-efficient new gas furnace complements the use of gas. We use these terms, substitutes and complements, to describe products that influence the demand for the primary good.

Substitute goods: when a price reduction (rise) for a related product reduces (increases) the demand for a primary product, it is a substitute for the primary product.

Complementary goods: when a price reduction (rise) for a related product increases (reduces) the demand for a primary product, it is a complement for the primary product.

Clearly electricity is a substitute for gas in the power market, whereas a gas furnace is a complement for gas as a fuel. The words substitutes and complements immediately suggest the nature of the relationships. Every product has complements and substitutes. As another example: Electronic readers and tablets are substitutes for paper-form books; a rise in the price of paper books should increase the demand for electronic readers at any given price for electronic readers. In graphical terms, the demand curve *shifts* in response to changes in the prices of other goods – an increase in the price of paper-form books shifts the demand for electronic readers outward, because more electronic readers will be demanded at any price.

Buyer incomes – which goods to buy

The demand for most goods increases in response to income growth. Given this, the demand curve for gas will shift outward if household incomes in the economy increase. Household incomes may increase either because there are more households in the economy or because the incomes of the existing households grow.

Most goods are demanded in greater quantity in response to higher incomes at any given price. But there are exceptions. For example, public transit demand may decline at any price when household incomes rise, because some individuals move to cars. Or the demand for laundromats may decline in response to higher incomes, as households purchase more of their own consumer

durables – washers and driers. We use the term inferior good to define these cases: An inferior good is one whose demand declines in response to increasing incomes, whereas a normal good experiences an increase in demand in response to rising incomes.

An **inferior good** is one whose demand falls in response to higher incomes.

A **normal good** is one whose demand increases in response to higher incomes.

There is a further sense in which consumer incomes influence demand, and this relates to how the incomes are *distributed* in the economy. In the discussion above we stated that higher total incomes shift demand curves outwards when goods are normal. But think of the difference in the demand for electronic readers between Portugal and Saudi Arabia. These economies have roughly the same average per-person income, but incomes are distributed more unequally in Saudi Arabia. It does not have a large middle class that can afford electronic readers or iPads, despite the huge wealth held by the elite. In contrast, Portugal has a relatively larger middle class that can afford such goods. Consequently, the *distribution of income* can be an important determinant of the demand for many commodities and services.

Tastes and networks – hemlines, lapels and homogeneity

While demand functions are drawn on the assumption that tastes are constant, in an evolving world they are not. We are all subject to peer pressure, the fashion industry, marketing, and a desire to maintain our image. If the fashion industry dictates that lapels on men's suits or long skirts are *de rigueur* for the coming season, some fashion-conscious individuals will discard a large segment of their wardrobe, even though the clothes may be in perfectly good condition: Their demand is influenced by the dictates of current fashion.

Correspondingly, the items that other individuals buy or use frequently determine our own purchases. Businesses frequently decide that all of their employees will have the same type of computer and software on account of *network economies*: It is easier to communicate if equipment is compatible, and it is less costly to maintain infrastructure where the variety is less.

Expectations – betting on the future

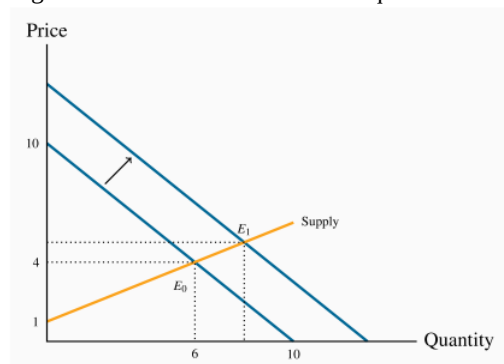
In our natural gas example, if households expected that the price of natural gas was going to stay relatively low for many years – perhaps on account of the discovery of large deposits – then they would be tempted to purchase a gas burning furnace rather than one based upon an alternative fuel. In this example, it is more than the current price that determines choices; *the prices that are expected to prevail in the future* also determine current demand.

Expectations are particularly important in stock markets. When investors anticipate that corporations will earn high rewards in the future they will buy a stock today. If enough people believe this, the price of the stock will be driven upward on the market, even before profitable earnings are registered.

Shifts in demand

The demand curve in Figure 3.2 is drawn for a given level of other prices, incomes, tastes, and expectations. Movements along the demand curve reflect solely the impact of different prices for the good in question, holding other influences constant. But changes in any of these other factors will change the position of the demand curve. Figure 3.3 illustrates a shift in the demand curve. This shift could result from a rise in household incomes that increase the quantity demanded *at every price*. This is illustrated by an outward shift in the demand curve. With supply conditions unchanged, there is a new equilibrium at E_1 , indicating a greater quantity of purchases accompanied by a higher price. The new equilibrium reflects a *change in quantity supplied and a change in demand*.

Figure 3.3 Demand shift and new equilibrium



The outward shift in demand leads to a new equilibrium E_1 .

We may well ask why so much emphasis in our diagrams and analysis is placed on the relationship between *price* and quantity, rather than on the relationship between quantity and its other determinants. The answer is that we could indeed draw diagrams with quantity on the horizontal axis and a measure of one of these other influences on the vertical axis. But the price mechanism plays a very important role. *Variations in price are what equilibrate the market.* By focusing primarily upon the price, we see the self-correcting mechanism by which the market reacts to excess supply or excess demand.

In addition, this analysis illustrates the method of comparative statics—examining the impact of changing one of the other things that are assumed constant in the supply and demand diagrams.

Comparative static analysis compares an initial equilibrium with a new equilibrium, where the difference is due to a change in one of the other things that lie behind the demand curve or the supply curve.

'Comparative' obviously denotes the idea of a comparison, and static means that we are not in a state of motion. Hence we use these words in conjunction to indicate that we compare one outcome with another, without being concerned too much about the transition from an initial equilibrium to a new equilibrium. The transition would be concerned with dynamics rather than statics. In Figure 3.3 we explain the difference between the points E_0 and E_1 by indicating that there has been a change in incomes or in the price of a substitute good. We do not attempt to analyze the details of this move or the exact path from E_0 to E_1 .

Application Box 3.1 Corn prices and demand shifts

In the middle of its second mandate, the Bush Administration in the US decided to encourage the production of ethanol – a fuel that is less polluting than gasoline. The target production was 35 billion for 2017 – from a base of 1 billion gallons in 2000. Corn is the principal input in ethanol production. It is also used as animal feed, as a sweetener and as a food for humans. The target was to be met with the help of a subsidy to producers and a tariff on imports of Brazil's sugar-cane based ethanol.

The impact on corn prices was immediate; from a farm-gate price of \$2 per bushel in 2005, the price reached the \$4 range two years later. In 2012 the price rose temporarily to \$7. While other factors were in play - growing incomes and possibly speculation by commodity investors, ethanol is seen as the main price driver: demand for corn increased and the supply could not be increased to keep up with the demand without an increase in price.

The wider impact of these developments was that the prices of virtually all grains increased in tandem with corn: the prices of sorghum and barley increased because of a switch in land use towards corn on account of its profitability.

While farmers benefited from the price rise, consumers – particularly those in less developed economies – experienced a dramatic increase in their basic living costs. Visit the site of the United Nations' Food and Agricultural Organization for an assessment. Since hitting \$7 per bushel in 2012, the price has dropped and averaged \$3.50 in 2016.

In terms of supply and demand shifts: the demand side has dominated, particularly in the short run. The ethanol drive, combined with secular growth in the demand for food, means that the demand for grains shifted outward faster than the supply. In the period 2013–2016, supply has increased and the price has moderated.

3.5 Non-price influences on supply

To date we have drawn supply curves with an upward slope. Is this a reasonable representation of supply in view of what is frequently observed in markets? We suggested earlier that the various producers of a particular good or service may have different levels of efficiency. If so, only the more efficient producers can make a profit at a low price, whereas at higher prices more producers or suppliers enter the market – producers who may not be as lean and efficient as those who can survive in a lower-price environment. This view of the world yields a positively-sloping supply curve.

As a second example, consider *Uber* or *Lyft* taxi drivers. Some drivers may be in serious need of income and may be willing to drive for a low hourly rate. For other individuals driving may be a secondary source of income, and such drivers are less likely to want to drive unless the hourly wage is higher. Consequently if these ride sharing services need a large number of drivers at any one time it may be necessary to pay a higher wage – *and charge a higher fare to passengers*, to induce more drivers to take their taxis onto the road. This phenomenon corresponds to a positively-sloped supply curve.

In contrast to these two examples, some suppliers simply choose a unique price and let buyers purchase as much as they want at that price. This is the practice of most retailers. For example, the price of *Samsung's Galaxy* is typically fixed, no matter how many are purchased – and tens of millions are sold at a fixed price when a new model is launched. *Apple* also sets a price, and buyers purchase as many as they desire at that price. This practice corresponds to a horizontal supply curve: The price does not vary and the market equilibrium occurs where the demand curve intersects this supply curve.

In yet other situations supply is fixed. This happens in auctions. Bidders at the auction simply determine the price to be paid. At a real estate auction a given property is put on the market and the price is determined by the bidding process. In this case the supply of a single property is represented by a vertical supply at a quantity of 1 unit.

Regardless of the type of market we encounter, however, it is safe to assume that supply curves rarely slope downward. So, for the moment, we adopt the stance that supply curves are generally upward sloping – somewhere between the extremes of being vertical or horizontal – as we have drawn them to this point.

Next, we examine those other influences that underlie supply curves. Technology, input costs, the prices of competing goods, expectations and the number of suppliers are the most important.

Technology – computers and fracking

A technological advance may involve an idea that allows more output to be produced with the same inputs, or an equal output with fewer inputs. A good example is *just-in-time* technology. Before the modern era, virtually all manufacturers kept large stocks of components in their production facilities, but developments in communications and computers at that time made it possible for manufacturers to link directly with their input suppliers. Nowadays auto assembly plants place their order for, say, seat delivery to their local seat supplier well ahead of assembly time. The seats swing into the assembly area hours or minutes before assembly—just in time. The result is that the assembler reduces her seat inventory (an input) and thereby reduces production cost.

Such a technology-induced cost saving is represented by moving the supply curve downward or outward: The supplier is now able and willing to supply the same quantity at a lower price because of the technological innovation. Or, saying the same thing slightly differently, suppliers will supply more at a given price than before.

A second example relates to the extraction of natural gas. The development of 'fracking' means that companies involved in gas recovery can now do so at a lower cost. Hence they are willing to supply any given quantity at a lower price. A third example concerns aluminum cans. Today they weigh a fraction of what they weighed 20 years ago. This is a technology-based cost saving.

Input costs

Input costs can vary independently of technology. For example, a wage negotiation that grants workers a substantial pay raise will increase the cost of production. This is reflected in a leftward, or upward, supply shift: Any quantity supplied is now priced higher; alternatively, suppliers are willing to supply less at the going price.

Production costs may increase as a result of higher required standards in production. As governments implement new safety or product-stress standards, costs may increase. In this instance the increase in costs is not a 'bad' outcome for the buyer. She may be purchasing a higher quality good as a result.

Competing products – Airbnb versus hotels

If competing products improve in quality or fall in price, a supplier may be forced to follow suit. For example, *Asus* and *Dell* are constantly watching each other's pricing policies. If *Dell* brings out a new generation of computers at a lower price, *Asus* may lower its prices in turn—which is to say that *Asus*' supply curve will shift downward. Likewise, *Samsung* and *Apple* each responds to the other's pricing and technology behaviours. The arrival of new products in the marketplace also impacts the willingness of suppliers to supply goods at a given price. New intermediaries such as *Airbnb* and *Vacation Rentals by Owner* have shifted the supply curves of hotel rooms downward.

These are some of the many factors that influence the position of the supply curve in a given market.

Application Box 3.2 The price of light

Technological developments have had a staggering impact on many price declines. Professor William Nordhaus of Yale University is an expert on measuring technological change. He has examined the trend in the real price of lighting. Originally, light was provided by whale oil and gas lamps and these sources of lumens (the scientific measure of the amount of light produced) were costly. In his research, Professor Nordhaus pieced together evidence on the actual historic cost of light produced at various times, going all the way back to 1800. He found that light in 1800 cost about 100 times more than in 1900, and light in the year 2000 was a fraction of its cost in 1900. A rough calculation suggests that light was five hundred times more expensive at the start of this 200-year period than at the end, and this was before the arrival of LEDs.

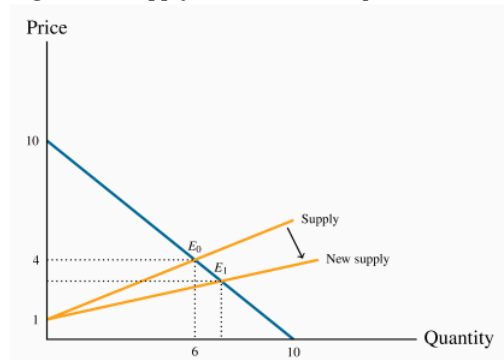
In terms of supply and demand analysis, light has been subject to very substantial downward supply shifts. Despite the long-term growth in demand, the technologically-induced supply changes have been the dominant factor in its price determination.

For further information, visit Professor Nordhaus's website in the Department of Economics at Yale University.

Shifts in supply

Whenever technology changes, or the costs of production change, or the prices of competing products adjust, then one of our *ceteris paribus* assumptions is violated. Such changes are generally reflected by shifting the supply curve. Figure 3.4 illustrates the impact of the arrival of just-in-time technology. The supply curve shifts, reflecting the ability of suppliers to supply the same output at a reduced price. The resulting new equilibrium price is lower, since production costs have fallen. At this reduced price more gas is traded at a lower price.

Figure 3.4 Supply shift and new equilibrium



The supply curve shifts due to lower production costs. A new equilibrium E_1 is attained in the market at a lower price.

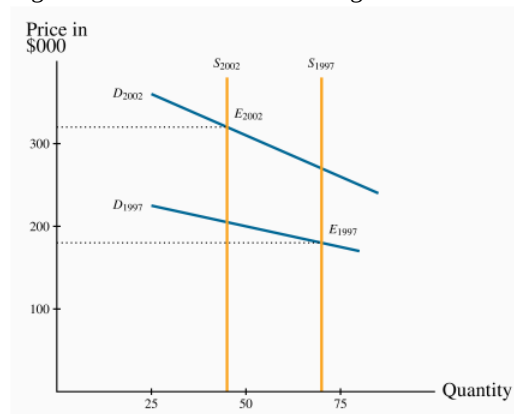
3.6 Simultaneous supply and demand impacts

In the real world, demand and supply frequently shift at the same time. We present such a case in Figure 3.5. It is based upon real estate data describing the housing market in a small Montreal municipality. Vertical curves define the supply side of the market. Such vertical curves mean that a given number of homeowners decide to put their homes on the market, and these suppliers just take whatever price results in the market. In this example, fewer houses were offered for sale in 2002 (less than 50) than in 1997 (more than 70). We are assuming in this market that the houses traded were similar; that is, we are not lumping together mansions with row houses.

During this time period household incomes increased substantially and, also, mortgage rates fell. Both of these developments shifted the demand curve upward/outward: Buyers were willing to pay more for housing in 2002 than in 1997, both because their incomes were on average higher and because they could borrow more cheaply.

The shifts on both sides of the market resulted in a higher average price. And each of these shifts compounded the other: The outward shift in demand would lead to a higher price on its own, and a reduction in supply would do likewise. Hence both forces acted to push up the price in 2002. If, instead, the supply had been greater in 2002 than in 1997 this would have acted to reduce the equilibrium price. And with the demand and supply shifts operating in opposing directions, it is not possible to say in general whether the price would increase or decrease. If the demand shift were strong and the supply shift weak then the demand forces would have dominated and led to a higher price. Conversely, if the supply forces were stronger than the demand forces,

Figure 3.5 A model of the housing market with shifts in demand and supply



The vertical supply denotes a fixed number of houses supplied each year. Demand was stronger in 2002 than in 1997 both on account of higher incomes and lower mortgage rates. Thus the higher price in 2002 is due to both a reduction in supply and an increase in demand.

3.7 Market interventions – governments and interest groups

The freely functioning markets that we have developed certainly do not describe all markets. For example, minimum wages characterize the labour market, most agricultural markets have supply restrictions, apartments are subject to rent controls, and blood is not a freely traded market commodity in Canada. In short, price controls and quotas characterize many markets. Price controls are government rules or laws that inhibit the formation of market-determined prices. Quotas are physical restrictions on how much output can be brought to the market.

Price controls are government rules or laws that inhibit the formation of market-determined prices.

Quotas are physical restrictions on output.

Price controls come in the form of either *floors* or *ceilings*. Price floors are frequently accompanied by *marketing boards*.

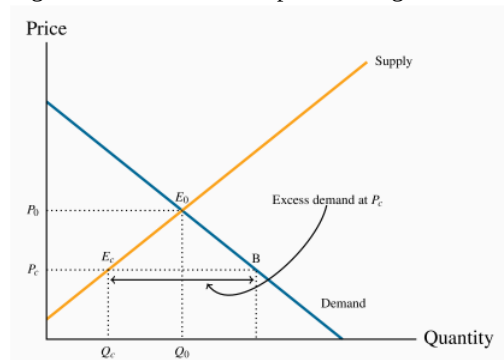
Price ceilings – rental boards

Ceilings mean that suppliers cannot legally charge more than a specific price. Limits on apartment rents are one form of ceiling. In times of emergency – such as flooding or famine, price controls are frequently imposed on foodstuffs, in conjunction with rationing, to ensure that access is not determined by who has the most income. The problem with price ceilings, however, is that they leave demand unsatisfied, and therefore they must be accompanied by some other allocation mechanism.

Consider an environment where, for some reason – perhaps a sudden and unanticipated growth in population – rents increase. Let the resulting equilibrium be defined by the point E_0 in Figure 3.6. If the government were to decide that this is an unfair price because it places hardships on low- and middle-income households, it might impose a price limit, or ceiling, of P_c . The problem with such a limit is that excess demand results: Individuals want to rent more apartments than are available in the city. In a free market the price would adjust upward to eliminate the excess demand, but in this controlled environment it cannot. So some other way of allocating the available supply between demanders must evolve.

In reality, most apartments are allocated to those households already occupying them. But what happens when such a resident household decides to purchase a home or move to another city? In a free market, the landlord could increase the rent in accordance with market pressures. But in a controlled market a city's rental tribunal may restrict the annual rent increase to just a couple of percent and the demand may continue to outstrip supply. So how does the stock of apartments get allocated between the potential renters? One allocation method is well known: The existing tenant informs her friends of her plan to move, and the friends are the first to apply to the landlord to occupy the apartment. But that still leaves much unmet demand. If this is a student rental market, students whose parents live nearby may simply return 'home'. Others may choose to move to a part of the city where rents are more affordable.

Figure 3.6 The effect of a price ceiling



The free market equilibrium occurs at E_0 . A price ceiling at P_c holds down the price but leads to excess demand E_cB , because Q_c is the quantity traded. A price ceiling above P_0 is irrelevant since the free market equilibrium E_0 can still be attained.

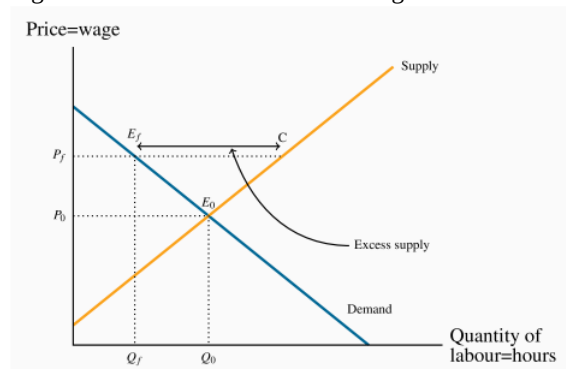
However, rent controls sometimes yield undesirable outcomes. Rent controls are widely studied in economics, and the consequences are well understood: Landlords tend not to repair or maintain their rental units in good condition if they cannot obtain the rent they believe they are entitled to. Accordingly, the residential rental stock deteriorates. In addition, builders realize that more money is to be made in building condominium units than rental units, or in *converting rental units to condominiums*. The

frequent consequence is thus a *reduction* in supply and a reduced quality. Market forces are hard to circumvent because, as we emphasized in Chapter 1, economic players react to the incentives they face. These outcomes are examples of what we call the *law of unintended consequences*.

Price floors – minimum wages

An effective price floor sets the price *above* the market-clearing price. A minimum wage is the most widespread example in the Canadian economy. Provinces each set their own minimum, and it is seen as a way of protecting the well-being of low-skill workers. Such a floor is illustrated in Figure 3.7. The free-market equilibrium is again E_0 , but the effective market outcome is the combination of price and quantity corresponding to the point E_f at the price floor, P_f . In this instance, there is excess supply equal to the amount E_fC .

Figure 3.7 Price floor – minimum wage



In a free market the equilibrium is E_0 . A minimum wage of P_f raises the hourly wage, but reduces the hours demanded to Q_f . Thus E_fC is the excess supply.

Note that there is a similarity between the outcomes defined in the floor and ceiling cases: The quantity actually traded is *the lesser of the supply quantity and demand quantity at the going price: The short side dominates*.

If price floors, in the form of minimum wages, result in some workers going unemployed, why do governments choose to put them in place? The excess supply in this case corresponds to unemployment – more individuals are willing to work for the going wage than buyers (employers) wish to employ. The answer really depends upon the magnitude of the excess supply. In particular, suppose, in Figure 3.7 that the supply and demand curves going through the equilibrium E_0 were more 'vertical'. This would result in a smaller excess supply than is represented with the existing supply and demand curves. This would mean in practice that a higher wage could go to workers, making them better off, without causing substantial unemployment. This is the trade off that governments face: With a view to increasing the purchasing power of generally lower-skill individuals, a minimum wage is set, hoping that the negative impact on employment will be small. We will return to this in the next chapter, where we examine the responsiveness of supply and demand curves to different prices.

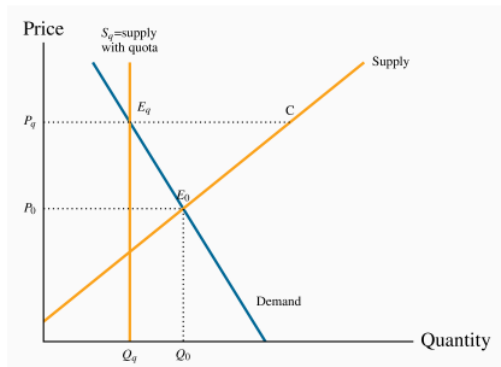
Quotas – agricultural supply

A quota represents the right to supply a specified quantity of a good to the market. It is a means of keeping prices higher than the free-market equilibrium price. As an alternative to imposing a price floor, the government can generate a high price by restricting supply.

Agricultural markets abound with examples. In these markets, farmers can supply only what they are permitted by the quota they hold, and there is usually a market for these quotas. For example, in several Canadian provinces it currently costs in the region of \$30,000 to purchase a quota granting the right to sell the milk of one cow. The cost of purchasing quotas can thus easily outstrip the cost of a farm and herd. Canadian cheese importers must pay for the right to import cheese from abroad. Restrictions also apply to poultry. The impact of all of these restrictions is to raise the domestic price above the free market price.

In Figure 3.8, the free-market equilibrium is at E_0 . In order to raise the price above P_0 , the government restricts supply to Q_q by granting quotas, which permit producers to supply a limited amount of the good in question. This supply is purchased at the price equal to P_q . From the standpoint of farmers, a higher price might be beneficial, even if they get to supply a smaller quantity, provided the amount of revenue they get as a result is as great as the revenue in the free market.

Figure 3.8 The effect of a quota



The government decides that the equilibrium price P_0 is too low. It decides to boost price by reducing supply from Q_0 to Q_q . It achieves this by requiring producers to have a production quota. This is equivalent to fixing supply at S_q .

Marketing boards – milk and maple syrup

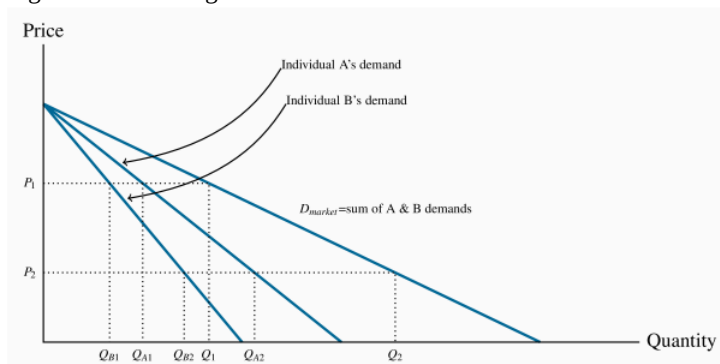
A marketing board is a means of insuring that a quota or price floor can be maintained. Quotas are frequent in the agriculture sector of the economy. One example is maple syrup in Quebec. The Federation of Maple Syrup Producers of Quebec has the sole right to market maple syrup. All producers must sell their syrup through this marketing board. The board thus has a particular type of power in the market: it has control of the market at the wholesale end, because it is a sole buyer. The Federation increases the total revenue going to producers by artificially restricting the supply to the market. The Federation calculates that by reducing supply and selling it at a higher price, more revenue will accrue to the producers. This is illustrated in Figure 3.8. The market equilibrium is given by E_0 , but the Federation restricts supply to the quantity Q_q , which is sold to buyers at price P_q . To make this possible the total supply must be restricted; otherwise producers would supply the amount given by the point C on the supply curve, and this would result in excess supply in the amount E_qC . In order to restrict supply to Q_q in total, individual producers are limited in what they can sell to the Federation; they have a quota, which gives them the right to produce and sell no more than a specified amount. This system of quotas is necessary to eliminate the excess supply that would emerge at the above-equilibrium price P_q .

We will return to this topic in Chapter 4. For the moment, to see that this type of revenue-increasing outcome is possible, examine Table 3.1 again. At this equilibrium price of \$4 the quantity traded is 6 units, yielding a total expenditure by buyers (revenue to suppliers) of \$24. However, if the supply were restricted and a price of \$5 were set, the expenditure by buyers (revenue to suppliers) would rise to \$25.

3.8 Individual and market functions

Markets are made up of many individual participants on the demand and supply side. The supply and demand functions that we have worked with in this chapter are those for the total of all participants on each side of the market. But how do we arrive at such market functions when the economy is composed of individuals? We can illustrate how, with the help of Figure 3.9.

Figure 3.9 Summing individual demands



At P_1 individual A purchases Q_{A1} and B purchases Q_{B1} . The total demand is the sum of these individual demands at this price (Q_1). At P_2 individual demands are summed to Q_2 . Since the points Q_1 and Q_2 define the demands of the market participants it follows that market demand is the horizontal sum of these curves.

To concentrate on the essentials, imagine that there are just two buyers of chocolate cookies in the economy. A has a stronger preference for cookies than B, so his demand is greater. To simplify, let the two demands have the same intercept on the vertical

axis. The curves D_A and D_B indicate how many cookies A and B, respectively, will buy at each price. The market demand indicates how much they buy *together* at any price. Accordingly, at P_1 , A and B purchase the quantities Q_{A1} and Q_{B1} respectively. Thus $Q_1 = Q_{A1} + Q_{B1}$. At a price P_2 , they purchase Q_{A2} and Q_{B2} . Thus $Q_2 = Q_{A2} + Q_{B2}$. The market demand is therefore the horizontal sum of the individual demands at these prices. In the figure this is defined by D_{market} .

Market demand: the horizontal sum of individual demands.

3.9 Useful techniques – demand and supply equations

The supply and demand functions, or equations, underlying Table 3.1 and Figure 3.2 can be written in their mathematical form:

Demand: $P = 10 - Q$

Supply: $P = 1 + (1/2)Q$

A straight line is represented completely by the intercept and slope. In particular, if the variable P is on the vertical axis and Q on the horizontal axis, the straight-line equation relating P and Q is defined by $P = a + bQ$. Where the line is negatively sloped, as in the demand equation, the parameter b must take a negative value. By observing either the data in Table 3.1 or Figure 3.2 it is clear that the vertical intercept, a , takes a value of \$10. The vertical intercept corresponds to a zero-value for the Q variable. Next we can see from Figure 3.2 that the slope (given by the rise over the run) is 10/10 and hence has a value of -1 . Accordingly the demand equation takes the form $P = 10 - Q$.

On the supply side the price-axis intercept, from either the figure or the table, is clearly 1. The slope is one half, because a two-unit change in quantity is associated with a one-unit change in price. This is a positive relationship obviously so the supply curve can be written as $P = 1 + (1/2)Q$.

Where the supply and demand curves intersect is the market equilibrium; that is, the price-quantity combination is the same for both supply and demand where the supply curve takes on the same values as the demand curve. This unique price-quantity combination is obtained by equating the two curves: If Demand=Supply, then

$$10 - Q = 1 + (1/2)Q.$$

Gathering the terms involving Q to one side and the numerical terms to the other side of the equation results in $9 = 1.5Q$. This implies that the equilibrium quantity must be 6 units. And this quantity must trade at a price of \$4. That is, when the price is \$4 both the quantity demanded and the quantity supplied take a value of 6 units.

Modelling market interventions using equations

To illustrate the impact of market interventions examined in Section 3.7 on our numerical market model for natural gas, suppose that the government imposes a minimum price of \$6 – above the equilibrium price obviously. We can easily determine the quantity supplied and demanded at such a price. Given the supply equation

$$P = 1 + (1/2)Q,$$

it follows that at $P = 6$ the quantity supplied is 10. This follows by solving the relationship $6 = 1 + (1/2)Q$ for the value of Q . Accordingly, suppliers *would like to supply* 10 units at this price.

Correspondingly on the demand side, given the demand curve

$$P = 10 - Q,$$

with a price given by $P = \$6$, it must be the case that $Q = 4$. So buyers *would like to buy* 4 units at that price: There is excess supply. But we know that the short side of the market will win out, and so the actual amount traded at this restricted price will be 4 units.

Conclusion

We have covered a lot of ground in this chapter. It is intended to open up the vista of economics to the new student in the discipline. Economics is powerful and challenging, and the ideas we have developed here will serve as conceptual foundations for our exploration of the subject. Our next chapter deals with measurement and responsiveness.

Key Terms

Demand is the quantity of a good or service that buyers wish to purchase at each possible price, with all other influences on demand remaining unchanged.

Supply is the quantity of a good or service that sellers are willing to sell at each possible price, with all other influences on supply remaining unchanged.

Quantity demanded defines the amount purchased at a particular price.

Quantity supplied refers to the amount supplied at a particular price.

Equilibrium price: equilibrates the market. It is the price at which quantity demanded equals the quantity supplied.

Excess supply exists when the quantity supplied exceeds the quantity demanded at the going price.

Excess demand exists when the quantity demanded exceeds quantity supplied at the going price.

Short side of the market determines outcomes at prices other than the equilibrium.

Demand curve is a graphical expression of the relationship between price and quantity demanded, with other influences remaining unchanged.

Supply curve is a graphical expression of the relationship between price and quantity supplied, with other influences remaining unchanged.

Substitute goods: when a price reduction (rise) for a related product reduces (increases) the demand for a primary product, it is a substitute for the primary product.

Complementary goods: when a price reduction (rise) for a related product increases (reduces) the demand for a primary product, it is a complement for the primary product.

Inferior good is one whose demand falls in response to higher incomes.

Normal good is one whose demand increases in response to higher incomes.

Comparative static analysis compares an initial equilibrium with a new equilibrium, where the difference is due to a change in one of the other things that lie behind the demand curve or the supply curve.

Price controls are government rules or laws that inhibit the formation of market-determined prices.

Quotas are physical restrictions on output.

Market demand: the horizontal sum of individual demands.

Exercises for Chapter 3

EXERCISE 3.1

The supply and demand for concert tickets are given in the table below.

Price (\$)	0	4	8	12	16	20	24	28	32	36	40
Quantity demanded	15	14	13	12	11	10	9	8	7	6	5
Quantity supplied	0	0	0	0	0	1	3	5	7	9	11

1. Plot the supply and demand curves to scale and establish the equilibrium price and quantity.
2. What is the excess supply or demand when price is \$24? When price is \$36?
3. Describe the market adjustments in price induced by these two prices.
4. *Optional:* The functions underlying the example in the table are linear and can be presented as $P=18+2Q$ (supply) and $P=60-4Q$ (demand). Solve the two equations for the equilibrium price and quantity values.

EXERCISE 3.2

Illustrate in a supply/demand diagram, by shifting the demand curve appropriately, the effect on the demand for flights between Calgary and Winnipeg as a result of:

1. Increasing the annual government subsidy to *Via Rail*.

2. Improving the Trans-Canada highway between the two cities.
3. The arrival of a new budget airline on the scene.

EXERCISE 3.3

A new trend in US high schools is the widespread use of chewing tobacco. A recent survey indicates that 15 percent of males in upper grades now use it – a figure not far below the use rate for cigarettes. This development came about in response to the widespread implementation by schools of regulations that forbade cigarette smoking on and around school property. Draw a supply-demand equilibrium for each of the cigarette and chewing tobacco markets before and after the introduction of the regulations.

EXERCISE 3.4

The following table describes the demand and supply conditions for labour.

Price (\$) = wage rate	0	10	20	30	40	50	60	70	80	90	100	110	120	130	140	150	160	170
Quantity demanded	1020	960	900	840	780	720	660	600	540	480	420	360	300	240	180	120	60	0
Quantity supplied	0	0	0	0	0	0	30	60	90	120	150	180	210	240	270	300	330	360

1. Graph the functions and find the equilibrium price and quantity by equating demand and supply.
2. Suppose a price ceiling is established by the government at a price of \$120. This price is below the equilibrium price that you have obtained in part (a). Calculate the amount that would be demanded and supplied and then calculate the excess demand.

EXERCISE 3.5

In Exercise 3.4, suppose that the supply and demand describe an agricultural market rather than a labour market, and the government implements a price floor of \$140. This is greater than the equilibrium price.

1. Estimate the quantity supplied and the quantity demanded at this price, and calculate the excess supply.
2. Suppose the government instead chose to maintain a price of \$140 by implementing a system of quotas. What quantity of quotas should the government make available to the suppliers?

EXERCISE 3.6

In Exercise 3.5, suppose that, at the minimum price, the government buys up all of the supply that is not demanded, and exports it at a price of \$80 per unit. Compute the cost to the government of this operation.

EXERCISE 3.7

Let us sum two demand curves to obtain a 'market' demand curve. We will suppose there are just two buyers in the market. Each of the individual demand curves has a price intercept of \$42. One has a quantity intercept of 126, the other 84.

1. Draw the demands either to scale or in an Excel spreadsheet, and label the intercepts on both the price and quantity axes.
2. Determine how much would be purchased in the market at prices \$10, \$20, and \$30.
3. *Optional:* Since you know the intercepts of the market (total) demand curve, can you write an equation for it?

EXERCISE 3.8

In Exercise 3.7 the demand curves had the same price intercept. Suppose instead that the first demand curve has a price intercept of \$36 and a quantity intercept of 126; the other individual has a demand curve defined by a price intercept of \$42 and a quantity intercept of 84. Graph these curves and illustrate the market demand curve.

EXERCISE 3.9

Here is an example of a demand curve that is not linear:

Price (\$)	4	3	2	1	0
Quantity demanded	25	100	225	400	625

1. Plot this demand curve to scale or in Excel.
2. If the supply function in this market is $P=2$, plot this function in the same diagram.
3. Determine the equilibrium quantity traded in this market.

EXERCISE 3.10

The football stadium of the University of the North West Territories has 30 seats. The demand curve for tickets has a price intercept of \$36 and a quantity intercept of 72.

1. Draw the supply and demand curves to scale in a graph or in Excel. (This demand curve has the form $P = 36 - 0.5 \times Q$.)
2. Determine the equilibrium admission price, and the amount of revenue generated from ticket sales for each game.
3. A local alumnus and benefactor offers to install 6 more seats at no cost to the University. Compute the price that would be charged with this new supply and compute the revenue that would accrue at this new equilibrium price. Should the University accept the offer to install the seats?
4. Redo the previous part of this question, assuming that the initial number of seats is 40, and the University has the option to increase capacity to 46 at no cost to itself. Should the University accept the offer in this case?

EXERCISE 3.11

Suppose farm workers in Mexico are successful in obtaining a substantial wage increase. Illustrate the effect of this on the price of lettuce in the Canadian winter, using a supply and demand diagram, on the assumption that all lettuce in Canada is imported during its winter.

This page titled [3: The classical marketplace – demand and supply](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.1: The marketplace - trading

The marketplace in today's economy has evolved from earlier times. It no longer has a unique form – one where buyers and sellers physically come together for the purpose of exchange. Indeed, supermarkets usually require individuals to be physically present to make their purchases. But when purchasing an airline ticket, individuals simply go online and interact with perhaps a number of different airlines (suppliers) simultaneously. Or again, individuals may simply give an instruction to their stock broker, who will execute a purchase on their behalf – the broker performs the role of a middleman, who may additionally give advice to the purchaser. Or a marketing agency may decide to subcontract work to a translator or graphic artist who resides in Mumbai. The advent of the coronavirus has shifted grocery purchases from in-store presence to home delivery for many buyers. In pure auctions (where a single work of art or a single residence is offered for sale) buyers compete one against the other for the single item supplied. Accommodations in private homes are supplied to potential visitors (buyers) through *Airbnb*. Taxi rides are mediated through *Lyft* or *Uber*. These institutions are all different types of markets; they serve the purpose of facilitating exchange and trade.

Not all goods and services in the modern economy are obtained through the marketplace. Schooling and health care are allocated in Canada primarily by government decree. In some instances the market plays a supporting role: Universities and colleges may levy fees, and most individuals must pay, at least in part, for their pharmaceuticals. In contrast, broadcasting services may carry a price of zero; Buzzfeed or other news and social media come free of payment. Furthermore, some markets have no price, yet they find a way of facilitating an exchange. For example, graduating medical students need to be matched with hospitals for their residencies. Matching mechanisms are a form of market in that they bring together suppliers and demanders. We explore their operation in Chapter 11.

The importance of the marketplace springs from its role as an allocating mechanism. Elevated prices effectively send a signal to suppliers that the buyers in the market place a high value on the product being traded; conversely when prices are low. Accordingly, suppliers may decide to cease supplying markets where prices do not remunerate them sufficiently, and redirect their energies and the productive resources under their control to other markets – markets where the product being traded is more highly valued, and where the buyer is willing to pay more.

Whatever their form, the marketplace is central to the economy we live in. Not only does it facilitate trade, it also provides a means of earning a livelihood. Suppliers must hire resources – human and non-human – in order to bring their supplies to market and these resources must be paid a return: income is generated.

In this chapter we will examine the process of price formation – how the prices that we observe in the marketplace come to be what they are. We will illustrate that the price for a good is inevitably linked to the quantity of a good; price and quantity are different sides of the same coin and cannot generally be analyzed separately. To understand this process more fully, we need to *model* a typical market. The essentials are demand and supply.

This page titled [3.1: The marketplace - trading](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.2: The market's building blocks

In economics we use the terminology that describes trade in a particular manner. Non-economists frequently describe microeconomics by saying "it's all about supply and demand". While this is largely true we need to define exactly what we mean by these two central words. Demand is the quantity of a good or service that buyers wish to purchase at each conceivable price, with all other influences on demand remaining unchanged. It reflects a multitude of values, not a single value. It is not a single or unique quantity such as two cell phones, but rather a full description of the quantity of a good or service that buyers would purchase at various prices.

Demand is the quantity of a good or service that buyers wish to purchase at each possible price, with all other influences on demand remaining unchanged.

As a hypothetical example, the first column of Table 3.1 shows the price of natural gas per cubic foot. The second column shows the quantity that would be purchased in a given time period at each price. It is therefore a schedule of quantities demanded at various prices. For example, at a price \$6 per unit, buyers would like to purchase 4 units, whereas at the lower price of \$3 buyers would like to purchase 7 units. Note also that this is a homogeneous good. A cubic foot of natural gas is considered to be the same product no matter which supplier brings it to the market. In contrast, accommodations supplied through *Airbnb* are heterogeneous – they vary in size and quality.

Table 3.1 Demand and supply for natural gas

Price (\$)	Demand (thousands of cu feet)	Supply (thousands of cu feet)	Excess
10	0	18	Excess Supply
9	1	16	
8	2	14	
7	3	12	
6	4	10	
5	5	8	Equilibrium
4	6	6	
3	7	4	Excess Demand
2	8	2	
1	9	0	
0	10	0	

Supply is interpreted in a similar manner. It is not a single value; we say that supply is the quantity of a good or service that sellers are willing to sell at each possible price, with all other influences on supply remaining unchanged. Such a supply schedule is defined in the third column of the table. It is assumed that no supplier can make a profit (on account of their costs) unless the price is at least \$2 per unit, and therefore a zero quantity is supplied below that price. The higher price is more profitable, and therefore induces a greater quantity supplied, perhaps by attracting more suppliers. This is reflected in the data. For example, at a price of \$3 suppliers are willing to supply 4 units, whereas with a price of \$7 they are willing to supply 12 units. There is thus a positive relationship between price and quantity for the supplier – a higher price induces a greater quantity; whereas on the demand side of the market a higher price induces a lower quantity demanded – a negative relationship.

Supply is the quantity of a good or service that sellers are willing to sell at each possible price, with all other influences on supply remaining unchanged.

We can now identify a key difference in terminology – between the words demand and quantity demanded, and between supply and quantity supplied. While the words demand and supply refer to the complete schedules of demand and supply, the terms quantity demanded and quantity supplied each define a single value of demand or supply at a particular price.

Quantity demanded defines the amount purchased at a particular price.

Quantity supplied refers to the amount supplied at a particular price.

Thus while the non-economist may say that when some fans did not get tickets to the Stanley Cup it was a case of demand exceeding supply, as economists we say that the quantity demanded exceeded the quantity supplied *at the going price of tickets*. In this instance, had every ticket been offered at a sufficiently high price, the market could have generated an excess supply rather than an excess demand. A higher ticket price would reduce the *quantity demanded*; yet would not change *demand*, because demand refers to the whole schedule of possible quantities demanded at different prices.

Other things equal – *ceteris paribus*

The demand and supply schedules rest on the assumption that all other influences on supply and demand remain the same as we move up and down the possible price values. The expression *other things being equal*, or its Latin counterpart *ceteris paribus*, describes this constancy of other influences. For example, we assume on the demand side that the prices of other goods remain constant, and that tastes and incomes are unchanging. On the supply side we assume, for example, that there is no technological change in production methods. If any of these elements change then the market supply or demand schedules will reflect such changes. For example, if coal or oil prices increase (decline) then some buyers may switch to (away from) gas or solar power. This will be reflected in the data: At any given price more (or less) will be demanded. We will illustrate this in graphic form presently.

Market equilibrium

Let us now bring the demand and supply schedules together in an attempt to analyze what the marketplace will produce – will a single price emerge that will equate supply and demand? We will keep other things constant for the moment, and explore what materializes at different prices. At low prices, the data in Table 3.1 indicate that the quantity demanded exceeds the quantity supplied – for example, verify what happens when the price is \$3 per unit. The opposite occurs when the price is high – what would happen if the price were \$8? Evidently, there exists an intermediate price, where the quantity demanded equals the quantity supplied. At this point we say that the market is in equilibrium. The equilibrium price equates demand and supply – it clears the market.

The **equilibrium price** equilibrates the market. It is the price at which quantity demanded equals the quantity supplied.

In Table 3.1 the equilibrium price is \$4, and the equilibrium quantity is 6 thousand cubic feet of gas (we use the notation 'k' to denote thousands). At higher prices there is an excess supply—suppliers wish to sell more than buyers wish to buy. Conversely, at lower prices there is an excess demand. Only at the equilibrium price is the quantity supplied equal to the quantity demanded.

Excess supply exists when the quantity supplied exceeds the quantity demanded at the going price.

Excess demand exists when the quantity demanded exceeds the quantity supplied at the going price.

Does the market automatically reach equilibrium? To answer this question, suppose initially that the sellers choose a price of \$10. Here suppliers would like to supply 18k cubic feet, but there are no buyers—a situation of extreme excess supply. At the price of \$7 the excess supply is reduced to 9k, because both the quantity demanded is now higher at 3k units, and the quantity supplied is lower at 12k. But excess supply means that there are suppliers willing to supply at a lower price, and this willingness exerts continual downward pressure on any price above the price that equates demand and supply.

At prices below the equilibrium there is, conversely, an excess demand. In this situation, suppliers could force the price upward, knowing that buyers will continue to buy at a price at which the suppliers are willing to sell. Such upward pressure would continue until the excess demand is eliminated.

In general then, above the equilibrium price excess supply exerts downward pressure on price, and below the equilibrium excess demand exerts upward pressure on price. This process implies that the buyers and sellers have information on the various elements that make up the marketplace.

We will explore later in this chapter some specific circumstances in which trading could take place at prices above or below the equilibrium price. In such situations the quantity actually traded always corresponds to the short side of the market: this means that at high prices the quantity demanded is less than the quantity supplied, and it is the quantity demanded that is traded because buyers will not buy the amount suppliers would like to supply. Correspondingly, at low prices the quantity demanded exceeds quantity supplied, and it is the amount that suppliers are willing to sell that is traded. In sum, when trading takes place at prices

other than the equilibrium price it is always the lesser of the quantity demanded or supplied that is traded. Hence we say that at non-equilibrium prices the short side dominates. We will return to this in a series of examples later in this chapter.

The **short side of the market** determines outcomes at prices other than the equilibrium.

Supply and the nature of costs

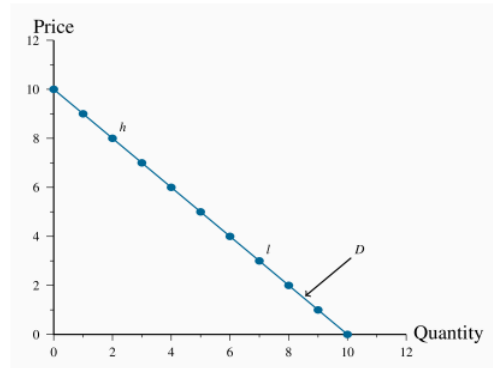
Before progressing to a graphical analysis, we should add a word about costs. The supply schedules are based primarily on the cost of producing the product in question, and we frequently assume that all of the costs associated with supply are incorporated in the supply schedules. In Chapter 6 we will explore cases where costs additional to those incurred by producers may be relevant. For example, coal burning power plants emit pollutants into the atmosphere; but the individual supplier may not take account of these pollutants, which are costs to society at large, in deciding how much to supply at different prices. Stated another way, the private costs of production would not reflect the total, or full social costs of production. Conversely, if some individuals immunize themselves against a rampant virus, other individuals gain from that action because they become less likely to contract the virus - the social value thus exceeds the private value. For the moment the assumption is that no such additional costs are associated with the markets we analyze.

This page titled [3.2: The market's building blocks](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.3: Demand and supply curves

The demand curve is a graphical expression of the relationship between price and quantity demanded, holding other things constant. Figure 3.1 measures price on the vertical axis and quantity on the horizontal axis. The curve D represents the data from the first two columns of Table 3.1. Each combination of price and quantity demanded lies on the curve. In this case the curve is *linear*—it is a straight line. The demand curve slopes downward (technically we say that its slope is negative), reflecting the fact that buyers wish to purchase more when the price is less.

Figure 3.1 Measuring price & quantity



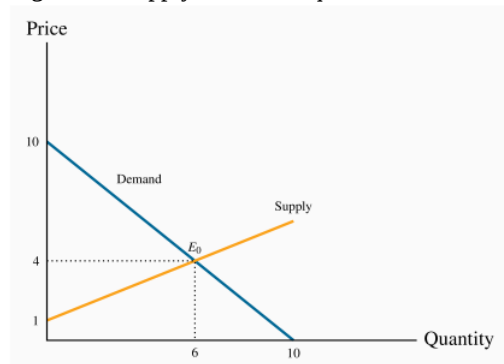
To derive this demand curve we take each price-quantity combination from the demand schedule in Table 3.1 and insert a point that corresponds to those combinations. For example, point h defines the combination $\{P = \$8, Q_d = 2\}$, the point l denotes the combination $\{P = \$3, Q_d = 7\}$. If we join all such points we obtain the demand curve in Figure 3.2. The same process yields the supply curve in Figure 3.2. In this example the supply and the demand curves are each linear. There is no reason why this linear property characterizes demand and supply curves in the real world; they are frequently found to have curvature. But straight lines are easier to work with, so we continue with them for the moment.

The **demand curve** is a graphical expression of the relationship between price and quantity demanded, with other influences remaining unchanged.

The supply curve is a graphical representation of the relationship between price and quantity supplied, holding other things constant. The supply curve S in Figure 3.2 is based on the data from columns 1 and 3 in Table 3.1. It has a positive slope indicating that suppliers wish to supply more at higher prices.

The **supply curve** is a graphical expression of the relationship between price and quantity supplied, with other influences remaining unchanged.

Figure 3.2 Supply, demand, equilibrium



The demand and supply curves intersect at point E_0 , corresponding to a price of \$4 which, as illustrated above, is the equilibrium price for this market. At any price below this the horizontal distance between the supply and demand curves represents excess demand, because demand exceeds supply. Conversely, at any price above \$4 there is an excess supply that is again measured by the horizontal distance between the two curves. Market forces tend to eliminate excess demand and excess supply as we explained

above. In the final section of the chapter we illustrate how the supply and demand curves can be 'solved' for the equilibrium price and quantity.

This page titled [3.3: Demand and supply curves](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.4: Non-price influences on demand

We have emphasized several times the importance of the *ceteris paribus* assumption when exploring the impact of different prices on the quantity demanded: We assume all other influences on the purchase decision are unchanged (at least momentarily). These other influences fall into several broad categories: The prices of related goods; the incomes of buyers; buyer tastes; and expectations about the future. Before proceeding, note that we are dealing with *market* demand rather than demand by one *individual* (the precise relationship between the two is developed later in this chapter).

The prices of related goods – oil and gas, Kindle and paperbacks

We expect that the price of other forms of energy would impact the price of natural gas. For example, if hydro-electricity, oil or solar becomes less expensive we would expect some buyers to switch to these other products. Alternatively, if gas-burning furnaces experience a technological breakthrough that makes them more efficient and cheaper we would expect some users of other fuels to move to gas. Among these examples, oil and electricity are substitute fuels for gas; in contrast a more fuel-efficient new gas furnace complements the use of gas. We use these terms, substitutes and complements, to describe products that influence the demand for the primary good.

Substitute goods: when a price reduction (rise) for a related product reduces (increases) the demand for a primary product, it is a substitute for the primary product.

Complementary goods: when a price reduction (rise) for a related product increases (reduces) the demand for a primary product, it is a complement for the primary product.

Clearly electricity is a substitute for gas in the power market, whereas a gas furnace is a complement for gas as a fuel. The words substitutes and complements immediately suggest the nature of the relationships. Every product has complements and substitutes. As another example: Electronic readers and tablets are substitutes for paper-form books; a rise in the price of paper books should increase the demand for electronic readers at any given price for electronic readers. In graphical terms, the demand curve *shifts* in response to changes in the prices of other goods – an increase in the price of paper-form books shifts the demand for electronic readers outward, because more electronic readers will be demanded at any price.

Buyer incomes – which goods to buy

The demand for most goods increases in response to income growth. Given this, the demand curve for gas will shift outward if household incomes in the economy increase. Household incomes may increase either because there are more households in the economy or because the incomes of the existing households grow.

Most goods are demanded in greater quantity in response to higher incomes at any given price. But there are exceptions. For example, public transit demand may decline at any price when household incomes rise, because some individuals move to cars. Or the demand for laundromats may decline in response to higher incomes, as households purchase more of their own consumer durables – washers and driers. We use the term inferior good to define these cases: An inferior good is one whose demand declines in response to increasing incomes, whereas a normal good experiences an increase in demand in response to rising incomes.

An **inferior good** is one whose demand falls in response to higher incomes.

A **normal good** is one whose demand increases in response to higher incomes.

There is a further sense in which consumer incomes influence demand, and this relates to how the incomes are *distributed* in the economy. In the discussion above we stated that higher total incomes shift demand curves outwards when goods are normal. But think of the difference in the demand for electronic readers between Portugal and Saudi Arabia. These economies have roughly the same average per-person income, but incomes are distributed more unequally in Saudi Arabia. It does not have a large middle class that can afford electronic readers or iPads, despite the huge wealth held by the elite. In contrast, Portugal has a relatively larger middle class that can afford such goods. Consequently, the *distribution of income* can be an important determinant of the demand for many commodities and services.

Tastes and networks – hemlines, lapels and homogeneity

While demand functions are drawn on the assumption that tastes are constant, in an evolving world they are not. We are all subject to peer pressure, the fashion industry, marketing, and a desire to maintain our image. If the fashion industry dictates that lapels on men's suits or long skirts are *de rigueur* for the coming season, some fashion-conscious individuals will discard a large segment of

their wardrobe, even though the clothes may be in perfectly good condition: Their demand is influenced by the dictates of current fashion.

Correspondingly, the items that other individuals buy or use frequently determine our own purchases. Businesses frequently decide that all of their employees will have the same type of computer and software on account of *network economies*: It is easier to communicate if equipment is compatible, and it is less costly to maintain infrastructure where the variety is less.

Expectations – betting on the future

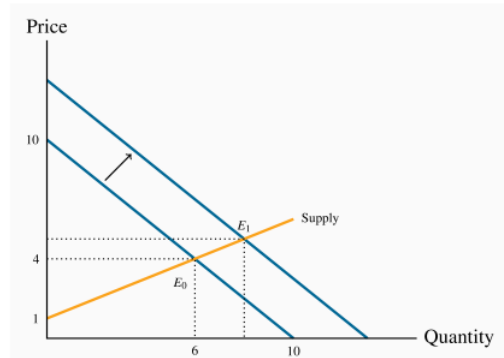
In our natural gas example, if households expected that the price of natural gas was going to stay relatively low for many years – perhaps on account of the discovery of large deposits – then they would be tempted to purchase a gas burning furnace rather than one based upon an alternative fuel. In this example, it is more than the current price that determines choices; *the prices that are expected to prevail in the future* also determine current demand.

Expectations are particularly important in stock markets. When investors anticipate that corporations will earn high rewards in the future they will buy a stock today. If enough people believe this, the price of the stock will be driven upward on the market, even before profitable earnings are registered.

Shifts in demand

The demand curve in Figure 3.2 is drawn for a given level of other prices, incomes, tastes, and expectations. Movements along the demand curve reflect solely the impact of different prices for the good in question, holding other influences constant. But changes in any of these other factors will change the position of the demand curve. Figure 3.3 illustrates a shift in the demand curve. This shift could result from a rise in household incomes that increase the quantity demanded *at every price*. This is illustrated by an outward shift in the demand curve. With supply conditions unchanged, there is a new equilibrium at E_1 , indicating a greater quantity of purchases accompanied by a higher price. The new equilibrium reflects a *change in quantity supplied and a change in demand*.

Figure 3.3 Demand shift and new equilibrium



The outward shift in demand leads to a new equilibrium E_1 .

We may well ask why so much emphasis in our diagrams and analysis is placed on the relationship between *price* and quantity, rather than on the relationship between quantity and its other determinants. The answer is that we could indeed draw diagrams with quantity on the horizontal axis and a measure of one of these other influences on the vertical axis. But the price mechanism plays a very important role. *Variations in price are what equilibrate the market*. By focusing primarily upon the price, we see the self-correcting mechanism by which the market reacts to excess supply or excess demand.

In addition, this analysis illustrates the method of comparative statics—examining the impact of changing one of the other things that are assumed constant in the supply and demand diagrams.

Comparative static analysis compares an initial equilibrium with a new equilibrium, where the difference is due to a change in one of the other things that lie behind the demand curve or the supply curve.

'Comparative' obviously denotes the idea of a comparison, and static means that we are not in a state of motion. Hence we use these words in conjunction to indicate that we compare one outcome with another, without being concerned too much about the transition from an initial equilibrium to a new equilibrium. The transition would be concerned with dynamics rather than statics. In Figure 3.3 we explain the difference between the points E_0 and E_1 by indicating that there has been a change in incomes or in the price of a substitute good. We do not attempt to analyze the details of this move or the exact path from E_0 to E_1 .

Application Box 3.1 Corn prices and demand shifts

In the middle of its second mandate, the Bush Administration in the US decided to encourage the production of ethanol – a fuel that is less polluting than gasoline. The target production was 35 billion for 2017 – from a base of 1 billion gallons in 2000. Corn is the principal input in ethanol production. It is also used as animal feed, as a sweetener and as a food for humans. The target was to be met with the help of a subsidy to producers and a tariff on imports of Brazil's sugar-cane based ethanol.

The impact on corn prices was immediate; from a farm-gate price of \$2 per bushel in 2005, the price reached the \$4 range two years later. In 2012 the price rose temporarily to \$7. While other factors were in play - growing incomes and possibly speculation by commodity investors, ethanol is seen as the main price driver: demand for corn increased and the supply could not be increased to keep up with the demand without an increase in price.

The wider impact of these developments was that the prices of virtually all grains increased in tandem with corn: the prices of sorghum and barley increased because of a switch in land use towards corn on account of its profitability.

While farmers benefited from the price rise, consumers – particularly those in less developed economies – experienced a dramatic increase in their basic living costs. Visit the site of the United Nations' Food and Agricultural Organization for an assessment. Since hitting \$7 per bushel in 2012, the price has dropped and averaged \$3.50 in 2016.

In terms of supply and demand shifts: the demand side has dominated, particularly in the short run. The ethanol drive, combined with secular growth in the demand for food, means that the demand for grains shifted outward faster than the supply. In the period 2013–2016, supply has increased and the price has moderated.

This page titled [3.4: Non-price influences on demand](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.5: Non-price influences on supply

To date we have drawn supply curves with an upward slope. Is this a reasonable representation of supply in view of what is frequently observed in markets? We suggested earlier that the various producers of a particular good or service may have different levels of efficiency. If so, only the more efficient producers can make a profit at a low price, whereas at higher prices more producers or suppliers enter the market – producers who may not be as lean and efficient as those who can survive in a lower-price environment. This view of the world yields a positively-sloping supply curve.

As a second example, consider *Uber* or *Lyft* taxi drivers. Some drivers may be in serious need of income and may be willing to drive for a low hourly rate. For other individuals driving may be a secondary source of income, and such drivers are less likely to want to drive unless the hourly wage is higher. Consequently if these ride sharing services need a large number of drivers at any one time it may be necessary to pay a higher wage – *and charge a higher fare to passengers*, to induce more drivers to take their taxis onto the road. This phenomenon corresponds to a positively-sloped supply curve.

In contrast to these two examples, some suppliers simply choose a unique price and let buyers purchase as much as they want at that price. This is the practice of most retailers. For example, the price of *Samsung's Galaxy* is typically fixed, no matter how many are purchased – and tens of millions are sold at a fixed price when a new model is launched. *Apple* also sets a price, and buyers purchase as many as they desire at that price. This practice corresponds to a horizontal supply curve: The price does not vary and the market equilibrium occurs where the demand curve intersects this supply curve.

In yet other situations supply is fixed. This happens in auctions. Bidders at the auction simply determine the price to be paid. At a real estate auction a given property is put on the market and the price is determined by the bidding process. In this case the supply of a single property is represented by a vertical supply at a quantity of 1 unit.

Regardless of the type of market we encounter, however, it is safe to assume that supply curves rarely slope downward. So, for the moment, we adopt the stance that supply curves are generally upward sloping – somewhere between the extremes of being vertical or horizontal – as we have drawn them to this point.

Next, we examine those other influences that underlie supply curves. Technology, input costs, the prices of competing goods, expectations and the number of suppliers are the most important.

Technology – computers and fracking

A technological advance may involve an idea that allows more output to be produced with the same inputs, or an equal output with fewer inputs. A good example is *just-in-time* technology. Before the modern era, virtually all manufacturers kept large stocks of components in their production facilities, but developments in communications and computers at that time made it possible for manufacturers to link directly with their input suppliers. Nowadays auto assembly plants place their order for, say, seat delivery to their local seat supplier well ahead of assembly time. The seats swing into the assembly area hours or minutes before assembly—just in time. The result is that the assembler reduces her seat inventory (an input) and thereby reduces production cost.

Such a technology-induced cost saving is represented by moving the supply curve downward or outward: The supplier is now able and willing to supply the same quantity at a lower price because of the technological innovation. Or, saying the same thing slightly differently, suppliers will supply more at a given price than before.

A second example relates to the extraction of natural gas. The development of 'fracking' means that companies involved in gas recovery can now do so at a lower cost. Hence they are willing to supply any given quantity at a lower price. A third example concerns aluminum cans. Today they weigh a fraction of what they weighed 20 years ago. This is a technology-based cost saving.

Input costs

Input costs can vary independently of technology. For example, a wage negotiation that grants workers a substantial pay raise will increase the cost of production. This is reflected in a leftward, or upward, supply shift: Any quantity supplied is now priced higher; alternatively, suppliers are willing to supply less at the going price.

Production costs may increase as a result of higher required standards in production. As governments implement new safety or product-stress standards, costs may increase. In this instance the increase in costs is not a 'bad' outcome for the buyer. She may be purchasing a higher quality good as a result.

Competing products – Airbnb versus hotels

If competing products improve in quality or fall in price, a supplier may be forced to follow suit. For example, *Asus* and *Dell* are constantly watching each other's pricing policies. If *Dell* brings out a new generation of computers at a lower price, *Asus* may lower its prices in turn—which is to say that *Asus*' supply curve will shift downward. Likewise, *Samsung* and *Apple* each responds to the other's pricing and technology behaviours. The arrival of new products in the marketplace also impacts the willingness of suppliers to supply goods at a given price. New intermediaries such as *Airbnb* and *Vacation Rentals by Owner* have shifted the supply curves of hotel rooms downward.

These are some of the many factors that influence the position of the supply curve in a given market.

Application Box 3.2 The price of light

Technological developments have had a staggering impact on many price declines. Professor William Nordhaus of Yale University is an expert on measuring technological change. He has examined the trend in the real price of lighting. Originally, light was provided by whale oil and gas lamps and these sources of lumens (the scientific measure of the amount of light produced) were costly. In his research, Professor Nordhaus pieced together evidence on the actual historic cost of light produced at various times, going all the way back to 1800. He found that light in 1800 cost about 100 times more than in 1900, and light in the year 2000 was a fraction of its cost in 1900. A rough calculation suggests that light was five hundred times more expensive at the start of this 200-year period than at the end, and this was before the arrival of LEDs.

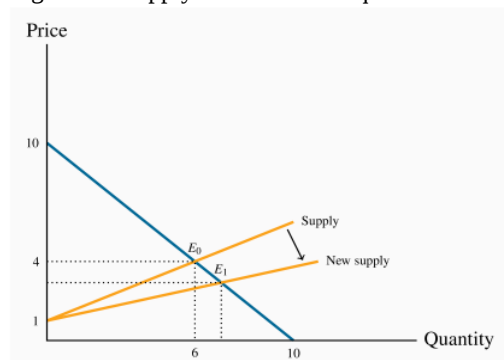
In terms of supply and demand analysis, light has been subject to very substantial downward supply shifts. Despite the long-term growth in demand, the technologically-induced supply changes have been the dominant factor in its price determination.

For further information, visit Professor Nordhaus's website in the Department of Economics at Yale University.

Shifts in supply

Whenever technology changes, or the costs of production change, or the prices of competing products adjust, then one of our *ceteris paribus* assumptions is violated. Such changes are generally reflected by shifting the supply curve. Figure 3.4 illustrates the impact of the arrival of just-in-time technology. The supply curve shifts, reflecting the ability of suppliers to supply the same output at a reduced price. The resulting new equilibrium price is lower, since production costs have fallen. At this reduced price more gas is traded at a lower price.

Figure 3.4 Supply shift and new equilibrium



The supply curve shifts due to lower production costs. A new equilibrium E_1 is attained in the market at a lower price.

This page titled [3.5: Non-price influences on supply](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

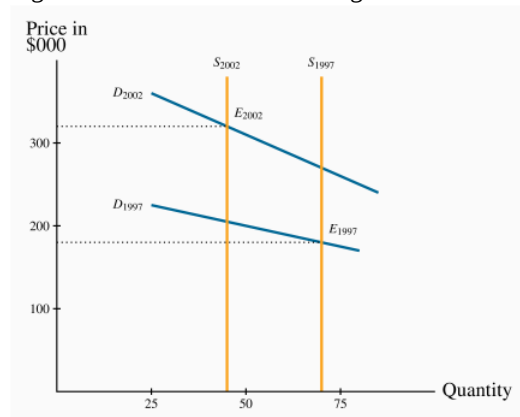
3.6: Simultaneous supply and demand impacts

In the real world, demand and supply frequently shift at the same time. We present such a case in Figure 3.5. It is based upon real estate data describing the housing market in a small Montreal municipality. Vertical curves define the supply side of the market. Such vertical curves mean that a given number of homeowners decide to put their homes on the market, and these suppliers just take whatever price results in the market. In this example, fewer houses were offered for sale in 2002 (less than 50) than in 1997 (more than 70). We are assuming in this market that the houses traded were similar; that is, we are not lumping together mansions with row houses.

During this time period household incomes increased substantially and, also, mortgage rates fell. Both of these developments shifted the demand curve upward/outward: Buyers were willing to pay more for housing in 2002 than in 1997, both because their incomes were on average higher and because they could borrow more cheaply.

The shifts on both sides of the market resulted in a higher average price. And each of these shifts compounded the other: The outward shift in demand would lead to a higher price on its own, and a reduction in supply would do likewise. Hence both forces acted to push up the price in 2002. If, instead, the supply had been greater in 2002 than in 1997 this would have acted to reduce the equilibrium price. And with the demand and supply shifts operating in opposing directions, it is not possible to say in general whether the price would increase or decrease. If the demand shift were strong and the supply shift weak then the demand forces would have dominated and led to a higher price. Conversely, if the supply forces were stronger than the demand forces.

Figure 3.5 A model of the housing market with shifts in demand and supply



The vertical supply denotes a fixed number of houses supplied each year. Demand was stronger in 2002 than in 1997 both on account of higher incomes and lower mortgage rates. Thus the higher price in 2002 is due to both a reduction in supply and an increase in demand.

This page titled [3.6: Simultaneous supply and demand impacts](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.7: Market interventions - governments and interest groups

The freely functioning markets that we have developed certainly do not describe all markets. For example, minimum wages characterize the labour market, most agricultural markets have supply restrictions, apartments are subject to rent controls, and blood is not a freely traded market commodity in Canada. In short, price controls and quotas characterize many markets. Price controls are government rules or laws that inhibit the formation of market-determined prices. Quotas are physical restrictions on how much output can be brought to the market.

Price controls are government rules or laws that inhibit the formation of market-determined prices.

Quotas are physical restrictions on output.

Price controls come in the form of either *floors* or *ceilings*. Price floors are frequently accompanied by *marketing boards*.

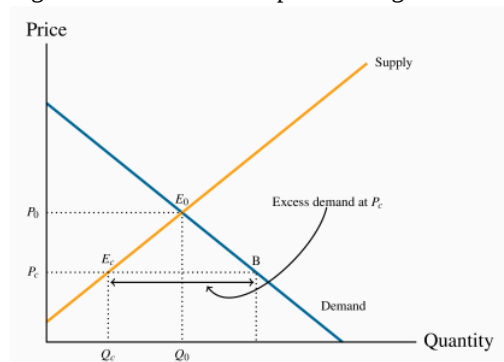
Price ceilings – rental boards

Ceilings mean that suppliers cannot legally charge more than a specific price. Limits on apartment rents are one form of ceiling. In times of emergency – such as flooding or famine, price controls are frequently imposed on foodstuffs, in conjunction with rationing, to ensure that access is not determined by who has the most income. The problem with price ceilings, however, is that they leave demand unsatisfied, and therefore they must be accompanied by some other allocation mechanism.

Consider an environment where, for some reason – perhaps a sudden and unanticipated growth in population – rents increase. Let the resulting equilibrium be defined by the point E_0 in Figure 3.6. If the government were to decide that this is an unfair price because it places hardships on low- and middle-income households, it might impose a price limit, or ceiling, of P_c . The problem with such a limit is that excess demand results: Individuals want to rent more apartments than are available in the city. In a free market the price would adjust upward to eliminate the excess demand, but in this controlled environment it cannot. So some other way of allocating the available supply between demanders must evolve.

In reality, most apartments are allocated to those households already occupying them. But what happens when such a resident household decides to purchase a home or move to another city? In a free market, the landlord could increase the rent in accordance with market pressures. But in a controlled market a city's rental tribunal may restrict the annual rent increase to just a couple of percent and the demand may continue to outstrip supply. So how does the stock of apartments get allocated between the potential renters? One allocation method is well known: The existing tenant informs her friends of her plan to move, and the friends are the first to apply to the landlord to occupy the apartment. But that still leaves much unmet demand. If this is a student rental market, students whose parents live nearby may simply return 'home'. Others may choose to move to a part of the city where rents are more affordable.

Figure 3.6 The effect of a price ceiling



The free market equilibrium occurs at E_0 . A price ceiling at P_c holds down the price but leads to excess demand E_cB , because Q_c is the quantity traded. A price ceiling above P_0 is irrelevant since the free market equilibrium E_0 can still be attained.

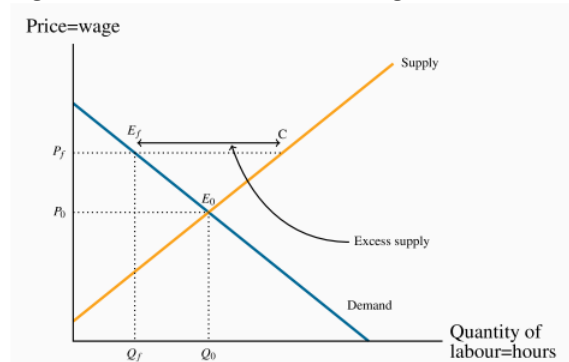
However, rent controls sometimes yield undesirable outcomes. Rent controls are widely studied in economics, and the consequences are well understood: Landlords tend not to repair or maintain their rental units in good condition if they cannot obtain the rent they believe they are entitled to. Accordingly, the residential rental stock deteriorates. In addition, builders realize that more money is to be made in building condominium units than rental units, or in *converting rental units to condominiums*. The frequent consequence is thus a *reduction* in supply and a reduced quality. Market forces are hard to circumvent because, as we

emphasized in Chapter 1, economic players react to the incentives they face. These outcomes are examples of what we call the *law of unintended consequences*.

Price floors – minimum wages

An effective price floor sets the price *above* the market-clearing price. A minimum wage is the most widespread example in the Canadian economy. Provinces each set their own minimum, and it is seen as a way of protecting the well-being of low-skill workers. Such a floor is illustrated in Figure 3.7. The free-market equilibrium is again E_0 , but the effective market outcome is the combination of price and quantity corresponding to the point E_f at the price floor, P_f . In this instance, there is excess supply equal to the amount E_fC .

Figure 3.7 Price floor – minimum wage



In a free market the equilibrium is E_0 . A minimum wage of P_f raises the hourly wage, but reduces the hours demanded to Q_f . Thus E_fC is the excess supply.

Note that there is a similarity between the outcomes defined in the floor and ceiling cases: The quantity actually traded is *the lesser of the supply quantity and demand quantity at the going price: The short side dominates*.

If price floors, in the form of minimum wages, result in some workers going unemployed, why do governments choose to put them in place? The excess supply in this case corresponds to unemployment – more individuals are willing to work for the going wage than buyers (employers) wish to employ. The answer really depends upon the magnitude of the excess supply. In particular, suppose, in Figure 3.7 that the supply and demand curves going through the equilibrium E_0 were more 'vertical'. This would result in a smaller excess supply than is represented with the existing supply and demand curves. This would mean in practice that a higher wage could go to workers, making them better off, without causing substantial unemployment. This is the trade off that governments face: With a view to increasing the purchasing power of generally lower-skill individuals, a minimum wage is set, hoping that the negative impact on employment will be small. We will return to this in the next chapter, where we examine the responsiveness of supply and demand curves to different prices.

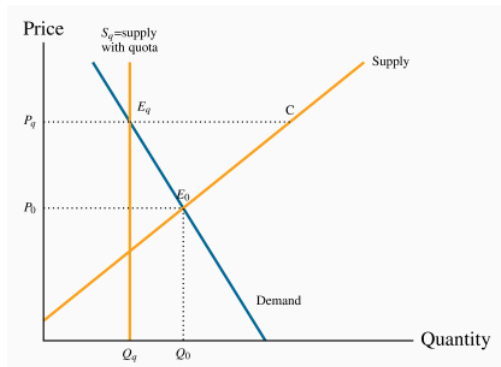
Quotas – agricultural supply

A quota represents the right to supply a specified quantity of a good to the market. It is a means of keeping prices higher than the free-market equilibrium price. As an alternative to imposing a price floor, the government can generate a high price by restricting supply.

Agricultural markets abound with examples. In these markets, farmers can supply only what they are permitted by the quota they hold, and there is usually a market for these quotas. For example, in several Canadian provinces it currently costs in the region of \$30,000 to purchase a quota granting the right to sell the milk of one cow. The cost of purchasing quotas can thus easily outstrip the cost of a farm and herd. Canadian cheese importers must pay for the right to import cheese from abroad. Restrictions also apply to poultry. The impact of all of these restrictions is to raise the domestic price above the free market price.

In Figure 3.8, the free-market equilibrium is at E_0 . In order to raise the price above P_0 , the government restricts supply to Q_q by granting quotas, which permit producers to supply a limited amount of the good in question. This supply is purchased at the price equal to P_q . From the standpoint of farmers, a higher price might be beneficial, even if they get to supply a smaller quantity, provided the amount of revenue they get as a result is as great as the revenue in the free market.

Figure 3.8 The effect of a quota



The government decides that the equilibrium price P_0 is too low. It decides to boost price by reducing supply from Q_0 to Q_q . It achieves this by requiring producers to have a production quota. This is equivalent to fixing supply at S_q .

Marketing boards – milk and maple syrup

A marketing board is a means of insuring that a quota or price floor can be maintained. Quotas are frequent in the agriculture sector of the economy. One example is maple syrup in Quebec. The Federation of Maple Syrup Producers of Quebec has the sole right to market maple syrup. All producers must sell their syrup through this marketing board. The board thus has a particular type of power in the market: it has control of the market at the wholesale end, because it is a sole buyer. The Federation increases the total revenue going to producers by artificially restricting the supply to the market. The Federation calculates that by reducing supply and selling it at a higher price, more revenue will accrue to the producers. This is illustrated in Figure 3.8. The market equilibrium is given by E_0 , but the Federation restricts supply to the quantity Q_q , which is sold to buyers at price P_q . To make this possible the total supply must be restricted; otherwise producers would supply the amount given by the point C on the supply curve, and this would result in excess supply in the amount E_qC . In order to restrict supply to Q_q in total, individual producers are limited in what they can sell to the Federation; they have a quota, which gives them the right to produce and sell no more than a specified amount. This system of quotas is necessary to eliminate the excess supply that would emerge at the above-equilibrium price P_q .

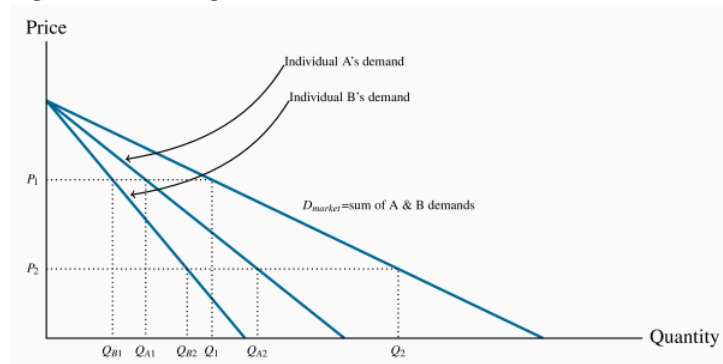
We will return to this topic in Chapter 4. For the moment, to see that this type of revenue-increasing outcome is possible, examine Table 3.1 again. At this equilibrium price of \$4 the quantity traded is 6 units, yielding a total expenditure by buyers (revenue to suppliers) of \$24. However, if the supply were restricted and a price of \$5 were set, the expenditure by buyers (revenue to suppliers) would rise to \$25.

This page titled [3.7: Market interventions - governments and interest groups](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.8: Individual and market functions

Markets are made up of many individual participants on the demand and supply side. The supply and demand functions that we have worked with in this chapter are those for the total of all participants on each side of the market. But how do we arrive at such market functions when the economy is composed of individuals? We can illustrate how, with the help of Figure 3.9.

Figure 3.9 Summing individual demands



At P_1 individual A purchases Q_{A1} and B purchases Q_{B1} . The total demand is the sum of these individual demands at this price (Q_1). At P_2 individual demands are summed to Q_2 . Since the points Q_1 and Q_2 define the demands of the market participants it follows that market demand is the horizontal sum of these curves.

To concentrate on the essentials, imagine that there are just two buyers of chocolate cookies in the economy. A has a stronger preference for cookies than B, so his demand is greater. To simplify, let the two demands have the same intercept on the vertical axis. The curves D_A and D_B indicate how many cookies A and B, respectively, will buy at each price. The market demand indicates how much they buy *together* at any price. Accordingly, at P_1 , A and B purchase the quantities Q_{A1} and Q_{B1} respectively. Thus $Q_1 = Q_{A1} + Q_{B1}$. At a price P_2 , they purchase Q_{A2} and Q_{B2} . Thus $Q_2 = Q_{A2} + Q_{B2}$. The market demand is therefore the horizontal sum of the individual demands at these prices. In the figure this is defined by D_{market} .

Market demand: the horizontal sum of individual demands.

This page titled [3.8: Individual and market functions](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.9: Useful techniques - demand and supply equations

The supply and demand functions, or equations, underlying Table 3.1 and Figure 3.2 can be written in their mathematical form:

Demand: $P = 10 - Q$

Supply: $P = 1 + (1/2)Q$

A straight line is represented completely by the intercept and slope. In particular, if the variable P is on the vertical axis and Q on the horizontal axis, the straight-line equation relating P and Q is defined by $P = a + bQ$. Where the line is negatively sloped, as in the demand equation, the parameter b must take a negative value. By observing either the data in Table 3.1 or Figure 3.2 it is clear that the vertical intercept, a , takes a value of \$10. The vertical intercept corresponds to a zero-value for the Q variable. Next we can see from Figure 3.2 that the slope (given by the rise over the run) is 10/10 and hence has a value of -1 . Accordingly the demand equation takes the form $P = 10 - Q$.

On the supply side the price-axis intercept, from either the figure or the table, is clearly 1. The slope is one half, because a two-unit change in quantity is associated with a one-unit change in price. This is a positive relationship obviously so the supply curve can be written as $P = 1 + (1/2)Q$.

Where the supply and demand curves intersect is the market equilibrium; that is, the price-quantity combination is the same for both supply and demand where the supply curve takes on the same values as the demand curve. This unique price-quantity combination is obtained by equating the two curves: If Demand = Supply, then

$$10 - Q = 1 + (1/2)Q.$$

Gathering the terms involving Q to one side and the numerical terms to the other side of the equation results in $9 = 1.5Q$. This implies that the equilibrium quantity must be 6 units. And this quantity must trade at a price of \$4. That is, when the price is \$4 both the quantity demanded and the quantity supplied take a value of 6 units.

Modelling market interventions using equations

To illustrate the impact of market interventions examined in Section 3.7 on our numerical market model for natural gas, suppose that the government imposes a minimum price of \$6 – above the equilibrium price obviously. We can easily determine the quantity supplied and demanded at such a price. Given the supply equation

$$P = 1 + (1/2)Q,$$

it follows that at $P = 6$ the quantity supplied is 10. This follows by solving the relationship $6 = 1 + (1/2)Q$ for the value of Q . Accordingly, suppliers *would like to supply* 10 units at this price.

Correspondingly on the demand side, given the demand curve

$$P = 10 - Q,$$

with a price given by $P = \$6$, it must be the case that $Q = 4$. So buyers *would like to buy* 4 units at that price: There is excess supply. But we know that the short side of the market will win out, and so the actual amount traded at this restricted price will be 4 units.

This page titled 3.9: Useful techniques - demand and supply equations is shared under a [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.10: Conclusion

We have covered a lot of ground in this chapter. It is intended to open up the vista of economics to the new student in the discipline. Economics is powerful and challenging, and the ideas we have developed here will serve as conceptual foundations for our exploration of the subject. Our next chapter deals with measurement and responsiveness.

This page titled [3.10: Conclusion](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.11: Key Terms

Demand is the quantity of a good or service that buyers wish to purchase at each possible price, with all other influences on demand remaining unchanged.

Supply is the quantity of a good or service that sellers are willing to sell at each possible price, with all other influences on supply remaining unchanged.

Quantity demanded defines the amount purchased at a particular price.

Quantity supplied refers to the amount supplied at a particular price.

Equilibrium price: equilibrates the market. It is the price at which quantity demanded equals the quantity supplied.

Excess supply exists when the quantity supplied exceeds the quantity demanded at the going price.

Excess demand exists when the quantity demanded exceeds quantity supplied at the going price.

Short side of the market determines outcomes at prices other than the equilibrium.

Demand curve is a graphical expression of the relationship between price and quantity demanded, with other influences remaining unchanged.

Supply curve is a graphical expression of the relationship between price and quantity supplied, with other influences remaining unchanged.

Substitute goods: when a price reduction (rise) for a related product reduces (increases) the demand for a primary product, it is a substitute for the primary product.

Complementary goods: when a price reduction (rise) for a related product increases (reduces) the demand for a primary product, it is a complement for the primary product.

Inferior good is one whose demand falls in response to higher incomes.

Normal good is one whose demand increases in response to higher incomes.

Comparative static analysis compares an initial equilibrium with a new equilibrium, where the difference is due to a change in one of the other things that lie behind the demand curve or the supply curve.

Price controls are government rules or laws that inhibit the formation of market-determined prices.

Quotas are physical restrictions on output.

Market demand: the horizontal sum of individual demands.

This page titled [3.11: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

3.12: Exercises for Chapter 3

EXERCISE 3.1

The supply and demand for concert tickets are given in the table below.

Price (\$)	0	4	8	12	16	20	24	28	32	36	40
Quantity demanded	15	14	13	12	11	10	9	8	7	6	5
Quantity supplied	0	0	0	0	0	1	3	5	7	9	11

- Plot the supply and demand curves to scale and establish the equilibrium price and quantity.
- What is the excess supply or demand when price is \$24? When price is \$36?
- Describe the market adjustments in price induced by these two prices.
- Optional:* The functions underlying the example in the table are linear and can be presented as $P=18+2Q$ (supply) and $P=60-4Q$ (demand). Solve the two equations for the equilibrium price and quantity values.

EXERCISE 3.2

Illustrate in a supply/demand diagram, by shifting the demand curve appropriately, the effect on the demand for flights between Calgary and Winnipeg as a result of:

- Increasing the annual government subsidy to *Via Rail*.
- Improving the Trans-Canada highway between the two cities.
- The arrival of a new budget airline on the scene.

EXERCISE 3.3

A new trend in US high schools is the widespread use of chewing tobacco. A recent survey indicates that 15 percent of males in upper grades now use it – a figure not far below the use rate for cigarettes. This development came about in response to the widespread implementation by schools of regulations that forbade cigarette smoking on and around school property. Draw a supply-demand equilibrium for each of the cigarette and chewing tobacco markets before and after the introduction of the regulations.

EXERCISE 3.4

The following table describes the demand and supply conditions for labour.

Price (\$) = wage rate	0	10	20	30	40	50	60	70	80	90	100	110	120	130	140	150	160	170
Quantity demanded	1020	960	900	840	780	720	660	600	540	480	420	360	300	240	180	120	60	0
Quantity supplied	0	0	0	0	0	0	30	60	90	120	150	180	210	240	270	300	330	360

- Graph the functions and find the equilibrium price and quantity by equating demand and supply.

- b. Suppose a price ceiling is established by the government at a price of \$120. This price is below the equilibrium price that you have obtained in part (a). Calculate the amount that would be demanded and supplied and then calculate the excess demand.

EXERCISE 3.5

In Exercise 3.4, suppose that the supply and demand describe an agricultural market rather than a labour market, and the government implements a price floor of \$140. This is greater than the equilibrium price.

- Estimate the quantity supplied and the quantity demanded at this price, and calculate the excess supply.
- Suppose the government instead chose to maintain a price of \$140 by implementing a system of quotas. What quantity of quotas should the government make available to the suppliers?

EXERCISE 3.6

In Exercise 3.5, suppose that, at the minimum price, the government buys up all of the supply that is not demanded, and exports it at a price of \$80 per unit. Compute the cost to the government of this operation.

EXERCISE 3.7

Let us sum two demand curves to obtain a 'market' demand curve. We will suppose there are just two buyers in the market. Each of the individual demand curves has a price intercept of \$42. One has a quantity intercept of 126, the other 84.

- Draw the demands either to scale or in an Excel spreadsheet, and label the intercepts on both the price and quantity axes.
- Determine how much would be purchased in the market at prices \$10, \$20, and \$30.
- Optional:* Since you know the intercepts of the market (total) demand curve, can you write an equation for it?

EXERCISE 3.8

In Exercise 3.7 the demand curves had the same price intercept. Suppose instead that the first demand curve has a price intercept of \$36 and a quantity intercept of 126; the other individual has a demand curve defined by a price intercept of \$42 and a quantity intercept of 84. Graph these curves and illustrate the market demand curve.

EXERCISE 3.9

Here is an example of a demand curve that is not linear:

Price (\$)	4	3	2	1	0
Quantity demanded	25	100	225	400	625

- Plot this demand curve to scale or in Excel.
- If the supply function in this market is $P=2$, plot this function in the same diagram.
- Determine the equilibrium quantity traded in this market.

EXERCISE 3.10

The football stadium of the University of the North West Territories has 30 seats. The demand curve for tickets has a price intercept of \$36 and a quantity intercept of 72.

- Draw the supply and demand curves to scale in a graph or in Excel. (This demand curve has the form $P = 36 - 0.5 \times Q$.)
- Determine the equilibrium admission price, and the amount of revenue generated from ticket sales for each game.
- A local alumnus and benefactor offers to install 6 more seats at no cost to the University. Compute the price that would be charged with this new supply and compute the revenue that would accrue at this new equilibrium price. Should the University accept the offer to install the seats?
- Redo the previous part of this question, assuming that the initial number of seats is 40, and the University has the option to increase capacity to 46 at no cost to itself. Should the University accept the offer in this case?

EXERCISE 3.11

Suppose farm workers in Mexico are successful in obtaining a substantial wage increase. Illustrate the effect of this on the price of lettuce in the Canadian winter, using a supply and demand diagram, on the assumption that all lettuce in Canada is imported during

its winter.

This page titled [3.12: Exercises for Chapter 3](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

SECTION OVERVIEW

Unit 2: Responsiveness and the Value of Markets

The degree to which individuals or firms, or any economic agent, respond to incentives is important to ascertain for pricing and policy purposes: If prices change, to what degree will suppliers and buyers respond? How will markets respond to taxes? Chapter 4 explores and develops the concept of elasticity, which is the word economists use to define responsiveness. A meaningful metric, one formulated in percentage terms, that is applicable to virtually any market or incentive means that behaviours can be compared in different environments.

In Chapter 5 we explore how markets allocate resources and how the well-being of society's members is impacted by uncontrolled and controlled markets. A central theme of this chapter is that markets are very useful environments, but, if they are to serve the social interest, need to be controlled in many circumstances.

4: Measures of response- Elasticities

- 4.1: Price responsiveness of demand
- 4.2: Price elasticities and public policy
- 4.3: The time horizon and inflation
- 4.4: Cross-price elasticities - cable or satellite
- 4.5: The income elasticity of demand
- 4.6: Elasticity of supply
- 4.7: Elasticities and tax incidence
- 4.8: Technical tricks with elasticities
- 4.9: Key Terms
- 4.10: Exercises for Chapter 4

5: Welfare economics and externalities

- 5.1: Equity and efficiency
- 5.2: Consumer and producer surplus
- 5.3: Efficient market outcomes
- 5.4: Taxation, surplus and efficiency
- 5.5: Market failures - externalities
- 5.6: Other market failures
- 5.7: Environmental policy and climate change
- 5.8: Conclusion
- 5.9: Key Terms
- 5.10: Exercises for Chapter 5

This page titled [Unit 2: Responsiveness and the Value of Markets](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4: Measures of response- Elasticities

Chapter 4: Measures of response: Elasticities

In this chapter we will explore:

4.1	Responsiveness as elasticities
4.2	Demand elasticities and public policy
4.3	The time horizon and inflation
4.4	Cross-price elasticities
4.5	Income elasticity of demand
4.6	Supply side responses
4.7	Tax incidence
4.8	Technical tricks with elasticities

4.1 Price responsiveness of demand

Put yourself in the position of an entrepreneur. One of your many challenges is to price your product appropriately. You may be Michael Dell choosing a price for your latest computer, or the local restaurant owner pricing your table d'hôte, or you may be pricing your part-time snow-shoveling service. A key component of the pricing decision is to know how *responsive* your market is to variations in your pricing. How we measure responsiveness is the subject matter of this chapter.

We begin by analyzing the responsiveness of consumers to price changes. For example, consumers tend not to buy much more or much less food in response to changes in the general price level of food. This is because food is a pretty basic item for our existence. In contrast, if the price of new textbooks becomes higher, students may decide to search for a second-hand copy, or make do with lecture notes from their friends or downloads from the course web site. In the latter case students have ready alternatives to the new text book, and so their expenditure patterns can be expected to reflect these options, whereas it is hard to find alternatives to food. In the case of food consumers are not very responsive to price changes; in the case of textbooks they are. The word 'elasticity' that appears in this chapter title is just another term for this concept of responsiveness. Elasticity has many different uses and interpretations, and indeed more than one way of being measured in any given situation. Let us start by developing a suitable numerical measure.

The slope of the demand curve suggests itself as one measure of responsiveness: If we lowered the price of a good by \$1, for example, how many more units would we sell? The difficulty with this measure is that it does not serve us well when comparing different products. One dollar may be a substantial part of the price of your morning coffee and croissant, but not very important if buying a computer or tablet. Accordingly, when goods and services are measured in different units (croissants versus tablets), or when their prices are very different, it is often best to use a *percentage* change measure, which is *unit-free*.

The price elasticity of demand is measured as the percentage change in quantity demanded, divided by the percentage change in price. Although we introduce several other elasticity measures later, when economists speak of the demand elasticity they invariably mean the price elasticity of demand defined in this way.

The **price elasticity of demand** is measured as the percentage change in quantity demanded, divided by the percentage change in price.

The price elasticity of demand can be written in different forms. We will use the Greek letter epsilon, ϵ , as a shorthand symbol, with a subscript d to denote demand, and the capital delta, Δ , to denote a change. Therefore, we can write

$$\text{Price elasticity of demand} = \epsilon_d = \frac{\text{Percentage change in quantity demanded}}{\text{Percentage change in price}}$$

or, using a shortened expression,

$\epsilon_d = \frac{\% \Delta Q}{\% \Delta P}$	(4.1)	
--	-------	--

Calculating the value of the elasticity is not difficult. If we are told that a 10 percent price increase reduces the quantity demanded by 20 percent, then the elasticity value is $-20\%/10\% = -2$. The negative sign denotes that price and quantity move in opposite directions, but for brevity the negative sign is often omitted.

Consider now the data in Table 4.1 and the accompanying Figure 4.1. These data reflect the demand relation for natural gas that we introduced in Chapter 3. Note first that, when the price and quantity change, we must decide what *reference price and quantity* to use in the percentage change calculation in the definition above. We could use the initial or final price-quantity combination, or an average of the two. Each choice will yield a slightly different numerical value for the elasticity. The best convention is to *use the midpoint of the price values and the corresponding midpoint of the quantity values*. This ensures that the elasticity value is the same regardless of whether we start at the higher price or the lower price. Using the subscript 1 to denote the initial value and 2 the final value:

$$\text{Average quantity } \bar{Q} = (Q_1 + Q_2)/2$$

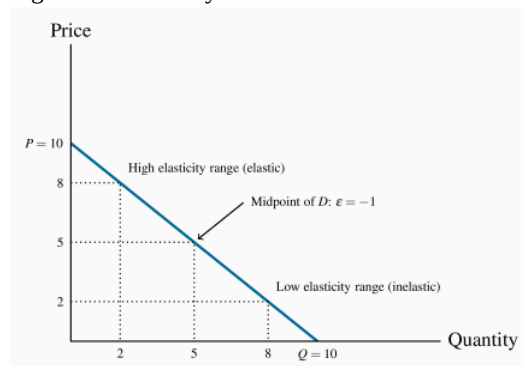
$$\text{Average price } \bar{P} = (P_1 + P_2)/2$$

Table 4.1 The demand for natural gas: Elasticities and revenue

Price (\$)	Quantity demanded	Elasticity value	Total revenue (\$)
10	0		0
9	1	-9.0	9
8	2		16
7	3	-2.33	21
6	4		24
5	5	-1.0	25
4	6		24
3	7	-0.43	21
2	8		16
1	9	-0.11	9
0	10		0

Elasticity calculations are based upon \$2 price changes.

Figure 4.1 Elasticity variation with linear demand



In the high-price region of the demand curve the elasticity takes on a high value. At the midpoint of a linear demand curve the elasticity takes on a value of one, and at lower prices the elasticity value continues to fall.

Using this rule, consider now the value of ϵ_d when price drops from \$10.00 to \$8.00. The change in price is \$2.00 and the average price is therefore \$9.00 $[= (\$10.00 + \$8.00)/2]$. On the quantity side, demand goes from zero to 2 units (measured in thousands of cubic feet), and the average quantity demanded is therefore $(0+2)/2=1$. Putting these numbers into the formula yields:

$$\epsilon_d = \frac{(Q_2 - Q_1)/\bar{Q}}{(P_2 - P_1)/\bar{P}} = \frac{(2/1)}{-(2/9)} = -\left(\frac{2}{1}\right) \times \left(\frac{9}{2}\right) = -9.$$

Note that the price has declined in this instance and thus the change in price is negative. Continuing down the table in this fashion yields the full set of elasticity values in the third column.

The demand elasticity is said to be *high* if it is a large negative number; the large number denotes a high degree of sensitivity. Conversely, the elasticity is *low* if it is a small negative number. High and low refer to the size of the number, ignoring the negative sign. The term *arc elasticity* is also used to define what we have just measured, indicating that it defines consumer responsiveness over a segment or *arc* of the demand curve.

It is helpful to analyze this numerical example by means of the corresponding demand curve that is plotted in Figure 4.1, and which we used in Chapter 3. It is a straight-line demand curve; but, despite this, the elasticity is not constant. At high prices the elasticity is high; at low prices it is low. The intuition behind this pattern is as follows. When the price is high, a given price change represents a small *percentage* change, because the average price in the price-term denominator is large. At high prices the quantity demanded is small and therefore the percentage quantity change tends to be large due to the small quantity value in its denominator. In sum, at high prices the elasticity is large; it contains a large numerator and a small denominator. By the same reasoning, at low prices the elasticity is small.

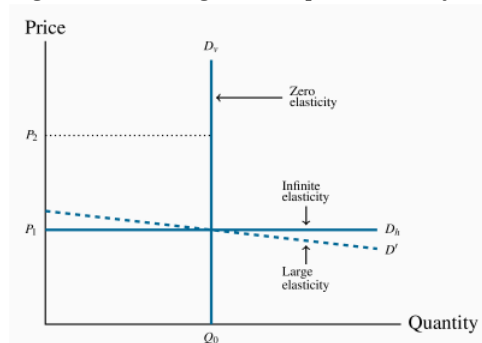
We can carry this reasoning one step further to see what happens when the demand curve intersects the axes. At the horizontal axis the average price is tending towards zero. Since this extremely small value appears in the denominator of the price term it means that the price term as a whole is extremely large. Accordingly, with an extremely large value in the denominator of the elasticity expression, the whole ratio is tending towards a zero value. By the same reasoning the elasticity value at the vertical intercept is tending towards an infinitely large value.

Extreme cases

The elasticity decreases in going from high prices to low prices. This is true for most non-linear demand curves also. Two exceptions are when the demand curve is horizontal and when it is vertical.

When the demand curve is vertical, no quantity change results from a change in price from P_1 to P_2 , as illustrated in Figure 4.2 using the demand curve D_v . Therefore, the numerator in Equation 4.1 is zero, and the elasticity has a zero value.

Figure 4.2 Limiting cases of price elasticity



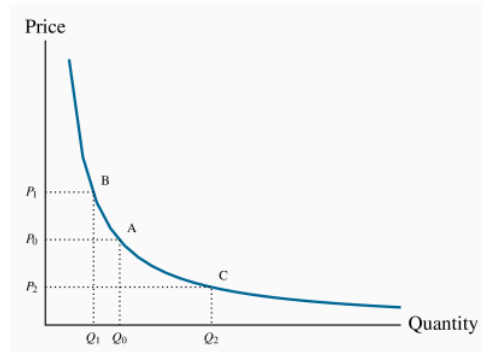
When the demand curve is vertical (D_v), the elasticity is zero: A change in price from P_1 to P_2 has no impact on the quantity demanded because the numerator in the elasticity formula has a zero value. When D becomes more horizontal the elasticity becomes larger and larger at P_1 , eventually becoming infinite.

In the horizontal case, we say that the elasticity is *infinite*, which means that any percentage price change brings forth an infinite quantity change! This case is also illustrated in Figure 4.2 using the demand curve D_h . As with the vertical demand curve, this is not immediately obvious. So consider a demand curve that is almost horizontal, such as D' instead of D_h . In this instance, we can achieve large changes in quantity demanded by implementing very small price changes. In terms of Equation 4.1, the numerator is large and the denominator small, giving rise to a large elasticity. Now imagine that this demand curve becomes ever more elastic (horizontal). The same quantity response can be obtained with a smaller price change, and hence the elasticity is larger. Pursuing this idea, we can say that, *as the demand curve becomes ever more elastic, the elasticity value tends towards infinity*.

A *non-linear demand curve* is illustrated in Figure 4.3. If price increases from P_0 to P_1 , the corresponding quantity change is given by $(Q_0 - Q_1)$. When the price declines to P_2 the quantity increases from Q_0 to Q_2 . When statisticians study data to determine how

responsive purchases are to price changes they do not always find a linear relationship between price and quantity. But a linear relationship is frequently a good approximation or representation of actual data and we will continue to analyze responsiveness in a linear framework in this chapter.

Figure 4.3 Non-linear demand curves



When the demand curve is non-linear the slope changes with the price. Hence, equal price changes do not lead to equal quantity changes: The quantity change associated with a change in price from P_0 to P_1 is smaller than the change in quantity associated with the same change in price from P_0 to P_2 .

Elastic and inelastic demands

While the elasticity value falls as we move down the demand curve, an important dividing line occurs at the value of -1 . This is illustrated in Table 4.1, and is a property of all straight-line demand curves. Disregarding the negative sign, demand is said to be elastic if the price elasticity is greater than unity, and inelastic if the value lies between unity and 0. It is unit elastic if the value is exactly one.

Demand is **elastic** if the price elasticity is greater than unity. It is **inelastic** if the value lies between unity and 0. It is **unit elastic** if the value is exactly one.

Economists frequently talk of goods as having a "high" or "low" demand elasticity. What does this mean, given that the elasticity varies throughout the length of a demand curve? It signifies that, *at the price usually charged*, the elasticity has a high or low value. For example, your weekly demand for regular coffee at Starbucks might be unresponsive to variations in price around the value of \$3.00, but if the price were \$6, you might be more responsive to price variations. Likewise, when we stated at the beginning of this chapter that the demand for food tends to be inelastic, we really meant that *at the price we customarily face for food*, demand is inelastic.

Determinants of price elasticity

Why is it that the price elasticities for some goods and services are high and for others low?

- One answer lies in *tastes*: If a good or service is a basic necessity in one's life, then price variations have a minimal effect on the quantity demanded, and these products thus have a relatively inelastic demand.
- A second answer lies in the *ease with which we can substitute* alternative goods or services for the product in question. If *Apple* Corporation had no serious competition in the smart-phone market, it could price its products even higher than in the presence of *Samsung* and *Google*, who also supply smart phones. A supplier who increases her price will lose more sales if there are ready substitutes to which buyers can switch, than if no such substitutes exist. It follows that a critical role for the marketing department in a firm is to convince buyers of the uniqueness of the firm's product.
- Where *product groups* are concerned, the price elasticity of demand for one product is necessarily higher than for the group as a whole: Suppose the price of one computer tablet brand alone falls. Buyers would be expected to substitute towards this product in large numbers – its manufacturer would find demand to be highly responsive. But if *all* brands are reduced in price, the increase in demand for any one will be more muted. In essence, the one tablet whose price falls has several close substitutes, but tablets in the aggregate do not.
- Finally, there is a *time dimension* to responsiveness, and this is explored in Section 4.3.

4.2 Price elasticities and public policy

In Chapter 3 we explored the implications of putting price floors and supply quotas in place. We saw that price floors can lead to excess supply. An important public policy question therefore is why these policies actually exist. It turns out that we can understand why with the help of elasticity concepts.

Price elasticity and expenditure

Let us return to Table 4.1 and explore what happens to total expenditure/revenue as the price varies. Since total revenue is simply the product of price times quantity it can be computed from the first two columns. The result is given in the final column. We see immediately that total expenditure on the good is highest at the midpoint of the demand curve corresponding to these data. At a price of \$5 expenditure is \$25. No other price yields more expenditure or revenue. Obviously the value \$5 is midway between the zero value and the price or quantity intercept of the demand curve in Figure 4.1. This is a general result for linear demand curves: Expenditure is greatest at the midpoint, and the mid-price corresponds to the mid-quantity on the horizontal axis.

Geometrically this can be seen from Figure 4.1. Since expenditure is the product of price and quantity, in geometric terms it is the area of the rectangle mapped out by any price-quantity combination. For example, at $\{P = \$8, Q = 2\}$ total expenditure is \$16 – the area of the rectangle bounded by these price and quantity values. Following this line of reasoning, if we were to compute the area bounded by a price of \$7 and a corresponding quantity of 3 units we get a larger rectangle – a value of \$21. This example indicates that the largest rectangle occurs at the midpoint of the demand curve. As a general geometric rule this is always the case. Hence we can conclude that the price that generates the greatest expenditure is the midpoint of a linear demand curve.

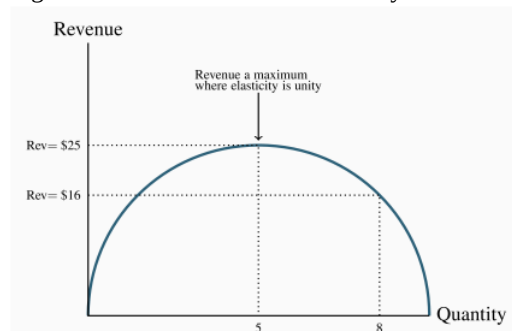
Let us now apply this rule to pricing in the market place. If our existing price is high and our goal is to generate more revenue, then we should reduce the price. Conversely, if our price is low and our goal is again to increase revenue we should raise the price. Starting from a high price let us see why this is so. By lowering the price we induce an increase in quantity demanded. Of course the lower price reduces the revenue obtained on the units already being sold at the initial high price. But since total expenditure increases at the new lower price, it must be the case that *the additional sales caused by the lower price more than compensate for this loss on the units being sold at the initial high price*. But there comes a point when this process ceases. Eventually the loss in revenue on the units being sold at the higher price is not offset by the revenue from additional quantity. We *lose a margin on so many existing units that the additional sales cannot compensate*. Accordingly revenue falls.

Note next that the top part of the demand curve is elastic and the lower part is inelastic. So, as a general rule we can state that:

A price decline (quantity increase) on an elastic segment of a demand curve necessarily increases revenue, and a price increase (quantity decline) on an inelastic segment also increases revenue.

The result is mapped in Figure 4.4, which plots total revenue as a function of the quantity demanded – columns 2 and 4 from Table 4.1. At low quantity values the price is high and the demand is elastic; at high quantity values the price is low and the demand is inelastic. The revenue maximizing point is the midpoint of the demand curve.

Figure 4.4 Total revenue and elasticity



Based upon the data in Table 4.1, revenue increases with quantity sold up to sales of 5 units. Beyond this output, the decline in price that must accompany additional sales causes revenue to decline.

We now have a general conclusion: In order to maximize the possible revenue from the sale of a good or service, it should be priced where the demand elasticity is unity.

Does this conclusion mean that every entrepreneur tries to find this magic region of the demand curve in pricing her product? Not necessarily: Most businesses seek to maximize their *profit* rather than their revenue, and so they have to focus on cost in addition to sales. We will examine this interaction in later chapters. Secondly, not every firm has control over the price they charge; the price corresponding to the unit elasticity may be too high relative to their competitors' price choices. Nonetheless, many firms, especially in the early phase of their life-cycle, focus on revenue growth rather than profit, and so, if they have any power over their price, the choice of the unit-elastic price may be appropriate.

The agriculture problem

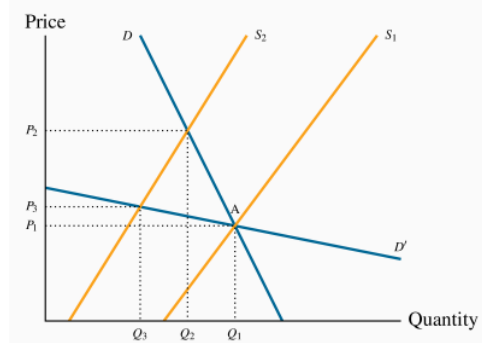
We are now in a position to address the question we posed above: Why are price floors frequently found in agricultural markets? The answer is that governments believe that the pressures of competition would force farm/food prices so low that many farmers would not be able to earn a reasonable income from farming. Accordingly, governments impose price floors. Keep in mind that price floors are prices above the market equilibrium and therefore lead to excess supply.

Since the demand for foodstuffs is inelastic we know that a higher price will induce more revenue, even with a lower quantity being sold. The government can force this outcome on the market by a policy of supply management. It can force farmers in the aggregate to bring only a specific amount of product to the market, and thus ensure that the price floor does not lead to excess supply. This is the system of supply management we observe in dairy markets in Canada, for example, and that we examined in the case of maple syrup in Chapter 3. Its supporters praise it because it helps farmers, its critics point out that higher food prices hurt lower-income households more than high-income households, and therefore it is not a good policy.

Elasticity values are frequently more informative than diagrams and figures. Our natural inclination is to view demand curves with a somewhat vertical profile as being inelastic, and demand curves with a flatter profile as elastic. But we must keep in mind that, as explained in Chapter 2, the vertical and horizontal axis of any diagram can be scaled in such a way as to change the visual impact of the data underlying the curves. But a numerical elasticity value will never deceive in this way. If its value is less than unity it is inelastic, regardless of the visual aspect of the demand curve.

At the same time, if we have two demand curves intersecting at a particular price-quantity combination, we can say that the curve with the more vertical profile is *relatively* more elastic, or less inelastic. This is illustrated in Figure 4.5. It is clear that, *at the price-quantity combination where they intersect*, the demand curve D' will yield a greater (percentage) quantity change than the demand curve D , for a given (percentage) price change. Hence, on the basis of diagrams, we can compare demand elasticities *in relative terms at a point where the two intersect*.

Figure 4.5 The impact of elasticity on quantity fluctuations



In the lower part of the demand curve D , demand is inelastic: At the point A, a shift in supply from S_1 to S_2 induces a large percentage increase in price, and a small percentage decrease in quantity demanded. In contrast, for the demand curve D' that goes through the original equilibrium, the region A is now an *elastic* region, and the impact of the supply shift is contrary: The $\% \Delta P$ is smaller and the $\% \Delta Q$ is larger.

4.3 The time horizon and inflation

The price elasticity of demand is frequently lower in the short run than in the long run. For example, a rise in the price of home heating oil may ultimately induce consumers to switch to natural gas or electricity, but such a transition may require a considerable amount of time. Time is required for decision-making and investment in new heating equipment. A further example is the elasticity of demand for tobacco. Some adults who smoke may be seriously dependent and find quitting almost impossible. Higher prices may provide a stronger incentive to reduce or quit, but successful quitters usually require several attempts before being successful. Several years may be required for the impact of a price increase to be fully apparent. Accordingly when we talk of the short run and

the long run, there is no simple rule for defining how long the long run actually is in terms of months or years. In some cases, adjustment may be complete in weeks, in other cases years.

In Chapter 2 we distinguished between real and nominal variables. The former adjust for inflation; the latter do not. Suppose all nominal variables double in value: Every good and service costs twice as much, wage rates double, dividends and rent double, etc. This implies that whatever bundle of goods was previously affordable is still affordable. Nothing has really changed. Demand behaviour is unaltered by this doubling of all prices and all incomes.

How do we reconcile this with the idea that own-price elasticities measure changes in quantity demanded as prices change? Keep in mind that elasticities measure the impact of changing one variable alone, *holding constant all of the others*. But when all variables are changing simultaneously, it is incorrect to think that the impact on quantity of one price or income change is a true measure of responsiveness or elasticity. The *price changes that go into measuring elasticities are therefore changes in prices relative to inflation*.

4.4 Cross-price elasticities – cable or satellite

The price elasticity of demand tells us about consumer responses to price changes in different regions of the demand curve, holding constant all other influences. One of those influences is the price of other goods and services. A cross-price elasticity indicates how demand is influenced by changes in the prices of other products.

The **cross-price elasticity of demand** is the percentage change in the quantity demanded of a product divided by the percentage change in the price of another.

We write the cross price elasticity of the demand for x due to a change in the price of y as

$$\epsilon_{d(x,y)} = \frac{\text{percentage change in quantity demanded of } x}{\text{percentage change in price of good } y} = \frac{\% \Delta Q_x}{\% \Delta P_y}$$

For example, if the price of cable-supply internet services declines, by how much will the demand for satellite-supply services change? The cross-price elasticity may be positive or negative. These particular goods are clearly *substitutable*, and this is reflected in a *positive* value of this cross-price elasticity: The percentage change in satellite subscribers will be negative in response to a decline in the price of cable; a negative divided by a negative is positive. In contrast, a change in the price of tablets or electronic readers should induce an opposing change in the quantity of e-books purchased: Lower tablet prices will induce greater e-book purchases. In this case the price and quantity movements are in opposite directions and the elasticity is therefore negative – the goods are complements.

Application Box 4.1 Cross-price elasticity of demand between legal and illegal marijuana

In November 2016 Canada's Parliamentary Budget Office produced a research paper on the challenges associated with pricing legalized marijuana. They proposed that taxes should be low rather than high on this product, surprising many health advocates. Specifically they argued that the legal price of marijuana should be just fractionally higher than the price in the illegal market. Otherwise marijuana users would avail of the illegal market supply, which is widely available and of high quality. Effectively their research pointed to a very high cross-price elasticity of demand. This recommendation may mean that tax revenue from marijuana sales will be small, but the size of the illegal market will decline substantially, thereby attaining a prime objective of legalization.

4.5 The income elasticity of demand

In Chapter 3 we stated that higher incomes tend to increase the quantity demanded at any price. To measure the responsiveness of demand to income changes, a unit-free measure exists: The income elasticity of demand. The income elasticity of demand is the percentage change in quantity demanded divided by a percentage change in income.

The **income elasticity of demand** is the percentage change in quantity demanded divided by a percentage change in income.

Let us use the Greek letter eta, η , to define the income elasticity of demand and I to denote income. Then,

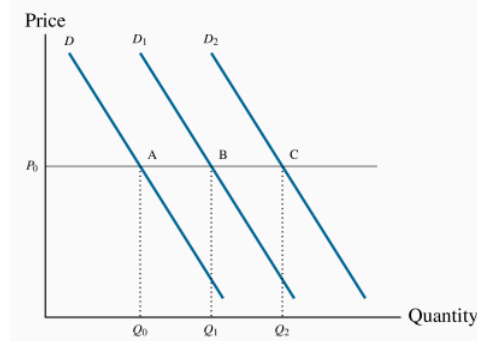
$$\eta_d = \frac{\text{percentage change in quantity demanded}}{\text{percentage change in income}} = \frac{\% \Delta Q}{\% \Delta I}$$

As an example, if monthly income increases by 10 percent, and the quantity of magazines purchased increases by 15 percent, then the income elasticity of demand for magazines is 1.5 in value ($= 15\%/10\%$). The income elasticity is generally positive, but not always – let us see why.

Normal, inferior, necessary, and luxury goods

The income elasticity of demand, in diagrammatic terms, is a percentage measure of how far the demand curve shifts in response to a change in income. Figure 4.6 shows two possible shifts. Suppose the demand curve is initially the one defined by D , and then income increases. In this example the supply curve is horizontal at the price P_0 . If the demand curve shifts to D_1 as a result, the change in quantity demanded at the existing price is $(Q_1 - Q_0)$. However, if instead the demand curve shifts to D_2 , that shift denotes a larger change in quantity $(Q_2 - Q_0)$. Since the shift in demand denoted by D_2 exceeds the shift to D_1 , the D_2 shift is more responsive to income, and therefore implies a higher income elasticity.

Figure 4.6 Income elasticity and shifts in demand



At the price P_0 , the income elasticity measures the percentage horizontal shift in demand caused by some percentage income increase. A shift from A to B reflects a lower income elasticity than a shift to C. A leftward shift in the demand curve in response to an income increase would denote a negative income elasticity – an inferior good.

In this example, the good is a *normal good*, as defined in Chapter 3, because the demand for it increases in response to income increases. If the demand curve were to shift back to the left in response to an increase in income, then the income elasticity would be negative. In such cases the goods or services are *inferior*, as defined in Chapter 3.

Finally, we distinguish between luxuries and necessities. A luxury good or service is one whose income elasticity equals or exceeds unity. A necessity is one whose income elasticity is greater than zero but less than unity. If quantity demanded is so responsive to an income increase that the percentage increase in quantity demanded exceeds the percentage increase in income, then the elasticity value is in excess of 1, and the good or service is called a luxury. In contrast, if the percentage change in quantity demanded is less than the percentage increase in income, the value is less than unity, and we call the good or service a necessity.

A **luxury** good or service is one whose income elasticity equals or exceeds unity.

A **necessity** is one whose income elasticity is greater than zero and less than unity.

Luxuries and necessities can also be defined in terms of their share of a typical budget. An income elasticity greater than unity means that the share of an individual's budget being allocated to the product is increasing. In contrast, if the elasticity is less than unity, the budget share is falling. This makes intuitive sense—luxury cars are luxury goods by this definition because they take up a larger share of the incomes of the rich than the non-rich.

Inferior goods are those for which there exist higher-quality, more expensive, substitutes. For example, lower-income households tend to satisfy their travel needs by using public transit. As income rises, households may reduce their reliance on public transit in favour of automobile use (despite the congestion and environmental impacts). Likewise, laundromats are inferior goods in the sense that, as income increases, individuals tend to purchase their own appliances and therefore use laundromat services less. Inferior goods, therefore, have a negative income elasticity: In the income elasticity equation definition, the numerator has a sign opposite to that of the denominator.

Inferior goods have negative income elasticity.

Empirical research indicates that goods like food and fuel have income elasticities less than 1; durable goods and services have elasticities slightly greater than 1; leisure goods and foreign holidays have elasticities very much greater than 1.

Income elasticities are useful in forecasting the demand for particular services and goods in a growing economy. Suppose real income is forecast to grow by 15% over the next five years. If we know that the income elasticity of demand for smart phones is 2.0, we could estimate the anticipated growth in demand by using the income elasticity formula: Since in this case $\eta = 2.0$ and $\% \Delta I = 15\%$ it follows that $2.0 = \% \Delta Q / 15\%$. Therefore the predicted demand change must be 30%.

4.6 Elasticity of supply

Now that we have developed the various dimensions of elasticity on the demand side, the analysis of elasticities on the supply side is straightforward. The elasticity of supply measures the responsiveness of the quantity supplied to a change in the price.

The **elasticity of supply** measures the responsiveness of quantity supplied to a change in the price.

$$\epsilon_s = \frac{\text{percentage change in quantity supplied}}{\text{percentage change in price}} = \frac{\% \Delta Q}{\% \Delta P}$$

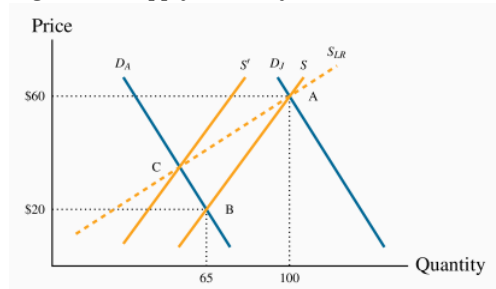
The subscript s denotes supply. This is exactly the same formula as for the demand curve, except that the quantities now come from a supply curve. Furthermore, and in contrast to the demand elasticity, the supply elasticity is generally a positive value because of the positive relationship between price and quantity supplied. The more elastic, or the more responsive, is supply to a given price change, the larger will be the elasticity value. In diagrammatic terms, this means that "flatter" supply curves have a greater elasticity than more "vertical" curves at a given price and quantity combination. Numerically the flatter curve has a larger value than the more vertical supply – try drawing a supply diagram similar to Figure 4.2. Technically, a completely vertical supply curve has a zero elasticity and a horizontal supply curve has an infinite elasticity – just as in the demand cases.

As always we keep in mind the danger of interpreting too much about the value of this elasticity from looking at the visual profiles of supply curves.

Application Box 4.2 The price of oil and the coronavirus

In January 2020, the price of oil in the US traded at \$60 per barrel. The coronavirus then struck the world economy, transport declined dramatically, and by the beginning of April the price had dropped to \$20 per barrel. The quantity of oil traded on world markets fell from approximately 100 million barrels per day in January to 65 million barrels by the start of April, a drop of 35% relative to its January level. This scenario is displayed in Figure 4.7 below.

Figure 4.7 Supply elasticity in the short run and the long run; the oil market in 2020



Demand drops suddenly between January and April. Equilibrium moves from point A to B. In the long run some producers exit and the supply curve shifts towards the origin. Following this, the equilibrium is C. Joining points such as A and C yields a long run supply curve; it is more elastic than the short run supply, when the number of suppliers is fixed.

The January equilibrium is at the point A, and the April equilibrium at point B. The move from A to B was caused by a collapse in demand as illustrated by the shift in the demand curve. We can compute the supply elasticity readily from this example. Note that it was demand that shifted rather than supply so we are observing two points on the supply curve. The supply elasticity, using arc values, is given by $(35/82.5)/(\$40/\$40) = 0.42$. So the supply curve is inelastic.

Can all oil producers survive in this bear market? The answer is no. It is inexpensive to pump oil in Saudi Arabia, but more costly to produce it from shale or tar sands, and thus some producers are not covering their costs with the price at \$20 per barrel. They do not want to shut their operations because closing down and reopening is expensive. They hang on in the hope that the price will return towards \$60. If demand does not recover, some suppliers exit the industry with the passage of time. With fewer suppliers the supply curve shifts towards the origin and ultimately another equilibrium price and quantity are established. Call this new equilibrium point C.

The supply elasticity takes on a different value in the short run than in the long run. Supply is more inelastic in the short run. A line through the points A and C would represent the long-run supply curve for the industry. It must be more elastic than the short run supply because the industry has had time to adjust - it is more flexible with the passage of time.

4.7 Elasticities and tax incidence

Elasticity values are critical in determining the impact of a government's taxation policies. The spending and taxing activities of the government influence the use of the economy's resources. By taxing cigarettes, alcohol and fuel, the government can restrict their use; by taxing income, the government influences the amount of time people choose to work. Taxes have a major impact on almost every sector of the Canadian economy.

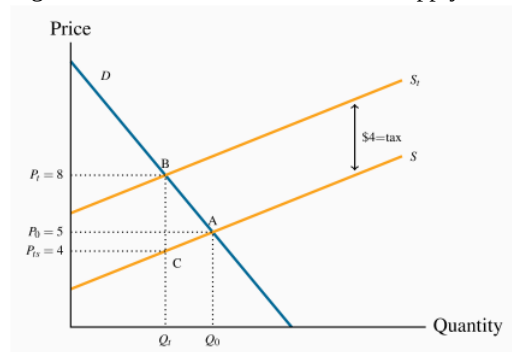
To illustrate the role played by demand and supply elasticities in tax analysis, we take the example of a sales tax. These can be of the *specific* or *ad valorem* type. A *specific* tax involves a fixed dollar levy per unit of a good sold (e.g., \$10 per airport departure). An *ad valorem* tax is a percentage levy, such as Canada's Goods and Services tax (e.g., 5 percent on top of the retail price of goods and services). The impact of each type of tax is similar, and we will use the specific tax in our example below.

A layperson's view of a sales tax is that the tax is borne by the consumer. That is to say, if no sales tax were imposed on the good or service in question, the price paid by the consumer would be the same net of tax price as exists when the tax is in place. Interestingly, this is not always the case. The study of the incidence of taxes is the study of who really bears the tax burden, and this in turn depends upon supply and demand elasticities.

Tax Incidence describes how the burden of a tax is shared between buyer and seller.

Consider Figures 4.8 and 4.9, which define an imaginary market for inexpensive wine. Let us suppose that, without a tax, the equilibrium price of a bottle of wine is \$5, and Q_0 is the equilibrium quantity traded. The pre-tax equilibrium is at the point A. The government now imposes a specific tax of \$4 per bottle. The impact of the tax is represented by an upward shift in supply of \$4: Regardless of the price that the consumer pays, \$4 of that price must be remitted to the government. As a consequence, the price paid to the supplier must be \$4 less than the consumer price, and this is represented by twin supply curves: One defines the price at which the supplier is willing to supply (S), and the other is the tax-inclusive supply curve that the consumer faces (S_t).

Figure 4.8 Tax incidence with elastic supply



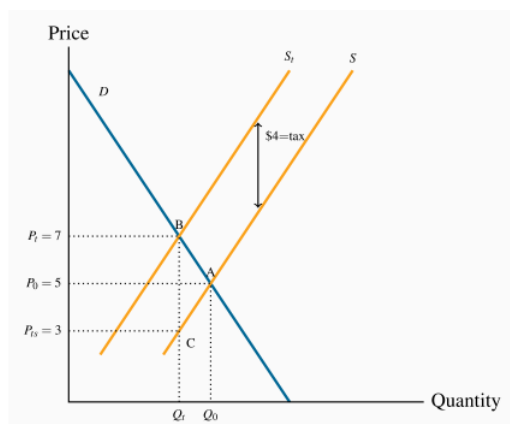
The imposition of a specific tax of \$4 shifts the supply curve vertically by \$4. The final price at B (P_t) increases by \$3 over the equilibrium price at A. At the new quantity traded, Q_t , the supplier gets \$4 per unit (P_s), the government gets \$4 also and the consumer pays \$8. The greater part of the incidence is upon the buyer, on account of the relatively elastic supply curve: His price increases by \$3 of the \$4 tax.

The introduction of the tax in Figure 4.8 means that consumers now face the supply curve S_t . The new equilibrium is at point B. Note that the price has increased by less than the full amount of the tax—in this example it has increased by \$3. This is because the reduced quantity at B is provided at a lower supply price: The supplier is willing to supply the quantity Q_t at a price defined by C (\$4), which is lower than the price at A (\$5).

So what is the incidence of the \$4 tax? Since the market price has increased from \$5 to \$8, and the price obtained by the supplier has fallen by \$1, we say that the incidence of the tax falls mainly on the consumer: The price to the consumer has risen by three dollars and the price received by the supplier has fallen by just one dollar.

Consider now Figure 4.9, where the supply curve is less elastic, and the demand curve is unchanged. Again the supply curve must shift upward with the imposition of the \$4 specific tax. But here the price received by the supplier is lower than in Figure 4.8, and the price paid by the consumer does not rise as much – the incidence is different. The consumer faces a price increase that is one-half, rather than three-quarters, of the tax value. The supplier faces a lower supply price, and bears a higher share of the tax.

Figure 4.9 Tax incidence with inelastic supply



The imposition of a specific tax of \$4 shifts the supply curve vertically by \$4. The final price at B (P_t) increases by \$2 over the no-tax price at A. At the new quantity traded, Q_t , the supplier gets \$3 per unit (P_s), the government gets \$4 also and the consumer pays \$7. The incidence is shared equally by suppliers and demanders.

We can conclude from this example that, for any given demand, *the more elastic is supply, the greater is the price increase* in response to a given tax. Furthermore, a *more elastic supply curve* means that the *incidence falls more on the consumer*; while a *less elastic supply curve* means the *incidence falls more on the supplier*. This conclusion can be verified by drawing a third version of Figure 4.8 and 4.9, in which the supply curve is horizontal – perfectly elastic. When the tax is imposed the price to the consumer increases by the full value of the tax, and the full incidence falls on the buyer. While this case corresponds to the layperson's intuition of the incidence of a tax, economists recognize it as a special case of the more general outcome, where the incidence falls on both the supply side and the demand side.

These are key results in the theory of taxation. It is equally the case that *the incidence of the tax depends upon the demand elasticity*. In Figure 4.8 and 4.9 we used the same demand curve. However, it is not difficult to see that, if we were to redo the exercise with a demand curve of a different elasticity, the incidence would not be identical. At the same time, the general result on supply elasticities still holds. We will return to this material in Chapter 5.

Statutory incidence

In the above example the tax is analyzed by means of shifting the supply curve. This implies that the supplier is obliged to charge the consumer a tax and then return this tax revenue to the government. But suppose the supplier did not bear the obligation to collect the revenue; instead the buyer is required to send the tax revenue to the government, as in the case of employers who are required to deduct income tax from their employees' pay packages (the employers here are the demanders). If this were the case we could analyze the impact of the tax by reducing the market *demand* curve by the \$4. This is because the demand curve reflects what buyers are willing to pay, and when suppliers are paid in the presence of the tax they will be paid the buyers' demand price minus the tax that the buyers must pay. It is not difficult to show that whether we move the supply curve upward (to reflect the responsibility of the supplier to pay the government) or move the demand curve downward, the outcome is the same – in the sense that the same price and quantity will be traded in each case. Furthermore the incidence of the tax, measured by how the price change is apportioned between the buyers and sellers is also unchanged.

Tax revenues and tax rates

It is useful to relate elasticity values to the policy question of the impact of higher or lower taxes on government tax revenue. Consider a situation in which a tax is already in place and the government considers increasing the rate of tax. Can an understanding of elasticities inform us on the likely outcome? The answer is yes. Suppose that at the initial tax-inclusive price demand is inelastic. We know immediately that a tax rate increase that increases the price must increase total expenditure. Hence the outcome is that the government will get a higher share of an increased total expenditure. In contrast, if demand is elastic at the initial tax-inclusive price a tax rate increase that leads to a higher price will *decrease* total expenditure. In this case the government will get a larger share of a smaller pie – not as valuable from a tax-revenue standpoint as a larger share of a larger pie.

4.8 Technical tricks with elasticities

We can easily compute elasticities at any *point* on a demand curve, rather than over a range or arc, by using the explicit formula for the demand curve. To see this note that we can rewrite Equation 4.1 as:

$$\epsilon_d = \frac{\% \Delta Q}{\% \Delta P} = \frac{\Delta Q / Q}{\Delta P / P} = \frac{\Delta Q}{\Delta P} \times \frac{P}{Q}$$

The first term in the final form of this expression ($\Delta Q / \Delta P$) is obtained from the slope of the demand curve, and the second term (P / Q) is defined by the point on the curve that interests us. For example if our demand curve is $P = 10 - 1 \times Q$, then $\Delta P / \Delta Q = -1$. Inverting this to get $\Delta Q / \Delta P$ yields -1 also. So, the point elasticity value at $\{P = \$6, Q = 4\}$ is $-1 \times 6 / 4 = -1.5$. This formula provides the elasticity value at a particular point on the demand curve, rather than over a range of values or an arc. Consequently it is called the point elasticity of demand. And obviously we could apply it to a demand curve that is not linear provided we know the mathematical form and are able to establish the slope.

Key Terms

Price elasticity of demand is measured as the percentage change in quantity demanded, divided by the percentage change in price.

Demand is elastic if the price elasticity is greater than unity. It is **inelastic** if the value lies between unity and 0. It is **unit elastic** if the value is exactly one.

Cross-price elasticity of demand is the percentage change in the quantity demanded of a product divided by the percentage change in the price of another.

Income elasticity of demand is the percentage change in quantity demanded divided by a percentage change in income.

Luxury good or service is one whose income elasticity equals or exceeds unity.

Necessity is one whose income elasticity is greater than zero and is less than unity.

Inferior goods have a negative income elasticity.

Elasticity of supply is defined as the percentage change in quantity supplied divided by the percentage change in price.

Tax Incidence describes how the burden of a tax is shared between buyer and seller.

Exercises for Chapter 4

EXERCISE 4.1

Consider the information in the table below that describes the demand for movie rentals from your on-line supplier Instant Flicks.

Price per movie (\$)	Quantity demanded	Total revenue	Elasticity of demand
2	1200		
3	1100		
4	1000		
5	900		
6	800		
7	700		
8	600		

1. Either on graph paper or a spreadsheet, map out the demand curve.
2. In column 3, insert the total revenue generated at each price.
3. At what price is total revenue maximized?
4. In column 4, compute the elasticity of demand corresponding to each \$1 price reduction, using the average price and quantity at each state.
5. Do you see a connection between your answers in parts (c) and (d)?

EXERCISE 4.2

Your fruit stall has 100 ripe bananas that must be sold today. Your supply curve is therefore vertical. From past experience, you know that these 100 bananas will all be sold if the price is set at 40 cents per unit.

1. Draw a supply and demand diagram illustrating the market equilibrium price and quantity.
2. The demand elasticity is -0.5 at the equilibrium price. But you now discover that 10 of your bananas are rotten and cannot be sold. Draw the new supply curve and calculate the percentage price increase that will be associated with the new equilibrium, on the basis of your knowledge of the demand elasticity.

EXERCISE 4.3

University fees in the State of Nirvana have been frozen in real terms for 10 years. During this period enrolments increased by 20 percent, reflecting an increase in demand. This means the supply curve is horizontal at a given price.

1. Draw a supply curve and two demand curves to represent the two equilibria described.
2. Can you estimate a price elasticity of demand for university education in this market?
3. In contrast, during the same time period fees in a neighbouring state (where supply is also horizontal) increased by 60 percent and enrolments increased by 15 percent. Illustrate this situation in a diagram, where supply is again horizontal.

EXERCISE 4.4

Consider the demand curve defined by the information in the table below.

Price of movies	Quantity demanded	Total revenue	Elasticity of demand
2	200		
3	150		
4	120		
5	100		

1. Plot the demand curve to scale and note that it is non-linear.
2. Compute the total revenue at each price.
3. Compute the arc elasticity of demand for each of the three price segments.

EXERCISE 4.5

Waterson Power Corporation's regulator has just allowed a rate increase from 9 to 11 cents per kilowatt hour of electricity. The short-run demand elasticity is -0.6 and the long-run demand elasticity is -1.2 at the current price.

1. What will be the percentage reduction in power demanded in the short run (use the midpoint 'arc' elasticity formula)?
2. What will be the percentage reduction in power demanded in the long run?
3. Will revenues increase or decrease in the short and long runs?

EXERCISE 4.6

Consider the own- and cross-price elasticity data in the table below.

		% change in price		
		CDs	Magazines	Cappuccinos
% change in quantity	CDs	-0.25	0.06	0.01
	Magazines	-0.13	-1.20	0.27
	Cappuccinos	0.07	0.41	-0.85

1. For which of the goods is demand elastic and for which is it inelastic?
2. What is the effect of an increase in the price of CDs on the purchase of magazines and cappuccinos? What does this suggest about the relationship between CDs and these other commodities; are they substitutes or complements?

3. In graphical terms, if the price of CDs or the price of cappuccinos increases, illustrate how the demand curve for magazines shifts.

EXERCISE 4.7

You are responsible for running the Speedy Bus Company and have information about the elasticity of demand for bus travel: The own-price elasticity is -1.4 at the current price. A friend who works in the competing railway company also tells you that she has estimated the cross-price elasticity of train-travel demand with respect to the price of bus travel to be 1.7 .

1. As an economic analyst, would you advocate an increase or decrease in the price of bus tickets if you wished to increase revenue for Speedy?
2. Would your price decision have any impact on train ridership?

EXERCISE 4.8

A household's income and restaurant visits are observed at different points in time. The table below describes the pattern.

Income (\$)	Restaurant visits	Income elasticity of demand
16,000	10	
24,000	15	
32,000	18	
40,000	20	
48,000	22	
56,000	23	
64,000	24	

1. Construct a scatter diagram showing quantity on the vertical axis and income on the horizontal axis.
2. Is there a positive or negative relationship between these variables?
3. Compute the income elasticity for each income increase, using midpoint values.
4. Are restaurant meals a normal or inferior good?

EXERCISE 4.9

The demand for bags of candy is given by $P=48-0.2Q$, and the supply by $P=Q$. The demand intercepts here are $P = \$48$ and $Q=240$; the supply curve is a 45 degree straight line through the origin.

1. Illustrate the resulting market equilibrium in a diagram knowing that the demand intercepts are $\{ \$48, 240 \}$, and that the supply curve is a 45 degree line through the origin.
2. If the government now puts a \$12 tax on all such candy bags, illustrate on a diagram how the supply curve will change.
3. Instead of the specific tax imposed in part (b), a percentage tax (ad valorem) equal to 30 percent is imposed. Illustrate how the supply curve would change.

EXERCISE 4.10

Optional: Consider the demand curve $P=100-2Q$. The supply curve is given by $P=30$.

1. Draw the supply and demand curves to scale, knowing that the demand curve intercepts are \$100 and 50, and compute the equilibrium price and quantity in this market.
2. If the government imposes a tax of \$10 per unit, draw the new equilibrium and compute the new quantity traded and the amount of tax revenue generated.
3. Is demand elastic or inelastic in this price range? [*Hint:* you should be able to answer this without calculations, by observing the figure you have constructed.]

EXERCISE 4.11

Optional: The supply of Henry's hamburgers is given by $P=2+0.5Q$; demand is given by $Q=20$.

1. Illustrate and compute the market equilibrium, knowing that the supply curve has an intercept of \$2 and a slope of 0.5.
2. A specific tax of \$3 per unit is subsequently imposed and that shifts the supply curve upwards and parallel by \$3, to become $P=5+0.5Q$. Solve for the equilibrium price and quantity after the tax.
3. Insert the post-tax supply curve along with the pre-tax supply curve, and determine who bears the burden of the tax.

This page titled [4: Measures of response- Elasticities](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.1: Price responsiveness of demand

Put yourself in the position of an entrepreneur. One of your many challenges is to price your product appropriately. You may be Michael Dell choosing a price for your latest computer, or the local restaurant owner pricing your table d'hôte, or you may be pricing your part-time snow-shoveling service. A key component of the pricing decision is to know how *responsive* your market is to variations in your pricing. How we measure responsiveness is the subject matter of this chapter.

We begin by analyzing the responsiveness of consumers to price changes. For example, consumers tend not to buy much more or much less food in response to changes in the general price level of food. This is because food is a pretty basic item for our existence. In contrast, if the price of new textbooks becomes higher, students may decide to search for a second-hand copy, or make do with lecture notes from their friends or downloads from the course web site. In the latter case students have ready alternatives to the new text book, and so their expenditure patterns can be expected to reflect these options, whereas it is hard to find alternatives to food. In the case of food consumers are not very responsive to price changes; in the case of textbooks they are. The word 'elasticity' that appears in this chapter title is just another term for this concept of responsiveness. Elasticity has many different uses and interpretations, and indeed more than one way of being measured in any given situation. Let us start by developing a suitable numerical measure.

The slope of the demand curve suggests itself as one measure of responsiveness: If we lowered the price of a good by \$1, for example, how many more units would we sell? The difficulty with this measure is that it does not serve us well when comparing different products. One dollar may be a substantial part of the price of your morning coffee and croissant, but not very important if buying a computer or tablet. Accordingly, when goods and services are measured in different units (croissants versus tablets), or when their prices are very different, it is often best to use a *percentage* change measure, which is *unit-free*.

The price elasticity of demand is measured as the percentage change in quantity demanded, divided by the percentage change in price. Although we introduce several other elasticity measures later, when economists speak of the demand elasticity they invariably mean the price elasticity of demand defined in this way.

The **price elasticity of demand** is measured as the percentage change in quantity demanded, divided by the percentage change in price.

The price elasticity of demand can be written in different forms. We will use the Greek letter epsilon, ϵ , as a shorthand symbol, with a subscript d to denote demand, and the capital delta, Δ , to denote a change. Therefore, we can write

$$\text{Price elasticity of demand} = \epsilon_d = \frac{\text{Percentage change in quantity demanded}}{\text{Percentage change in price}}$$

or, using a shortened expression,

$$\epsilon_d = \frac{\% \Delta Q}{\% \Delta P} \quad (4.1)$$

Calculating the value of the elasticity is not difficult. If we are told that a 10 percent price increase reduces the quantity demanded by 20 percent, then the elasticity value is $-20\%/10\% = -2$. The negative sign denotes that price and quantity move in opposite directions, but for brevity the negative sign is often omitted.

Consider now the data in Table 4.1 and the accompanying Figure 4.1. These data reflect the demand relation for natural gas that we introduced in Chapter 3. Note first that, when the price and quantity change, we must decide what *reference price and quantity* to use in the percentage change calculation in the definition above. We could use the initial or final price-quantity combination, or an average of the two. Each choice will yield a slightly different numerical value for the elasticity. The best convention is to *use the midpoint of the price values and the corresponding midpoint of the quantity values*. This ensures that the elasticity value is the same regardless of whether we start at the higher price or the lower price. Using the subscript 1 to denote the initial value and 2 the final value:

$$\text{Average quantity } \bar{Q} = (Q_1 + Q_2)/2$$

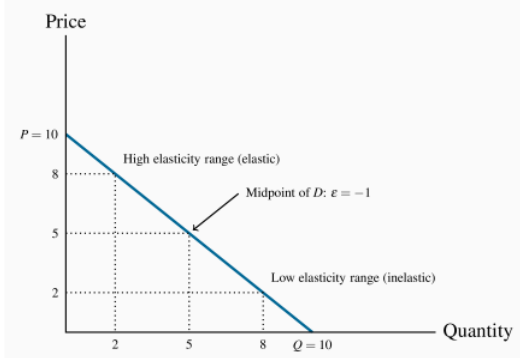
$$\text{Average price } \bar{P} = (P_1 + P_2)/2$$

Table 4.1 The demand for natural gas: Elasticities and revenue

Price (\$)	Quantity	Elasticity	Total
	demand	value	revenue (\$)

10	0		0
9	1	-9.0	9
8	2		16
7	3	-2.33	21
6	4		24
5	5	-1.0	25
4	6		24
3	7	-0.43	21
2	8		16
1	9	-0.11	9
0	10		0

Elasticity calculations are based upon \$2 price changes.
Figure 4.1 Elasticity variation with linear demand



In the high-price region of the demand curve the elasticity takes on a high value. At the midpoint of a linear demand curve the elasticity takes on a value of one, and at lower prices the elasticity value continues to fall.

Using this rule, consider now the value of ϵ_d when price drops from \$10.00 to \$8.00. The change in price is \$2.00 and the average price is therefore \$9.00 $[(\$10.00 + \$8.00)/2]$. On the quantity side, demand goes from zero to 2 units (measured in thousands of cubic feet), and the average quantity demanded is therefore $(0+2)/2=1$. Putting these numbers into the formula yields:

$$\epsilon_d = \frac{(Q_2 - Q_1)/\bar{Q}}{(P_2 - P_1)/\bar{P}} = \frac{(2/1)}{-(2/9)} = -\left(\frac{2}{1}\right) \times \left(\frac{9}{2}\right) = -9.$$

Note that the price has declined in this instance and thus the change in price is negative. Continuing down the table in this fashion yields the full set of elasticity values in the third column.

The demand elasticity is said to be *high* if it is a large negative number; the large number denotes a high degree of sensitivity. Conversely, the elasticity is *low* if it is a small negative number. High and low refer to the size of the number, ignoring the negative sign. The term *arc elasticity* is also used to define what we have just measured, indicating that it defines consumer responsiveness over a segment or *arc* of the demand curve.

It is helpful to analyze this numerical example by means of the corresponding demand curve that is plotted in Figure 4.1, and which we used in Chapter 3. It is a straight-line demand curve; but, despite this, the elasticity is not constant. At high prices the elasticity is high; at low prices it is low. The intuition behind this pattern is as follows. When the price is high, a given price change represents a small *percentage* change, because the average price in the price-term denominator is large. At high prices the quantity demanded is small and therefore the percentage quantity change tends to be large due to the small quantity value in its denominator. In sum, at high prices the elasticity is large; it contains a large numerator and a small denominator. By the same reasoning, at low prices the elasticity is small.

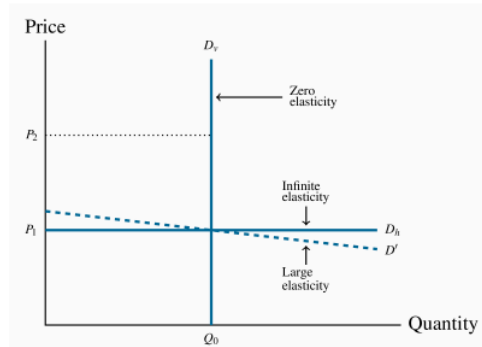
We can carry this reasoning one step further to see what happens when the demand curve intersects the axes. At the horizontal axis the average price is tending towards zero. Since this extremely small value appears in the denominator of the price term it means that the price term as a whole is extremely large. Accordingly, with an extremely large value in the denominator of the elasticity expression, the whole ratio is tending towards a zero value. By the same reasoning the elasticity value at the vertical intercept is tending towards an infinitely large value.

Extreme cases

The elasticity decreases in going from high prices to low prices. This is true for most non-linear demand curves also. Two exceptions are when the demand curve is horizontal and when it is vertical.

When the demand curve is vertical, no quantity change results from a change in price from P_1 to P_2 , as illustrated in Figure 4.2 using the demand curve D_v . Therefore, the numerator in Equation 4.1 is zero, and the elasticity has a zero value.

Figure 4.2 Limiting cases of price elasticity

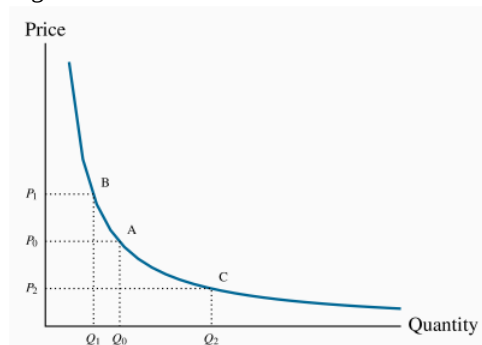


When the demand curve is vertical (D_v), the elasticity is zero: A change in price from P_1 to P_2 has no impact on the quantity demanded because the numerator in the elasticity formula has a zero value. When D becomes more horizontal the elasticity becomes larger and larger at P_1 , eventually becoming infinite.

In the horizontal case, we say that the elasticity is *infinite*, which means that any percentage price change brings forth an infinite quantity change! This case is also illustrated in Figure 4.2 using the demand curve D_h . As with the vertical demand curve, this is not immediately obvious. So consider a demand curve that is almost horizontal, such as D' instead of D_h . In this instance, we can achieve large changes in quantity demanded by implementing very small price changes. In terms of Equation 4.1, the numerator is large and the denominator small, giving rise to a large elasticity. Now imagine that this demand curve becomes ever more elastic (horizontal). The same quantity response can be obtained with a smaller price change, and hence the elasticity is larger. Pursuing this idea, we can say that, *as the demand curve becomes ever more elastic, the elasticity value tends towards infinity*.

A non-linear demand curve is illustrated in Figure 4.3. If price increases from P_0 to P_1 , the corresponding quantity change is given by $(Q_0 - Q_1)$. When the price declines to P_2 the quantity increases from Q_0 to Q_2 . When statisticians study data to determine how responsive purchases are to price changes they do not always find a linear relationship between price and quantity. But a linear relationship is frequently a good approximation or representation of actual data and we will continue to analyze responsiveness in a linear framework in this chapter.

Figure 4.3 Non-linear demand curves



When the demand curve is non-linear the slope changes with the price. Hence, equal price changes do not lead to equal quantity changes: The quantity change associated with a change in price from P_0 to P_1 is smaller than the change in quantity associated with

the same change in price from P_0 to P_2 .

Elastic and inelastic demands

While the elasticity value falls as we move down the demand curve, an important dividing line occurs at the value of -1 . This is illustrated in Table 4.1, and is a property of all straight-line demand curves. Disregarding the negative sign, demand is said to be elastic if the price elasticity is greater than unity, and inelastic if the value lies between unity and 0. It is unit elastic if the value is exactly one.

Demand is **elastic** if the price elasticity is greater than unity. It is **inelastic** if the value lies between unity and 0. It is **unit elastic** if the value is exactly one.

Economists frequently talk of goods as having a "high" or "low" demand elasticity. What does this mean, given that the elasticity varies throughout the length of a demand curve? It signifies that, *at the price usually charged*, the elasticity has a high or low value. For example, your weekly demand for regular coffee at Starbucks might be unresponsive to variations in price around the value of \$3.00, but if the price were \$6, you might be more responsive to price variations. Likewise, when we stated at the beginning of this chapter that the demand for food tends to be inelastic, we really meant that *at the price we customarily face for food*, demand is inelastic.

Determinants of price elasticity

Why is it that the price elasticities for some goods and services are high and for others low?

- One answer lies in *tastes*: If a good or service is a basic necessity in one's life, then price variations have a minimal effect on the quantity demanded, and these products thus have a relatively inelastic demand.
- A second answer lies in the *ease with which we can substitute* alternative goods or services for the product in question. If *Apple* Corporation had no serious competition in the smart-phone market, it could price its products even higher than in the presence of *Samsung* and *Google*, who also supply smart phones. A supplier who increases her price will lose more sales if there are ready substitutes to which buyers can switch, than if no such substitutes exist. It follows that a critical role for the marketing department in a firm is to convince buyers of the uniqueness of the firm's product.
- Where *product groups* are concerned, the price elasticity of demand for one product is necessarily higher than for the group as a whole: Suppose the price of one computer tablet brand alone falls. Buyers would be expected to substitute towards this product in large numbers – its manufacturer would find demand to be highly responsive. But if *all* brands are reduced in price, the increase in demand for any one will be more muted. In essence, the one tablet whose price falls has several close substitutes, but tablets in the aggregate do not.
- Finally, there is a *time dimension* to responsiveness, and this is explored in Section 4.3.

This page titled [4.1: Price responsiveness of demand](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.2: Price elasticities and public policy

In Chapter 3 we explored the implications of putting price floors and supply quotas in place. We saw that price floors can lead to excess supply. An important public policy question therefore is why these policies actually exist. It turns out that we can understand why with the help of elasticity concepts.

Price elasticity and expenditure

Let us return to Table 4.1 and explore what happens to total expenditure/revenue as the price varies. Since total revenue is simply the product of price times quantity it can be computed from the first two columns. The result is given in the final column. We see immediately that total expenditure on the good is highest at the midpoint of the demand curve corresponding to these data. At a price of \$5 expenditure is \$25. No other price yields more expenditure or revenue. Obviously the value \$5 is midway between the zero value and the price or quantity intercept of the demand curve in Figure 4.1. This is a general result for linear demand curves: Expenditure is greatest at the midpoint, and the mid-price corresponds to the mid-quantity on the horizontal axis.

Geometrically this can be seen from Figure 4.1. Since expenditure is the product of price and quantity, in geometric terms it is the area of the rectangle mapped out by any price-quantity combination. For example, at $\{P = \$8, Q = 2\}$ total expenditure is \$16 – the area of the rectangle bounded by these price and quantity values. Following this line of reasoning, if we were to compute the area bounded by a price of \$7 and a corresponding quantity of 3 units we get a larger rectangle – a value of \$21. This example indicates that the largest rectangle occurs at the midpoint of the demand curve. As a general geometric rule this is always the case. Hence we can conclude that the price that generates the greatest expenditure is the midpoint of a linear demand curve.

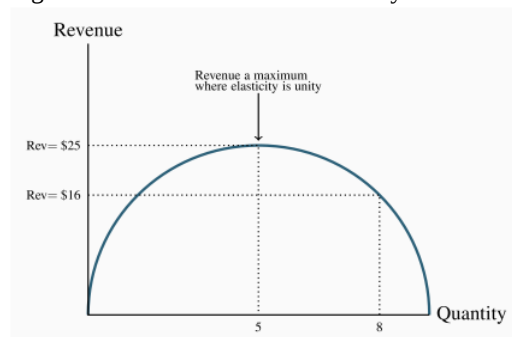
Let us now apply this rule to pricing in the market place. If our existing price is high and our goal is to generate more revenue, then we should reduce the price. Conversely, if our price is low and our goal is again to increase revenue we should raise the price. Starting from a high price let us see why this is so. By lowering the price we induce an increase in quantity demanded. Of course the lower price reduces the revenue obtained on the units already being sold at the initial high price. But since total expenditure increases at the new lower price, it must be the case that *the additional sales caused by the lower price more than compensate for this loss on the units being sold at the initial high price*. But there comes a point when this process ceases. Eventually the loss in revenue on the units being sold at the higher price is not offset by the revenue from additional quantity. We *lose a margin on so many existing units that the additional sales cannot compensate*. Accordingly revenue falls.

Note next that the top part of the demand curve is elastic and the lower part is inelastic. So, as a general rule we can state that:

A price decline (quantity increase) on an elastic segment of a demand curve necessarily increases revenue, and a price increase (quantity decline) on an inelastic segment also increases revenue.

The result is mapped in Figure 4.4, which plots total revenue as a function of the quantity demanded – columns 2 and 4 from Table 4.1. At low quantity values the price is high and the demand is elastic; at high quantity values the price is low and the demand is inelastic. The revenue maximizing point is the midpoint of the demand curve.

Figure 4.4 Total revenue and elasticity



Based upon the data in Table 4.1, revenue increases with quantity sold up to sales of 5 units. Beyond this output, the decline in price that must accompany additional sales causes revenue to decline.

We now have a general conclusion: In order to maximize the possible revenue from the sale of a good or service, it should be priced where the demand elasticity is unity.

Does this conclusion mean that every entrepreneur tries to find this magic region of the demand curve in pricing her product? Not necessarily: Most businesses seek to maximize their *profit* rather than their revenue, and so they have to focus on cost in addition to sales. We will examine this interaction in later chapters. Secondly, not every firm has control over the price they charge; the price corresponding to the unit elasticity may be too high relative to their competitors' price choices. Nonetheless, many firms, especially in the early phase of their life-cycle, focus on revenue growth rather than profit, and so, if they have any power over their price, the choice of the unit-elastic price may be appropriate.

The agriculture problem

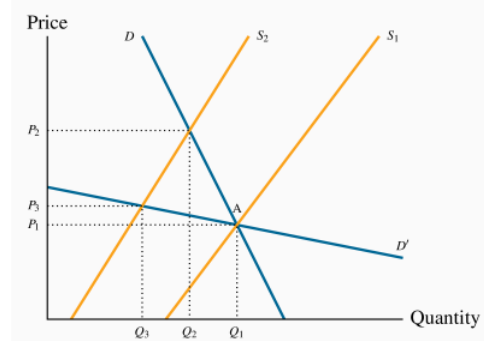
We are now in a position to address the question we posed above: Why are price floors frequently found in agricultural markets? The answer is that governments believe that the pressures of competition would force farm/food prices so low that many farmers would not be able to earn a reasonable income from farming. Accordingly, governments impose price floors. Keep in mind that price floors are prices above the market equilibrium and therefore lead to excess supply.

Since the demand for foodstuffs is inelastic we know that a higher price will induce more revenue, even with a lower quantity being sold. The government can force this outcome on the market by a policy of supply management. It can force farmers in the aggregate to bring only a specific amount of product to the market, and thus ensure that the price floor does not lead to excess supply. This is the system of supply management we observe in dairy markets in Canada, for example, and that we examined in the case of maple syrup in Chapter 3. Its supporters praise it because it helps farmers, its critics point out that higher food prices hurt lower-income households more than high-income households, and therefore it is not a good policy.

Elasticity values are frequently more informative than diagrams and figures. Our natural inclination is to view demand curves with a somewhat vertical profile as being inelastic, and demand curves with a flatter profile as elastic. But we must keep in mind that, as explained in Chapter 2, the vertical and horizontal axis of any diagram can be scaled in such a way as to change the visual impact of the data underlying the curves. But a numerical elasticity value will never deceive in this way. If its value is less than unity it is inelastic, regardless of the visual aspect of the demand curve.

At the same time, if we have two demand curves intersecting at a particular price-quantity combination, we can say that the curve with the more vertical profile is *relatively* more elastic, or less inelastic. This is illustrated in Figure 4.5. It is clear that, *at the price-quantity combination where they intersect*, the demand curve D' will yield a greater (percentage) quantity change than the demand curve D , for a given (percentage) price change. Hence, on the basis of diagrams, we can compare demand elasticities *in relative terms at a point where the two intersect*.

Figure 4.5 The impact of elasticity on quantity fluctuations



In the lower part of the demand curve D , demand is inelastic: At the point A, a shift in supply from S_1 to S_2 induces a large percentage increase in price, and a small percentage decrease in quantity demanded. In contrast, for the demand curve D' that goes through the original equilibrium, the region A is now an *elastic* region, and the impact of the supply shift is contrary: The $\% \Delta P$ is smaller and the $\% \Delta Q$ is larger.

This page titled [4.2: Price elasticities and public policy](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyrux\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.3: The time horizon and inflation

The price elasticity of demand is frequently lower in the short run than in the long run. For example, a rise in the price of home heating oil may ultimately induce consumers to switch to natural gas or electricity, but such a transition may require a considerable amount of time. Time is required for decision-making and investment in new heating equipment. A further example is the elasticity of demand for tobacco. Some adults who smoke may be seriously dependent and find quitting almost impossible. Higher prices may provide a stronger incentive to reduce or quit, but successful quitters usually require several attempts before being successful. Several years may be required for the impact of a price increase to be fully apparent. Accordingly when we talk of the short run and the long run, there is no simple rule for defining how long the long run actually is in terms of months or years. In some cases, adjustment may be complete in weeks, in other cases years.

In Chapter 2 we distinguished between real and nominal variables. The former adjust for inflation; the latter do not. Suppose all nominal variables double in value: Every good and service costs twice as much, wage rates double, dividends and rent double, etc. This implies that whatever bundle of goods was previously affordable is still affordable. Nothing has really changed. Demand behaviour is unaltered by this doubling of all prices and all incomes.

How do we reconcile this with the idea that own-price elasticities measure changes in quantity demanded as prices change? Keep in mind that elasticities measure the impact of changing one variable alone, *holding constant all of the others*. But when all variables are changing simultaneously, it is incorrect to think that the impact on quantity of one price or income change is a true measure of responsiveness or elasticity. The *price changes that go into measuring elasticities are therefore changes in prices relative to inflation*.

This page titled [4.3: The time horizon and inflation](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.4: Cross-price elasticities - cable or satellite

The price elasticity of demand tells us about consumer responses to price changes in different regions of the demand curve, holding constant all other influences. One of those influences is the price of other goods and services. A cross-price elasticity indicates how demand is influenced by changes in the prices of other products.

The **cross-price elasticity of demand** is the percentage change in the quantity demanded of a product divided by the percentage change in the price of another.

We write the cross price elasticity of the demand for x due to a change in the price of y as

$$\epsilon_{d(x,y)} = \frac{\text{percentage change in quantity demanded of } x}{\text{percentage change in price of good } y} = \frac{\% \Delta Q_x}{\% \Delta P_y}$$

For example, if the price of cable-supply internet services declines, by how much will the demand for satellite-supply services change? The cross-price elasticity may be positive or negative. These particular goods are clearly *substitutable*, and this is reflected in a *positive* value of this cross-price elasticity: The percentage change in satellite subscribers will be negative in response to a decline in the price of cable; a negative divided by a negative is positive. In contrast, a change in the price of tablets or electronic readers should induce an opposing change in the quantity of e-books purchased: Lower tablet prices will induce greater e-book purchases. In this case the price and quantity movements are in opposite directions and the elasticity is therefore negative – the goods are complements.

Application Box 4.1 Cross-price elasticity of demand between legal and illegal marijuana

In November 2016 Canada's Parliamentary Budget Office produced a research paper on the challenges associated with pricing legalized marijuana. They proposed that taxes should be low rather than high on this product, surprising many health advocates. Specifically they argued that the legal price of marijuana should be just fractionally higher than the price in the illegal market. Otherwise marijuana users would avail of the illegal market supply, which is widely available and of high quality. Effectively their research pointed to a very high cross-price elasticity of demand. This recommendation may mean that tax revenue from marijuana sales will be small, but the size of the illegal market will decline substantially, thereby attaining a prime objective of legalization.

This page titled [4.4: Cross-price elasticities - cable or satellite](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.5: The income elasticity of demand

In Chapter 3 we stated that higher incomes tend to increase the quantity demanded at any price. To measure the responsiveness of demand to income changes, a unit-free measure exists: The income elasticity of demand. The income elasticity of demand is the percentage change in quantity demanded divided by a percentage change in income.

The **income elasticity of demand** is the percentage change in quantity demanded divided by a percentage change in income.

Let us use the Greek letter eta, η , to define the income elasticity of demand and I to denote income. Then,

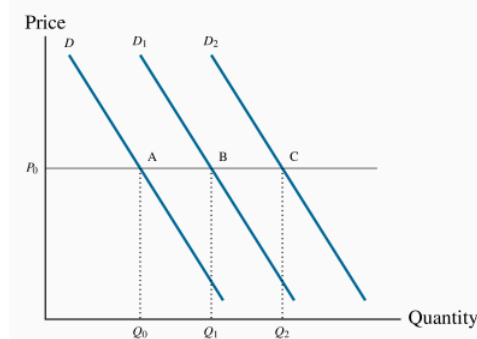
$$\eta_d = \frac{\text{percentage change in quantity demanded}}{\text{percentage change in income}} = \frac{\% \Delta Q}{\% \Delta I}$$

As an example, if monthly income increases by 10 percent, and the quantity of magazines purchased increases by 15 percent, then the income elasticity of demand for magazines is 1.5 in value ($= 15\% / 10\%$). The income elasticity is generally positive, but not always – let us see why.

Normal, inferior, necessary, and luxury goods

The income elasticity of demand, in diagrammatic terms, is a percentage measure of how far the demand curve shifts in response to a change in income. Figure 4.6 shows two possible shifts. Suppose the demand curve is initially the one defined by D , and then income increases. In this example the supply curve is horizontal at the price P_0 . If the demand curve shifts to D_1 as a result, the change in quantity demanded at the existing price is $(Q_1 - Q_0)$. However, if instead the demand curve shifts to D_2 , that shift denotes a larger change in quantity ($Q_2 - Q_0$). Since the shift in demand denoted by D_2 exceeds the shift to D_1 , the D_2 shift is more responsive to income, and therefore implies a higher income elasticity.

Figure 4.6 Income elasticity and shifts in demand



At the price P_0 , the income elasticity measures the percentage horizontal shift in demand caused by some percentage income increase. A shift from A to B reflects a lower income elasticity than a shift to C. A leftward shift in the demand curve in response to an income increase would denote a negative income elasticity – an inferior good.

In this example, the good is a *normal good*, as defined in Chapter 3, because the demand for it increases in response to income increases. If the demand curve were to shift back to the left in response to an increase in income, then the income elasticity would be negative. In such cases the goods or services are *inferior*, as defined in Chapter 3.

Finally, we distinguish between luxuries and necessities. A luxury good or service is one whose income elasticity equals or exceeds unity. A necessity is one whose income elasticity is greater than zero but less than unity. If quantity demanded is so responsive to an income increase that the percentage increase in quantity demanded exceeds the percentage increase in income, then the elasticity value is in excess of 1, and the good or service is called a luxury. In contrast, if the percentage change in quantity demanded is less than the percentage increase in income, the value is less than unity, and we call the good or service a necessity.

A **luxury** good or service is one whose income elasticity equals or exceeds unity.

A **necessity** is one whose income elasticity is greater than zero and less than unity.

Luxuries and necessities can also be defined in terms of their share of a typical budget. An income elasticity greater than unity means that the share of an individual's budget being allocated to the product is increasing. In contrast, if the elasticity is less than unity, the budget share is falling. This makes intuitive sense—luxury cars are luxury goods by this definition because they take up a larger share of the incomes of the rich than the non-rich.

Inferior goods are those for which there exist higher-quality, more expensive, substitutes. For example, lower-income households tend to satisfy their travel needs by using public transit. As income rises, households may reduce their reliance on public transit in favour of automobile use (despite the congestion and environmental impacts). Likewise, laundromats are inferior goods in the sense that, as income increases, individuals tend to purchase their own appliances and therefore use laundromat services less. Inferior goods, therefore, have a negative income elasticity: In the income elasticity equation definition, the numerator has a sign opposite to that of the denominator.

Inferior goods have negative income elasticity.

Empirical research indicates that goods like food and fuel have income elasticities less than 1; durable goods and services have elasticities slightly greater than 1; leisure goods and foreign holidays have elasticities very much greater than 1.

Income elasticities are useful in forecasting the demand for particular services and goods in a growing economy. Suppose real income is forecast to grow by 15% over the next five years. If we know that the income elasticity of demand for smart phones is 2.0, we could estimate the anticipated growth in demand by using the income elasticity formula: Since in this case $\eta = 2.0$ and $\% \Delta I = 15\%$ it follows that $2.0 = \% \Delta Q / 15\%$. Therefore the predicted demand change must be 30%.

This page titled [4.5: The income elasticity of demand](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.6: Elasticity of supply

Now that we have developed the various dimensions of elasticity on the demand side, the analysis of elasticities on the supply side is straightforward. The elasticity of supply measures the responsiveness of the quantity supplied to a change in the price.

The **elasticity of supply** measures the responsiveness of quantity supplied to a change in the price.

$$\epsilon_s = \frac{\text{percentage change in quantity supplied}}{\text{percentage change in price}} = \frac{\% \Delta Q}{\% \Delta P}$$

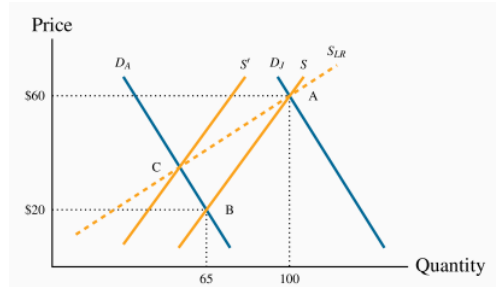
The subscript s denotes supply. This is exactly the same formula as for the demand curve, except that the quantities now come from a supply curve. Furthermore, and in contrast to the demand elasticity, the supply elasticity is generally a positive value because of the positive relationship between price and quantity supplied. The more elastic, or the more responsive, is supply to a given price change, the larger will be the elasticity value. In diagrammatic terms, this means that "flatter" supply curves have a greater elasticity than more "vertical" curves at a given price and quantity combination. Numerically the flatter curve has a larger value than the more vertical supply – try drawing a supply diagram similar to Figure 4.2. Technically, a completely vertical supply curve has a zero elasticity and a horizontal supply curve has an infinite elasticity – just as in the demand cases.

As always we keep in mind the danger of interpreting too much about the value of this elasticity from looking at the visual profiles of supply curves.

Application Box 4.2 The price of oil and the coronavirus

In January 2020, the price of oil in the US traded at \$60 per barrel. The coronavirus then struck the world economy, transport declined dramatically, and by the beginning of April the price had dropped to \$20 per barrel. The quantity of oil traded on world markets fell from approximately 100 million barrels per day in January to 65 million barrels by the start of April, a drop of 35% relative to its January level. This scenario is displayed in Figure 4.7 below.

Figure 4.7 Supply elasticity in the short run and the long run; the oil market in 2020



Demand drops suddenly between January and April. Equilibrium moves from point A to B. In the long run some producers exit and the supply curve shifts towards the origin. Following this, the equilibrium is C. Joining points such as A and C yields a long run supply curve; it is more elastic than the short run supply, when the number of suppliers is fixed.

The January equilibrium is at the point A, and the April equilibrium at point B. The move from A to B was caused by a collapse in demand as illustrated by the shift in the demand curve. We can compute the supply elasticity readily from this example. Note that it was demand that shifted rather than supply so we are observing two points on the supply curve. The supply elasticity, using arc values, is given by $(35/82.5)/(\$40/\$40) = 0.42$. So the supply curve is inelastic.

Can all oil producers survive in this bear market? The answer is no. It is inexpensive to pump oil in Saudi Arabia, but more costly to produce it from shale or tar sands, and thus some producers are not covering their costs with the price at \$20 per barrel. They do not want to shut their operations because closing down and reopening is expensive. They hang on in the hope that the price will return towards \$60. If demand does not recover, some suppliers exit the industry with the passage of time. With fewer suppliers the supply curve shifts towards the origin and ultimately another equilibrium price and quantity are established. Call this new equilibrium point C.

The supply elasticity takes on a different value in the short run than in the long run. Supply is more inelastic in the short run. A line through the points A and C would represent the long-run supply curve for the industry. It must be more elastic than the short run supply because the industry has had time to adjust - it is more flexible with the passage of time.

This page titled [4.6: Elasticity of supply](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.7: Elasticities and tax incidence

Elasticity values are critical in determining the impact of a government's taxation policies. The spending and taxing activities of the government influence the use of the economy's resources. By taxing cigarettes, alcohol and fuel, the government can restrict their use; by taxing income, the government influences the amount of time people choose to work. Taxes have a major impact on almost every sector of the Canadian economy.

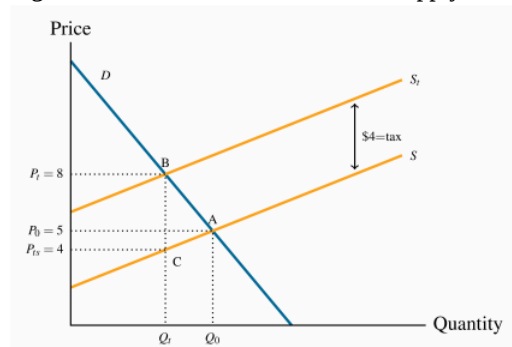
To illustrate the role played by demand and supply elasticities in tax analysis, we take the example of a sales tax. These can be of the *specific* or *ad valorem* type. A *specific* tax involves a fixed dollar levy per unit of a good sold (e.g., \$10 per airport departure). An *ad valorem* tax is a percentage levy, such as Canada's Goods and Services tax (e.g., 5 percent on top of the retail price of goods and services). The impact of each type of tax is similar, and we will use the specific tax in our example below.

A layperson's view of a sales tax is that the tax is borne by the consumer. That is to say, if no sales tax were imposed on the good or service in question, the price paid by the consumer would be the same net of tax price as exists when the tax is in place. Interestingly, this is not always the case. The study of the incidence of taxes is the study of who really bears the tax burden, and this in turn depends upon supply and demand elasticities.

Tax Incidence describes how the burden of a tax is shared between buyer and seller.

Consider Figures 4.8 and 4.9, which define an imaginary market for inexpensive wine. Let us suppose that, without a tax, the equilibrium price of a bottle of wine is \$5, and Q_0 is the equilibrium quantity traded. The pre-tax equilibrium is at the point A. The government now imposes a specific tax of \$4 per bottle. The impact of the tax is represented by an upward shift in supply of \$4: Regardless of the price that the consumer pays, \$4 of that price must be remitted to the government. As a consequence, the price paid to the supplier must be \$4 less than the consumer price, and this is represented by twin supply curves: One defines the price at which the supplier is willing to supply (S), and the other is the tax-inclusive supply curve that the consumer faces (S_t).

Figure 4.8 Tax incidence with elastic supply



The imposition of a specific tax of \$4 shifts the supply curve vertically by \$4. The final price at B (P_t) increases by \$3 over the equilibrium price at A. At the new quantity traded, Q_t , the supplier gets \$4 per unit (P_s), the government gets \$4 also and the consumer pays \$8. The greater part of the incidence is upon the buyer, on account of the relatively elastic supply curve: His price increases by \$3 of the \$4 tax.

The introduction of the tax in Figure 4.8 means that consumers now face the supply curve S_t . The new equilibrium is at point B. Note that the price has increased by less than the full amount of the tax—in this example it has increased by \$3. This is because the reduced quantity at B is provided at a lower supply price: The supplier is willing to supply the quantity Q_t at a price defined by C (\$4), which is lower than the price at A (\$5).

So what is the incidence of the \$4 tax? Since the market price has increased from \$5 to \$8, and the price obtained by the supplier has fallen by \$1, we say that the incidence of the tax falls mainly on the consumer: The price to the consumer has risen by three dollars and the price received by the supplier has fallen by just one dollar.

Consider now Figure 4.9, where the supply curve is less elastic, and the demand curve is unchanged. Again the supply curve must shift upward with the imposition of the \$4 specific tax. But here the price received by the supplier is lower than in Figure 4.8, and the price paid by the consumer does not rise as much – the incidence is different. The consumer faces a price increase that is one-half, rather than three-quarters, of the tax value. The supplier faces a lower supply price, and bears a higher share of the tax.

Figure 4.9 Tax incidence with inelastic supply



The imposition of a specific tax of \$4 shifts the supply curve vertically by \$4. The final price at B (P_t) increases by \$2 over the no-tax price at A. At the new quantity traded, Q_t , the supplier gets \$3 per unit (P_s), the government gets \$4 also and the consumer pays \$7. The incidence is shared equally by suppliers and demanders.

We can conclude from this example that, for any given demand, *the more elastic is supply, the greater is the price increase* in response to a given tax. Furthermore, a *more elastic supply curve* means that the *incidence falls more on the consumer*; while a *less elastic supply curve* means the *incidence falls more on the supplier*. This conclusion can be verified by drawing a third version of Figure 4.8 and 4.9, in which the supply curve is horizontal – perfectly elastic. When the tax is imposed the price to the consumer increases by the full value of the tax, and the full incidence falls on the buyer. While this case corresponds to the layperson's intuition of the incidence of a tax, economists recognize it as a special case of the more general outcome, where the incidence falls on both the supply side and the demand side.

These are key results in the theory of taxation. It is equally the case that *the incidence of the tax depends upon the demand elasticity*. In Figure 4.8 and 4.9 we used the same demand curve. However, it is not difficult to see that, if we were to redo the exercise with a demand curve of a different elasticity, the incidence would not be identical. At the same time, the general result on supply elasticities still holds. We will return to this material in Chapter 5.

Statutory incidence

In the above example the tax is analyzed by means of shifting the supply curve. This implies that the supplier is obliged to charge the consumer a tax and then return this tax revenue to the government. But suppose the supplier did not bear the obligation to collect the revenue; instead the buyer is required to send the tax revenue to the government, as in the case of employers who are required to deduct income tax from their employees' pay packages (the employers here are the demanders). If this were the case we could analyze the impact of the tax by reducing the market *demand* curve by the \$4. This is because the demand curve reflects what buyers are willing to pay, and when suppliers are paid in the presence of the tax they will be paid the buyers' demand price minus the tax that the buyers must pay. It is not difficult to show that whether we move the supply curve upward (to reflect the responsibility of the supplier to pay the government) or move the demand curve downward, the outcome is the same – in the sense that the same price and quantity will be traded in each case. Furthermore the incidence of the tax, measured by how the price change is apportioned between the buyers and sellers is also unchanged.

Tax revenues and tax rates

It is useful to relate elasticity values to the policy question of the impact of higher or lower taxes on government tax revenue. Consider a situation in which a tax is already in place and the government considers increasing the rate of tax. Can an understanding of elasticities inform us on the likely outcome? The answer is yes. Suppose that at the initial tax-inclusive price demand is inelastic. We know immediately that a tax rate increase that increases the price must increase total expenditure. Hence the outcome is that the government will get a higher share of an increased total expenditure. In contrast, if demand is elastic at the initial tax-inclusive price a tax rate increase that leads to a higher price will *decrease* total expenditure. In this case the government will get a larger share of a smaller pie – not as valuable from a tax-revenue standpoint as a larger share of a larger pie.

This page titled [4.7: Elasticities and tax incidence](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyrx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.8: Technical tricks with elasticities

We can easily compute elasticities at any *point* on a demand curve, rather than over a range or arc, by using the explicit formula for the demand curve. To see this note that we can rewrite Equation 4.1 as:

$$\epsilon_d = \frac{\% \Delta Q}{\% \Delta P} = \frac{\Delta Q / Q}{\Delta P / P} = \frac{\Delta Q}{\Delta P} \times \frac{P}{Q}.$$

The first term in the final form of this expression ($\Delta Q / \Delta P$) is obtained from the slope of the demand curve, and the second term (P / Q) is defined by the point on the curve that interests us. For example if our demand curve is $P = 10 - 1 \times Q$, then $\Delta P / \Delta Q = -1$. Inverting this to get $\Delta Q / \Delta P$ yields -1 also. So, the point elasticity value at $\{P = \$6, Q = 4\}$ is $-1 \times 6 / 4 = -1.5$. This formula provides the elasticity value at a particular point on the demand curve, rather than over a range of values or an arc. Consequently it is called the point elasticity of demand. And obviously we could apply it to a demand curve that is not linear provided we know the mathematical form and are able to establish the slope.

This page titled [4.8: Technical tricks with elasticities](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.9: Key Terms

Price elasticity of demand is measured as the percentage change in quantity demanded, divided by the percentage change in price.

Demand is elastic if the price elasticity is greater than unity. It is **inelastic** if the value lies between unity and 0. It is **unit elastic** if the value is exactly one.

Cross-price elasticity of demand is the percentage change in the quantity demanded of a product divided by the percentage change in the price of another.

Income elasticity of demand is the percentage change in quantity demanded divided by a percentage change in income.

Luxury good or service is one whose income elasticity equals or exceeds unity.

Necessity is one whose income elasticity is greater than zero and is less than unity.

Inferior goods have a negative income elasticity.

Elasticity of supply is defined as the percentage change in quantity supplied divided by the percentage change in price.

Tax Incidence describes how the burden of a tax is shared between buyer and seller.

This page titled [4.9: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

4.10: Exercises for Chapter 4

EXERCISE 4.1

Consider the information in the table below that describes the demand for movie rentals from your on-line supplier Instant Flicks.

Price per movie (\$)	Quantity demanded	Total revenue	Elasticity of demand
2	1200		
3	1100		
4	1000		
5	900		
6	800		
7	700		
8	600		

- Either on graph paper or a spreadsheet, map out the demand curve.
- In column 3, insert the total revenue generated at each price.
- At what price is total revenue maximized?
- In column 4, compute the elasticity of demand corresponding to each \$1 price reduction, using the average price and quantity at each state.
- Do you see a connection between your answers in parts (c) and (d)?

EXERCISE 4.2

Your fruit stall has 100 ripe bananas that must be sold today. Your supply curve is therefore vertical. From past experience, you know that these 100 bananas will all be sold if the price is set at 40 cents per unit.

- Draw a supply and demand diagram illustrating the market equilibrium price and quantity.
- The demand elasticity is -0.5 at the equilibrium price. But you now discover that 10 of your bananas are rotten and cannot be sold. Draw the new supply curve and calculate the percentage price increase that will be associated with the new equilibrium, on the basis of your knowledge of the demand elasticity.

EXERCISE 4.3

University fees in the State of Nirvana have been frozen in real terms for 10 years. During this period enrolments increased by 20 percent, reflecting an increase in demand. This means the supply curve is horizontal at a given price.

- Draw a supply curve and two demand curves to represent the two equilibria described.
- Can you estimate a price elasticity of demand for university education in this market?
- In contrast, during the same time period fees in a neighbouring state (where supply is also horizontal) increased by 60 percent and enrolments increased by 15 percent. Illustrate this situation in a diagram, where supply is again horizontal.

EXERCISE 4.4

Consider the demand curve defined by the information in the table below.

Price of movies	Quantity demanded	Total revenue	Elasticity of demand
2	200		
3	150		
4	120		
5	100		

- Plot the demand curve to scale and note that it is non-linear.
- Compute the total revenue at each price.
- Compute the arc elasticity of demand for each of the three price segments.

EXERCISE 4.5

Waterson Power Corporation's regulator has just allowed a rate increase from 9 to 11 cents per kilowatt hour of electricity. The short-run demand elasticity is -0.6 and the long-run demand elasticity is -1.2 at the current price..

- What will be the percentage reduction in power demanded in the short run (use the midpoint 'arc' elasticity formula)?
- What will be the percentage reduction in power demanded in the long run?
- Will revenues increase or decrease in the short and long runs?

EXERCISE 4.6

Consider the own- and cross-price elasticity data in the table below.

		% change in price		
		CDs	Magazines	Cappuccinos
% change in quantity	CDs	-0.25	0.06	0.01
	Magazines	-0.13	-1.20	0.27
	Cappuccinos	0.07	0.41	-0.85

- For which of the goods is demand elastic and for which is it inelastic?
- What is the effect of an increase in the price of CDs on the purchase of magazines and cappuccinos? What does this suggest about the relationship between CDs and these other commodities; are they substitutes or complements?
- In graphical terms, if the price of CDs or the price of cappuccinos increases, illustrate how the demand curve for magazines shifts.

EXERCISE 4.7

You are responsible for running the Speedy Bus Company and have information about the elasticity of demand for bus travel: The own-price elasticity is -1.4 at the current price. A friend who works in the competing railway company also tells you that she has estimated the cross-price elasticity of train-travel demand with respect to the price of bus travel to be 1.7.

- As an economic analyst, would you advocate an increase or decrease in the price of bus tickets if you wished to increase revenue for Speedy?
- Would your price decision have any impact on train ridership?

EXERCISE 4.8

A household's income and restaurant visits are observed at different points in time. The table below describes the pattern.

Income (\$)	Restaurant visits	Income elasticity of demand
16,000	10	
24,000	15	
32,000	18	
40,000	20	
48,000	22	
56,000	23	
64,000	24	

- Construct a scatter diagram showing quantity on the vertical axis and income on the horizontal axis.
- Is there a positive or negative relationship between these variables?
- Compute the income elasticity for each income increase, using midpoint values.
- Are restaurant meals a normal or inferior good?

EXERCISE 4.9

The demand for bags of candy is given by $P=48-0.2Q$, and the supply by $P=Q$. The demand intercepts here are $P = \$48$ and $Q=240$; the supply curve is a 45 degree straight line through the origin.

- Illustrate the resulting market equilibrium in a diagram knowing that the demand intercepts are $\{ \$48, 240 \}$, and that the supply curve is a 45 degree line through the origin.
- If the government now puts a \$12 tax on all such candy bags, illustrate on a diagram how the supply curve will change.
- Instead of the specific tax imposed in part (b), a percentage tax (ad valorem) equal to 30 percent is imposed. Illustrate how the supply curve would change.

EXERCISE 4.10

Optional: Consider the demand curve $P=100-2Q$. The supply curve is given by $P=30$.

- Draw the supply and demand curves to scale, knowing that the demand curve intercepts are \$100 and 50, and compute the equilibrium price and quantity in this market.
- If the government imposes a tax of \$10 per unit, draw the new equilibrium and compute the new quantity traded and the amount of tax revenue generated.
- Is demand elastic or inelastic in this price range? [*Hint:* you should be able to answer this without calculations, by observing the figure you have constructed.]

EXERCISE 4.11

Optional: The supply of Henry's hamburgers is given by $P=2+0.5Q$; demand is given by $Q=20$.

- Illustrate and compute the market equilibrium, knowing that the supply curve has an intercept of \$2 and a slope of 0.5.
- A specific tax of \$3 per unit is subsequently imposed and that shifts the supply curve upwards and parallel by \$3, to become $P=5+0.5Q$. Solve for the equilibrium price and quantity after the tax.
- Insert the post-tax supply curve along with the pre-tax supply curve, and determine who bears the burden of the tax.

This page titled [4.10: Exercises for Chapter 4](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5: Welfare economics and externalities

Chapter 5: Welfare economics and externalities

In this chapter we will explore:

5.1	Equity and efficiency
5.2	Consumer and producer surplus
5.3	Efficient market outcomes
5.4	Taxation, surplus and efficiency
5.5	Market failures – externalities
5.6	Other market failures
5.7	Environment and climate change

5.1 Equity and efficiency

In modern mixed economies, markets and governments together determine the output produced and also who benefits from that output. In this chapter we explore a very broad question that forms the core of welfare economics: Even if market forces drive efficiency, are they a good way to allocate scarce resources in view of the fact that they not only give rise to inequality and poverty, but also fail to capture the impacts of productive activity on non-market participants? Mining impacts the environment, traffic results in road fatalities, alcohol, tobacco and opioids cause premature deaths. These products all generate secondary impacts beyond their stated objective. We frequently call these *external* effects.

The analysis of markets in this larger sense involves not just economic efficiency; public policy additionally has a normative content because policies can impact the various participants in different ways and to different degrees. Welfare economics, therefore, deals with both normative and positive issues.

Welfare economics assesses how well the economy allocates its scarce resources in accordance with the goals of efficiency and equity.

Political parties on the left and right disagree on how well a market economy works. Canada's New Democratic Party emphasizes the market's failings and the need for government intervention, while the Progressive Conservative Party believes, broadly, that the market fosters choice, incentives, and efficiency. What lies behind this disagreement? The two principal factors are efficiency and equity. Efficiency addresses the question of how well the economy's resources are used and allocated. In contrast, equity deals with how society's goods and rewards are, and should be, distributed among its different members, and how the associated costs should be apportioned.

Equity deals with how society's goods and rewards are, and should be, distributed among its different members, and how the associated costs should be apportioned.

Efficiency addresses the question of how well the economy's resources are used and allocated.

Equity is also concerned with how different generations share an economy's productive capabilities: More investment today makes for a more productive economy tomorrow, but more greenhouse gases today will reduce environmental quality tomorrow. These are inter-generational questions.

Climate change caused by global warming forms one of the biggest challenges for humankind at the present time. As we shall see in this chapter, economics has much to say about appropriate policies to combat warming. Whether pollution-abatement policies should be implemented today or down the road involves considerations of equity between generations. Our first task is to develop an analytical tool which will prove vital in assessing and computing welfare benefits and costs – economic surplus.

5.2 Consumer and producer surplus

An understanding of economic efficiency is greatly facilitated as a result of understanding two related measures: Consumer surplus and producer surplus. Consumer surplus relates to the demand side of the market, producer surplus to the supply side. Producer

surplus is also termed supplier surplus. These measures can be understood with the help of a standard example, the market for city apartments.

The market for apartments

Table 5.1 and Figure 5.1 describe the hypothetical data. We imagine first a series of city-based students who are in the market for a standardized downtown apartment. These individuals are not identical; they value the apartment differently. For example, Alex enjoys comfort and therefore places a higher value on a unit than Brian. Brian, in turn, values it more highly than Cathy or Don. Evan and Frank would prefer to spend their money on entertainment, and so on. These valuations are represented in the middle column of the demand panel in Table 5.1, and also in Figure 5.1 with the highest valuations closest to the origin. The valuations reflect the willingness to pay of each consumer.

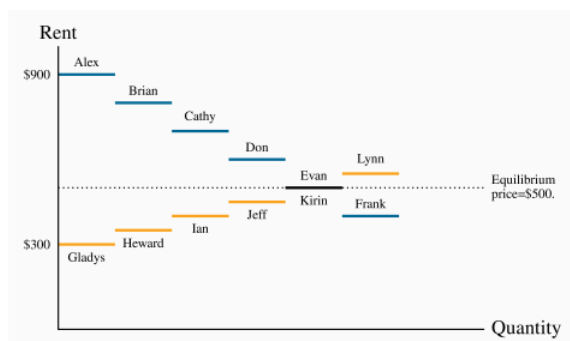
Table 5.1 Consumer and supplier surpluses

Demand		
Individual	Demand valuation	Surplus
Alex	900	400
Brian	800	300
Cathy	700	200
Don	600	100
Evan	500	0
Frank	400	0
Supply		
Individual	Reservation value	Surplus
Gladys	300	200
Heward	350	150
Ian	400	100
Jeff	450	50
Kirin	500	0
Lynn	550	0

On the supply side we imagine the market as being made up of different individuals or owners, who are willing to put their apartments on the market for different prices. Gladys will accept less rent than Heward, who in turn will accept less than Ian. The minimum prices that the suppliers are willing to accept are called *reservation* prices or values, and these are given in the lower part of Table 5.1. Unless the market price is greater than their reservation price, suppliers will hold back.

By definition, as stated in Chapter 3, the demand curve is made up of the valuations placed on the good by the various demanders. Likewise, the reservation values of the suppliers form the supply curve. If Alex is willing to pay \$900, then that is his demand price; if Heward is willing to put his apartment on the market for \$350, he is by definition willing to supply it for that price. Figure 5.1 therefore describes the demand and supply curves in this market. The steps reflect the willingness to pay of the buyers and the reservation valuations or prices of the suppliers.

Figure 5.1 The apartment market



Demanders and suppliers are ranked in order of the value they place on an apartment. The market equilibrium is where the marginal demand value of Evan equals the marginal supply value of Kirin at \$500. Five apartments are rented in equilibrium.

In this example, the equilibrium price for apartments will be \$500. Let us see why. At that price the value placed on the marginal unit supplied by Kirin equals Evan's willingness to pay. Five apartments will be rented. A sixth apartment will not be rented because Lynn will let her apartment only if the price reaches \$550. But the sixth potential demander is willing to pay only \$400. Note that, as usual, there is just a single price in the market. Each renter pays \$500, and therefore each supplier also receives \$500.

The consumer and supplier surpluses can now be computed. Note that, while Don is willing to pay \$600, he actually pays \$500. His consumer surplus is therefore \$100. In Figure 5.1, we can see that each consumer's surplus is the distance between the market price and the individual's valuation. These values are given in the final column of the top half of Table 5.1.

Consumer surplus is the excess of consumer willingness to pay over the market price.

Using the same reasoning, we can compute each supplier's surplus, which is the excess of the amount obtained for the rented apartment over the reservation price. For example, Heward obtains a surplus on the supply side of \$150, while Jeff gets \$50. Heward is willing to put his apartment on the market for \$350, but gets the equilibrium price/rent of \$500 for it. Hence his surplus is \$150.

Supplier or producer surplus is the excess of market price over the reservation price of the supplier.

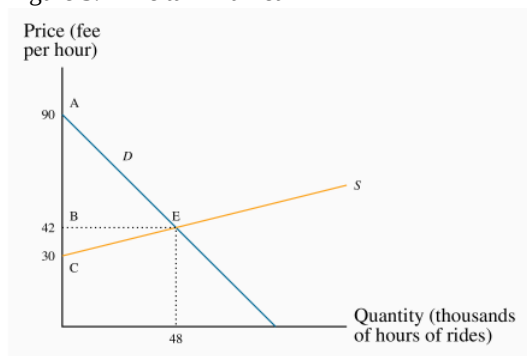
It should now be clear why these measures are called surpluses. *The suppliers and demanders are all willing to participate in this market because they earn this surplus.* It is a measure of their gain from being involved in the trading. The sum of each participant's surplus in the final column of Table 5.1 defines the total surplus in the market. Hence, on the demand side a total surplus arises of \$1,000 and on the supply side a value of \$500.

The taxi market

We do not normally think of demand and supply functions in terms of the steps illustrated in Figure 5.1. Usually there are so many participants in the market that the differences in reservation prices on the supply side and willingness to pay on the demand side are exceedingly small, and so the demand and supply curves are drawn as continuous lines. So our second example reflects this, and comes from the market for taxi rides. We might think of this as an *Uber-* or *Lyft-*type taxi operation.

Let us suppose that the demand and supply curves for taxi rides in a given city are given by the functions in Figure 5.2.

Figure 5.2 The taxi market



Consumer surplus is the area ABE, supplier surplus is the area BCE.

The demand curve represents the willingness to pay on the part of riders. The supply curve represents the willingness to supply on the part of drivers. The price per hour of rides defines the vertical axis; hours of rides (in thousands) are measured on the horizontal axis. The demand intercept of \$90 says that the person who values the ride most highly is willing to pay \$90 per hour. The downward slope of the demand curve states that other buyers are willing to pay less. On the supply side no driver is willing to supply his time and vehicle unless he obtains at least \$30 per hour. To induce additional suppliers a higher price must be paid, and this is represented by the upward sloping supply curve.

The intersection occurs at a price of \$42 per hour and the equilibrium number of ride-hours supplied is 48 thousand¹. Computing the surpluses is very straightforward. By definition the consumer surplus is the excess of the willingness to pay by each buyer above the uniform price. Buyers who value the ride most highly obtain the biggest surplus – the highest valuation rider gets a surplus of \$48 per hour – the difference between his willingness to pay of \$90 and the actual price of \$42. Each successive rider gets a slightly lower surplus until the final rider, who obtains zero. She pays \$42 and values the ride hours at \$42 also. On the supply side, the drivers who are willing to supply rides at the lowest reservation price (\$30 and above) obtain the biggest surplus. The 'marginal' supplier gets no surplus, because the price equals her reservation price.

From this discussion it follows that the consumer surplus is given by the area ABE and the supplier surplus by the area CBE. These are two triangular areas, and measured as half of the base by the perpendicular height. Therefore, in thousands of units:

Consumer Surplus	= (demand value - price) = area ABE
	= $(1/2) \times 48 \times \$48 = \$1,152$
Producer Surplus	= (price - reservation supply value) = area BEC
	= $(1/2) \times 48 \times \$12 = \288

The total surplus that arises in the market is the sum of producer and consumer surpluses, and since the units are in thousands of hours the total surplus here is $(\$1,152 + \$288) \times 1,000 = \$1,440,000$.

5.3 Efficient market outcomes

The definition and measurement of the surplus is straightforward provided the supply and demand functions are known. An important characteristic of the marketplace is that in certain circumstances it produces what we call an efficient outcome, or an efficient market. Such an outcome yields the highest possible sum of surpluses.

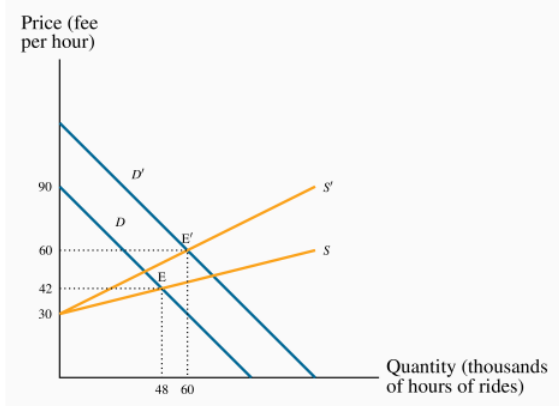
An **efficient market** maximizes the sum of producer and consumer surpluses.

To see that this outcome achieves the goal of maximizing the total surplus, consider what would happen if the quantity $Q=48$ in the taxi example were not supplied. Suppose that the city's taxi czar decreed that 50 units should be supplied, and the czar forced additional drivers on the road. If 2 additional units are to be traded in the market, consider the value of this at the margin. Suppliers value the supply more highly than the buyers are willing to pay. So on these additional 2 units negative surplus would accrue, thus reducing the total.

A second characteristic of the market equilibrium is that potential buyers who would like a cheaper ride and drivers who would like a higher hourly payment do not participate in the market. On the demand side those individuals who are unwilling to pay \$42/hour can take public transit, and on the supply side the those drivers who are unwilling to supply at \$42/hour can allocate their time to alternative activities. Obviously, only those who participate in the market benefit from a surplus.

One final characteristic of surplus measurement should be emphasized. That is, the surplus number is not unique, it depends upon the economic environment. We can illustrate this easily using the taxi example. A well recognized feature of *Uber* taxi rides is that the price varies with road and weather conditions. Poor weather conditions mean that there is an increased demand, and poor road or weather conditions mean that drivers are less willing to supply their services – their reservation payment increases. This situation is illustrated in Figure 5.3.

Figure 5.3 The taxi market



The curves represented by D' and S' represent the curves for bad weather: Taxi rides are more highly valued on the demand side, and drivers must be paid more to supply in less favourable work conditions.

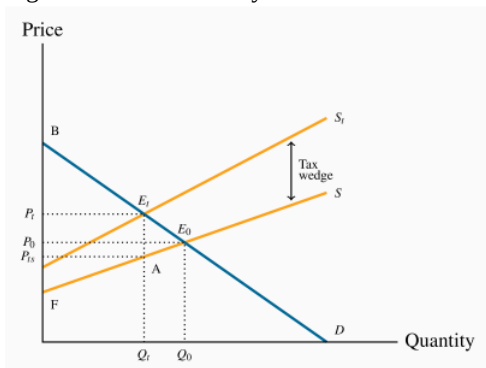
The demand curve has shifted upwards and the supply curve has also changed in such a way that any quantity will now be supplied at a higher price. The new equilibrium is given by E' rather than E .² There is a new equilibrium price-quantity combination that is *efficient in the new market conditions*. This illustrates that there is no such thing as a unique unchanging efficient outcome. When economic factors that influence the buyers' valuations (demand) or the suppliers' reservation prices (supply) change, then the efficient market outcome must be recomputed.

5.4 Taxation, surplus and efficiency

Despite enormous public interest in taxation and its impact on the economy, it is one of the least understood areas of public policy. In this section we will show how an understanding of two fundamental tools of analysis – elasticities and economic surplus – provides powerful insights into the field of taxation.

We begin with the simplest of cases: The federal government's goods and services tax (GST) or the provincial governments' sales taxes (PST). These taxes combined vary by province, but we suppose that a typical rate is 13 percent. In some provinces these two taxes are harmonized. Note that this is a *percentage*, or *ad valorem*, tax, not a *specific* tax of so many dollars per unit traded.

Figure 5.4 The efficiency cost of taxation



The tax shifts S to S_t and reduces the quantity traded from Q_0 to Q_t . At Q_t the demand value placed on an additional unit exceeds the supply valuation by E_tA . Since the tax keeps output at this lower level, the economy cannot take advantage of the additional potential surplus between Q_t and Q_0 . Excess burden = deadweight loss = AE_tE_0 .

Figure 5.4 illustrates the supply and demand curves for some commodity. In the absence of taxes, the equilibrium E_0 is defined by the combination (P_0, Q_0) .

A 13-percent tax is now imposed, and the new supply curve S_t lies 13 percent above the no-tax supply S . A tax wedge is therefore imposed between the price the consumer must pay and the price that the supplier receives. The new equilibrium is E_t , and the new market price is at P_t . The price received by the supplier is lower than that paid by the buyer by the amount of the tax wedge. The post-tax supply price is denoted by P_{ts} .

There are two *burdens* associated with this tax. The first is the revenue burden, the amount of tax revenue paid by the market participants and received by the government. On each of the Q_t units sold, the government receives the amount $(P_t - P_{ts})$. Therefore,

tax revenue is the amount $P_t E_t A P_t$. As illustrated in Chapter 4, the degree to which the market price P_t rises above the no-tax price P_0 depends on the supply and demand elasticities.

A **tax wedge** is the difference between the consumer and producer prices.

The **revenue burden** is the amount of tax revenue raised by a tax.

The second burden of the tax is called the *excess burden*. The concepts of consumer and producer surpluses help us comprehend this. The effect of the tax has been to reduce consumer surplus by $P_t E_t E_0 P_0$. This is the reduction in the pre-tax surplus given by the triangle $P_0 B E_0$. By the same reasoning, supplier surplus is reduced by the amount $P_0 E_0 A P_t$; prior to the tax it was $P_0 E_0 F$. Consumers and suppliers have therefore seen a reduction in their well-being that is measured by these dollar amounts. Nonetheless, the government has additional revenues amounting to $P_t E_t A P_t$, and this tax imposition therefore represents a *transfer* from the consumers and suppliers in the marketplace to the government. Ultimately, the citizens should benefit from this revenue when it is used by the government, and it is therefore not considered to be a net loss of surplus.

However, there remains a part of the surplus loss that is not transferred, the triangular area $E_t E_0 A$. This component is called the excess burden, for the reason that it represents the component of the economic surplus that is not transferred to the government in the form of tax revenue. It is also called the deadweight loss, DWL.

The **excess burden**, or **deadweight loss**, of a tax is the component of consumer and producer surpluses forming a net loss to the whole economy.

The intuition behind this concept is not difficult. At the output Q_t , the value placed by consumers on the last unit supplied is $P_t (= E_t)$, while the production cost of that last unit is $P_s (= A)$. But the potential surplus $(P_t - P_s)$ associated with producing an additional unit cannot be realized, because the tax dictates that the production equilibrium is at Q_t rather than any higher output. Thus, if output could be increased from Q_t to Q_0 , a surplus of value over cost would be realized on every additional unit equal to the vertical distance between the demand and supply functions D and S . Therefore, the loss associated with the tax is the area $E_t E_0 A$.

In public policy debates, this excess burden is rarely discussed. The reason is that notions of consumer and producer surpluses are not well understood by non-economists, despite the fact that the value of lost surpluses is frequently large. Numerous studies have estimated the excess burden associated with raising an additional dollar from the tax system. They rarely find that the excess burden is less than 25 percent of total expenditure. This is a sobering finding. It tells us that if the government wished to implement a new program by raising additional tax revenue, the benefits of the new program should be 25 percent greater than the amount expended on it!

The impact of taxes and other influences that result in an inefficient use of the economy's resources are frequently called distortions because they necessarily lead the economy away from the efficient output. The magnitude of the excess burden is determined by the elasticities of supply and demand in the markets where taxes are levied. To see this, return to Figure 5.4, and suppose that the demand curve through E_0 were more elastic (with the same supply curve, for simplicity). The post-tax equilibrium E_t would now yield a lower Q_t value and a price between P_t and P_0 . The resulting tax revenue raised and the magnitude of the excess burden would differ because of the new elasticity.

A **distortion** in resource allocation means that production is not at an efficient output.

5.5 Market failures – externalities

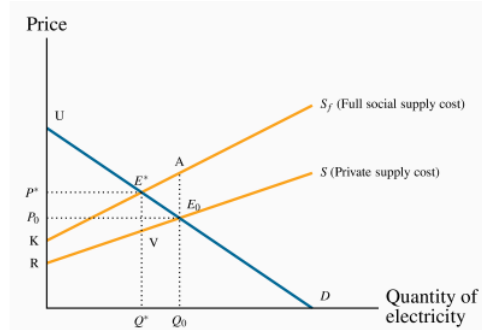
The consumer and producer surplus concepts we have developed are extremely powerful tools of analysis, but the world is not always quite as straightforward as simple models indicate. For example, many suppliers generate pollutants that adversely affect the health of the population, or damage the environment, or both. The term externality is used to denote such impacts. Externalities impact individuals who are not participants in the market in question, and the effects of the externalities may not be captured in the market price. For example, electricity-generating plants that use coal reduce air quality, which, in turn, adversely impacts individuals who suffer from asthma or other lung ailments. While this is an example of a negative externality, externalities can also be positive.

An **externality** is a benefit or cost falling on people other than those involved in the activity's market. It can create a difference between private costs or values and social costs or values.

We will now show why markets characterized by externalities are not efficient, and also show how these externalities might be corrected or reduced. The essence of an externality is that it creates a divergence between private costs/benefits and social

costs/benefits. If a steel producer pollutes the air, and the steel buyer pays only the costs incurred by the producer, then the buyer is not paying the full "social" cost of the product. The problem is illustrated in Figure 5.5.

Figure 5.5 Negative externalities and inefficiency



A negative externality is associated with this good. S reflects private costs, whereas S_f reflects the full social cost. The socially optimal output is Q^* , not the market outcome Q_0 . Beyond Q^* the real cost exceeds the demand value; therefore Q_0 is not an efficient output. A tax that increases P to P^* and reduces output is one solution to the externality.

Negative externalities

In Figure 5.5, the supply curve S represents the cost to the supplier, whereas S_f (the *full cost*) reflects, in addition, the cost of bad air to the population. Of course, we are assuming that this external cost is ascertainable, in order to be able to characterize S_f accurately. Note also that this illustration assumes that, as power output increases, the external cost *per unit* rises, because the difference between the two supply curves increases with output. This implies that low levels of pollution do less damage per unit: Perhaps the population has a natural tolerance for low levels, but higher levels cannot be tolerated easily and so the cost per unit is greater.

Despite the externality, *an efficient level of production can still be defined*. It is given by Q^* , not Q_0 . To see why, consider the impact of reducing output by one unit from Q_0 . At Q_0 the willingness of buyers to pay for the marginal unit supplied is E_0 . The (private) supply cost is also E_0 . But from a societal standpoint there is a pollution/health cost of AE_0 associated with that unit of production. The full cost, as represented by S_f , exceeds the buyer's valuation. Accordingly, if the last unit of output produced is cut, society gains by the amount AE_0 , because the cut in output reduces the excess of true cost over value.

Applying this logic to each unit of output between Q_0 and Q^* , it is evident that society can increase its well-being by the dollar amount equal to the area E^*AE_0 , as a result of reducing production.

Next, consider the consequences of reducing output further from Q^* . Note that some pollution is being created here, and environmentalists frequently advocate that pollution should be reduced to zero. However, an efficient outcome may not involve a zero level of pollution! If the production of power were reduced below Q^* , the loss in value to buyers, as a result of not being able to purchase the good, would exceed the full cost of its production.

If the government decreed that, instead of producing Q^* , no pollution would be tolerated, then society would forgo the possibility of earning the total real surplus equal to the area UE^*K . Economists do not advocate such a zero-pollution policy; rather, we advocate a policy that permits a "tolerable" pollution level – one that still results in net benefits to society. In this particular example, the total cost of the tolerated pollution equals the area between the private and full supply functions, KE^*VR .

As a matter of policy, how is this market influenced to produce the amount Q^* rather than Q_0 ? One option would be for the government to intervene directly with production quotas for each firm. An alternative would be to impose a corrective tax on the good whose production causes the externality: With an appropriate increase in the price, consumers will demand a reduced quantity. In Figure 5.5 a tax equal to the dollar value VE^* would shift the supply curve upward by that amount and result in the quantity Q^* being traded.

A **corrective tax** seeks to direct the market towards a more efficient output.

We are now venturing into the field of environmental policy, where a corrective tax is usually called a carbon tax, and this is explored in the following section. The key conclusion of the foregoing analysis is that an efficient working of the market continues to have meaning in the presence of externalities. An efficient output level still maximizes economic surplus where surplus is correctly defined.

Positive externalities

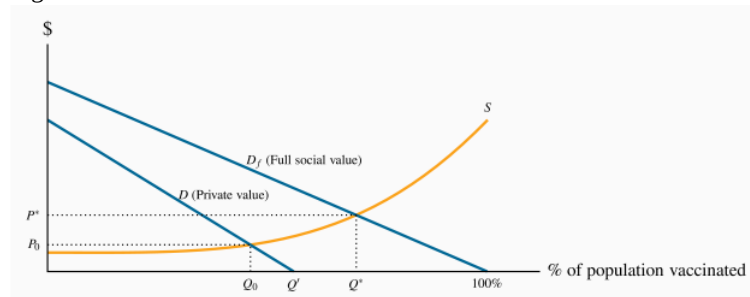
Externalities of the *positive* kind enable individuals or producers to get a type of 'free ride' on the efforts of others. Real world examples abound: When a large segment of the population is immunized against disease, the remaining individuals benefit on account of the reduced probability of transmission.

A less well recognized example is the benefit derived by many producers world-wide from research and development (R&D) undertaken in advanced economies and in universities and research institutes. The result is that society at large, including the corporate sector, gain from this enhanced understanding of science, the environment, or social behaviours.

The free market may not cope any better with these positive externalities than it does with negative externalities, and government intervention may be beneficial. Furthermore, firms that invest heavily in research and development would not undertake such investment if competitors could have a complete free ride and appropriate the fruits. This is why *patent laws* exist, as we shall see later in discussing Canada's competition policy. These laws prevent competitors from copying the product development of firms that invest in R&D. If such protection were not in place, firms would not allocate sufficient resources to R&D, which is a real engine of economic growth. In essence, the economy's research-directed resources would not be appropriately rewarded, and thus too little research would take place.

While patent protection is one form of corrective action, subsidies are another. We illustrated above that an appropriately formulated tax on a good that creates negative externalities can reduce demand for that good, and thereby reduce pollution. A subsidy can be thought of as a negative tax, and can stimulate the supply of goods and services that have positive externalities. Consider the example in Figure 5.6.

Figure 5.6 Positive externalities – the market for flu shots



The value to society of vaccinations exceeds the value to individuals: The greater the number of individuals vaccinated, the lower is the probability of others contracting the virus. D_f reflects this additional value. Consequently, the social optimum is Q^* which exceeds Q_0 .

Individuals have a demand for flu shots given by D . This reflects their private valuation – their personal willingness to pay. But the social value of flu shots is greater. When more individuals are vaccinated, the probability that others will be infected falls. Additionally, with higher rates of immunization, the health system will incur fewer costs in treating the infected. Therefore, the value to society of any quantity of flu shots is greater than the sum of the values that individuals place on them.

D_f reflects the full social value of any quantity of flu shots. In this instance the quantity axis measures the percentage of the population vaccinated, which has a maximum of 100%. If S is the supply curve, the socially optimal, efficient, market outcome is Q^* . The steeply upward-sloping section of S denotes that it may be very costly to vaccinate every last person – particularly those living in outlying communities. How can we influence the market to move from Q_0 towards Q^* ? One solution is a subsidy that would reduce the price to zero. In this case that gets us almost to the optimum, because the percentage of the population now choosing to be vaccinated is given by Q' . The zero price essentially makes the supply curve, as perceived by the population, to be running along the horizontal axis.

Note the social value of the improvement in moving from Q_0 to Q' ; the social value exceeds the social cost. But even at Q' further gains are available because at Q' the social value of additional vaccinations is greater than the social cost. Overall, at the point Q^* the social value is given by the area under the demand curve, and the social cost by the area under the supply curve.

5.6 Other market failures

There are other ways in which markets can fail to reflect accurately the social value or social cost of economic activity. Profit-seeking monopolies, which restrict output in order to increase profits, represent inefficient markets, and we will see why in the chapter on monopoly. Or the market may not deal very well with what are called public goods. These are goods, like radio and

television service, national defence, or health information: With such goods and services many individuals can be supplied with the same good at the same total cost as one individual. We will address this problem in our chapter on government. And, of course, there are international externalities that cannot be corrected by national governments because the interests of adjoining states may differ: One economy may wish to see cheap coal-based electricity being supplied to its consumers, even if this means acid rain or reduced air quality in a neighbouring state. Markets may fail to supply an "efficient" amount of a good or service in all of these situations. Global warming is perhaps the best, and most extreme, example of international externalities and market failure.

5.7 Environmental policy and climate change

Greenhouse gases

The greatest externality challenge in the modern world is to control our emissions of greenhouse gases. The emission of greenhouse gases (GHGs) is associated with a wide variety of economic activities such as coal-based power generation, oil-burning motors, wood-burning stoves, ruminant animals, *etc.* The most common GHG is carbon dioxide, methane is another. The gases, upon emission, circulate in the earth's atmosphere and, following an excessive build-up, prevent sufficient radiant heat from escaping. The result is a slow warming of the earth's surface and air temperatures. It is envisaged that such temperature increases will, in the long term, increase water temperatures and cause glacial melting, with the result that water levels worldwide will rise. In addition to the higher water levels, which the Intergovernmental Panel on Climate Change (IPCC) estimates will be between one foot and one metre by the end of the 21st century, oceans will become more acidic, weather patterns will change and weather events become more variable and severe. The changes will be latitude-specific and vary by economy and continent, and ultimately will impact the agricultural production abilities of certain economies.

Greenhouse gases that accumulate excessively in the earth's atmosphere prevent heat from escaping and lead to **global warming**.

While most scientific findings and predictions are subject to a degree of uncertainty, there is little disagreement in the scientific community on the long-term impact of increasing GHGs in the atmosphere. There is some skepticism as to whether the generally higher temperatures experienced in recent decades are completely attributable to anthropogenic activity since the industrial revolution, or whether they also reflect a natural cycle in the earth's temperature. But scientists agree that a continuance of the recent rate of GHG emissions is leading to serious climatic problems.

The major economic environmental challenge facing the world economy is this: Historically, GHG emissions have been strongly correlated with economic growth. The very high rate of economic growth in many large-population economies such as China and India that will be necessary to raise hundreds of millions out of poverty means that that historical pattern needs to be broken – GHG accumulation must be "decoupled" from economic growth.

GHGs as a common property

A critical characteristic of GHGs is that they are what we call in economics a 'common property': Every citizen in the world 'owns' them, every citizen has equal access to them, and it matters little where these GHGs originate. Consequently, if economy A reduces its GHG emissions, economy B may simply increase its emissions rather than incur the cost of reducing them. Hence, economy A's behaviour goes unrewarded. This is the crux of international agreements – or disagreements. Since GHGs are a common property, in order for A to have the incentive to reduce emissions, it needs to know that B will act correspondingly.

From the Kyoto Protocol to the Paris Accord

The world's first major response to climate concerns came in the form of the United Nations–sponsored Earth Summit in Rio de Janeiro in 1992. This was followed by the signing of the Kyoto Protocol in 1997, in which a group of countries committed themselves to reducing their GHG emissions relative to their 1990 emissions levels by the year 2012. Canada's Parliament subsequently ratified the Kyoto Protocol, and thereby agreed to meet Canada's target of a 6 percent reduction in GHGs relative to the amount emitted in 1990.

On a per-capita basis, Canada is one of the world's largest contributors to global warming, even though Canada's percentage of the total is just 2 percent. Many of the world's major economies refrained from signing the Protocol—most notably China, the United States, and India. Canada's emissions in 1990 amounted to approximately 600 giga tonnes (Gt) of carbon dioxide; but by the time we ratified the treaty in 2002, emissions were 25% above that level. Hence the signing was somewhat meaningless, in that Canada had virtually a zero possibility of attaining its target.

The target date of 2012 has come and gone and subsequent conferences in Copenhagen and Rio failed to yield an international agreement. But in Paris, December 2015, 195 economies committed to reduce their GHG emissions by specific amounts. Canada

was a party to that agreement. Target reductions varied by country. Canada committed itself to reduce GHG emissions by 30% by the year 2030 relative to 2005 emissions levels. To this end the Liberal government of Prime Minister Justin Trudeau announced in late 2016 that if individual Canadian provinces failed to implement a carbon tax, or equivalent, the federal government would impose one unilaterally. The program involves a carbon tax of \$10 per tonne in 2018, that increases by \$10 per annum until it attains a value of \$50 in 2022. Some provinces already have GHG limitation systems in place (cap and trade systems - developed below), and these provinces would not be subject to the federal carbon tax provided the province-level limitation is equivalent to the federal carbon tax.

Canada's GHG emissions

An excellent summary source of data on Canada's emissions and performance during the period 1990-2018 is available on Environment Canada's web site. See:

www.canada.ca/en/environment-climate-change/services/climate-change/greenhouse-gas-emissions/sources-sinks-executive-summary-2020.html#toc3

Canada, like many economies, has become more efficient in its use of energy (the main source of GHGs) in recent decades—its *use of energy per unit of total output* has declined steadily. Canada emitted 0.44 mega tonnes of CO₂ equivalent per billion dollars of GDP in 2005, and 0.36 mega tonnes in 2017. On a *per capita* basis Canada's emissions amounted to 22.9 tonnes in 2005, and dropped to 19.5 by 2017. This modest improvement in efficiency means that Canada's GDP is now *less energy intensive*. The critical challenge is to produce more output while using not just less energy per unit of output, but to use less energy in total.

While Canada's energy intensity (GHGs per unit of output) has dropped, overall emissions have increased by almost 20% since 1990. Furthermore, while developed economies have increased their efficiency, it is the *world's* efficiency that is ultimately critical. By outsourcing much of its manufacturing sector to China, Canada and the West have offloaded some of their most GHG-intensive activities. But GHGs are a common property resource.

Canada's GHG emissions also have a regional aspect: The *production* of oil and gas, which has created considerable wealth for all Canadians, is both energy intensive and concentrated in a limited number of provinces (Alberta, Saskatchewan and more recently Newfoundland and Labrador).

GHG measurement

GHG atmospheric concentrations are measured in parts per million (ppm). Current levels in the atmosphere are slightly above 400 ppm, and continued growth in concentration will lead to serious economic and social disruption. In the immediate pre-industrial revolution era concentrations were in the 280 ppm range. Hence, our world seems to be headed towards a doubling of GHG concentrations in the coming decades.

GHGs are augmented by the annual additions to the stock already in the atmosphere, and at the same time they decay—though very slowly. GHG-reduction strategies that propose an immediate reduction in emissions are more costly than those aimed at a more gradual reduction. For example, a slower investment strategy would permit in-place production and transportation equipment to reach the end of its economic life rather than be scrapped and replaced 'prematurely'. Policies that focus upon longer-term replacement are therefore less costly in this specific sense.

While not all economists and policy makers agree on the time scale for attacking the problem, the longer that GHG reduction is postponed, the greater the efforts will have to be in the long term—because GHGs will build up more rapidly in the near term.

A critical question in controlling GHG emissions relates to the cost of their control: How much of annual growth might need to be sacrificed in order to get emissions onto a sustainable path? Again estimates vary. The Stern Review (2006) proposed that, with an increase in technological capabilities, a strategy that focuses on the relative near-term implementation of GHG reduction measures might cost "only" a few percentage points of the value of world output. If correct, this is a low price to pay for risk avoidance in the longer term.

Nonetheless, such a reduction will require particular economic policies, and specific sectors will be impacted more than others.

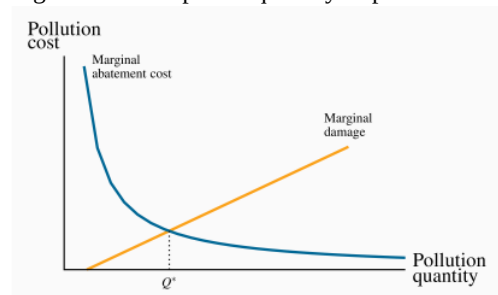
Economic policies for climate change

There are three main ways in which polluters can be controlled. One involves issuing direct controls; the other two involve incentives—in the form of pollution taxes, or on tradable "permits" to pollute.

To see how these different policies operate, consider first Figure 5.7. It is a standard diagram in environmental economics, and is somewhat similar to our supply and demand curves. On the horizontal axis is measured the quantity of environmental damage or pollution, and on the vertical axis its dollar value or cost. The upward-sloping damage curve represents the cost to society of each additional unit of pollution or gas, and it is therefore called a marginal damage curve. It is positively sloped to reflect the reality that, at low levels of emissions, the damage of one more unit is less than at higher levels. In terms of our earlier discussion, this means that an increase in GHGs of 10 ppm when concentrations are at 300 ppm may be less damaging than a corresponding increase when concentrations are at 500 ppm.

The **marginal damage curve** reflects the cost to society of an additional unit of pollution.

Figure 5.7 The optimal quantity of pollution



Q^* represents the optimal amount of pollution. More than this would involve additional social costs because damages exceed abatement costs. Conversely, less than Q^* would require an abatement cost that exceeds the reduction in damage.

The second curve is the abatement curve. It reflects the cost of reducing emissions by one unit, and is therefore called a marginal abatement curve. This curve has a negative slope indicating that, as we reduce the total quantity of pollution produced (moving towards the origin on the horizontal axis), the cost of further unit reductions rises. This shape corresponds to reality. For example, halving the emissions of pollutants and gases from automobiles may be achieved by adding a catalytic converter and reducing the amount of lead in gasoline. But reducing those emissions all the way to zero requires the development of major new technologies such as electric cars—an enormously more costly undertaking.

The **marginal abatement curve** reflects the cost to society of reducing the quantity of pollution by one unit.

If producers are unconstrained in the amount of pollution they produce, they will produce more than what we will show is the optimal amount – corresponding to Q^* . This amount is optimal in the sense that at levels greater than Q^* the damage exceeds the cost of reducing the emissions. However, reducing emissions below Q^* would mean incurring a cost per unit reduction that exceeds the benefit of that reduction. Another way of illustrating this is to observe that at a level of pollution above Q^* the cost of reducing it is less than the damage it inflicts, and therefore a net gain accrues to society as a result of the reduction. But to reduce pollution below Q^* would involve an abatement cost greater than the reduction in pollution damage and therefore no net gain to society. This constitutes a first rule in optimal pollution policy.

An optimal quantity of pollution occurs when the marginal cost of abatement equals the marginal damage.

A second guiding principle emerges by considering a situation in which some firms are relatively 'clean' and others are 'dirty'. More specifically, a clean firm A may have already invested in new equipment that uses less energy per unit of output produced, or emits fewer pollutants per unit of output. In contrast, the dirty firm B uses older dirtier technology. Suppose furthermore that these two firms form a particular sector of the economy and that the government sets a limit on total pollution from this sector, and that this limit is less than what the two firms are currently producing. What is the least costly method to meet the target?

The intuitive answer to this question goes as follows: In order to reduce pollution at least cost to the sector, calculate what it would cost each firm to reduce pollution from its present level. Then implement a system so that the firm with the least cost of reduction is the first to act. In this case the 'dirty' firm will likely have a lower cost of abatement since it has not yet upgraded its physical plant. This leads to a second rule in pollution policy:

With many polluters, the least cost policy to society requires producers with the lowest abatement costs to act first.

This principle implies that policies which impose the same emission limits on firms may not be the least costly manner of achieving a target level of pollution. Let us now consider the use of tradable permits and corrective/carbon taxes as policy instruments. These are market-based systems aimed at reducing GHGs.

Tradable permits and **corrective/carbon taxes** are market-based systems aimed at reducing GHGs.

Incentive mechanism I: Tradable permits

A system of tradable permits is frequently called a 'cap and trade' system, because it limits or caps the total permissible emissions, while at the same time allows a market to develop in permits. For illustrative purposes, consider the hypothetical two-firm sector we developed above, composed of firms A and B. Firm A has invested in clean technology, firm B has not. Thus it is less costly for B to reduce emissions than A if further reductions are required. Next suppose that each firm is allocated by the government a specific number of 'GHG emission permits'; and that the total of such permits is less than the amount of emissions at present, and that each firm is emitting more than its permits allow. How can these firms achieve the target set for this sector of the economy?

The answer is that they should be able to engage in mutually beneficial trade: If firm B has a lower cost of reducing emissions than A, then it may be in A's interest to pay B to reduce B's emissions heavily. Imagine that each firm is emitting 60 units of GHG, but they have permits to emit only 50 units each. And furthermore suppose it costs B \$20 to reduce GHGs by one unit, whereas it costs A \$30 to do this. In this situation A could pay B \$25 for several permits and this would benefit both firms. B can reduce GHGs at a cost of \$20 and is being paid \$25 to do this. In turn A would incur a cost of \$30 per unit to reduce his GHGs but he can buy permits from B for just \$25 and avoid the \$30 cost. Both firms gain, and the total cost to the economy is lower than if each firm had to reduce by the same amount.

The benefit of the cap 'n trade system is that it enables the marketplace to reduce GHGs at least cost.

The largest system of tradable permits currently operates in the European Union: The *EU Emissions Trading System*. It covers more than 10,000 large energy-using installations. Trading began in 2005. In North America a number of Western states and several Canadian provinces are joined, either as participants or observers, in the *Western Climate Initiative*, which is committed to reduce GHGs by means of tradable emissions permits. The longer-term goal of these systems is for the government to issue progressively fewer permits each year, and to include an ever larger share of GHG-emitting enterprises with the passage of time.

Policy in practice – international

In an ideal world, permits would be traded internationally, and such a system might be of benefit to developing economies: If the cost of reducing pollution is relatively low in developing economies because they have few controls in place, then developed economies, for whom the cost of GHG reduction is high, could induce firms in the developing world to undertake cost reductions. Such a trade would be mutually beneficial. For example, imagine in the above example that B is located in the developing world and A in the developed world. Both would obviously gain from such an arrangement, and because GHGs are a common property, the source of GHGs from a damage standpoint is immaterial.

Incentive mechanism II: Taxes

Corrective taxes are frequently called *Pigovian* taxes, after the economist Arthur Pigou. He advocated taxing activities that cause negative externalities. These taxes have been examined above in Section 5.4. Corrective taxes of this type can be implemented as part of a tax package *reform*. For example, taxpayers are frequently reluctant to see governments take 'yet more' of their money, in the form of new taxes. Such concerns can be addressed by reducing taxes in other sectors of the economy, in such a way that the package of tax changes maintains a 'revenue neutral' impact.

Revenues from taxes and permits

Taxes and tradable permits differ in that taxes generate revenue for the government from polluting producers, whereas permits may not generate revenue, or may generate less revenue. If the government simply *allocates* permits initially to all polluters, free of charge, and allows a market to develop, such a process generates no revenue to the government. While economists may advocate an *auction* of permits in the start-up phase of a tradable permits market, such a mechanism may run into political objections.

Setting taxes at the appropriate level requires knowledge of the cost and damage functions associated with GHGs. At the present time, economists and environmental scientists think that an appropriate price or tax on one tonne of GHG is in the \$40–\$50 range. Such a tax would reduce emissions to a point where the longer-term impact of GHGs would not be so severe as otherwise.

British Columbia introduced a carbon tax of \$10 per tonne of GHG on fuels in 2008, and has increased that price regularly. This tax was designed to be revenue neutral in order to make it more acceptable. This means that British Columbia reduced its income tax rates by an amount such that income tax payments would fall by an amount equal to the revenue captured by the carbon tax.

GHG policy at the federal level in Canada is embodied in the *Greenhouse Gas Pollution Pricing Act of 2018*. As detailed earlier, the Act imposes a yearly increasing levy on emissions. The system is intended to be revenue neutral, in that the revenues will be

returned to households in the form of a 'paycheck' by the federal government. Large emitters of GHGs are permitted a specific threshold number of tonnes of emission each year without being penalized. Beyond that threshold the above rates apply.

Will this amount of carbon taxation hurt consumers, and will it enable Canada to reach its 2030 GHG goal? As a specific example: the gasoline-pricing rule of thumb is that each \$10 in carbon taxation or pricing leads to an increase in the price of gasoline at the pump of about 2.5 cents. So a \$50 levy per tonne means gas at the pump should rise by 12.5 cents per litre. The proceeds are returned to households.

As for the goal of reaching the 2030 target announced at Paris: Environment Canada estimates that the pricing scheme will reduce GHG emissions by about 60 tonnes per annum. But Canada's goal stated in Paris is to reduce emissions in 2030 by approximately four times this amount. Under the Paris Accord, Canada stated that its 2030 goal would be to reduce emissions by 30% from their 2005 level of 725 MT, that is by an amount equal to approximately 220 tonnes.

Policy in practice – domestic large final emitters

Governments frequently focus upon quantities emitted by individual large firms, or *large final emitters (LFEs)*. In some economies, a relatively small number of producers are responsible for a disproportionate amount of an economy's total pollution, and limits are placed on those firms in the belief that significant economy-wide reductions can be achieved in this manner. One reason for concentrating on these LFEs is that the monitoring costs are relatively small compared to the costs associated with monitoring *all* firms in the economy. It must be kept in mind that pollution permits may be a legal requirement in some jurisdictions, but monitoring is still required, because firms could choose to risk polluting without owning a permit.

Conclusion

Welfare economics lies at the heart of public policy. Demand and supply curves can be interpreted as value curves and cost curves when there are no externalities involved. This is what enables us to define an efficient output of a product, and consequently an efficient use of the economy's resources. While efficiency is a central concept in economics, we must keep in mind that when the economic environment changes so too will the efficient use of resources, as we illustrated in Section 5.7.

In this chapter we have focused on equity issues through the lens of GHG emissions. The build-up of GHGs in our atmosphere invokes the concept of *intergenerational equity*: The current generation is damaging the environment and the costs of that damage will be borne by subsequent generations. Hence it is inequitable in the intergenerational sense for us to leave a negative legacy to succeeding generations. Equity arises within generations also. For example, how much more in taxes should the rich pay relative to the non-rich? We will explore this type of equity in our chapter on government.

Key Terms

Welfare economics assesses how well the economy allocates its scarce resources in accordance with the goals of efficiency and equity.

Efficiency addresses the question of how well the economy's resources are used and allocated.

Equity deals with how society's goods and rewards are, and should be, distributed among its different members, and how the associated costs should be apportioned.

Consumer surplus is the excess of consumer willingness to pay over the market price.

Supplier or producer surplus is the excess of market price over the reservation price of the supplier.

Efficient market: maximizes the sum of producer and consumer surpluses.

Tax wedge is the difference between the consumer and producer prices.

Revenue burden is the amount of tax revenue raised by a tax.

Excess burden of a tax is the component of consumer and producer surpluses forming a net loss to the whole economy.

Deadweight loss of a tax is the component of consumer and producer surpluses forming a net loss to the whole economy.

Distortion in resource allocation means that production is not at an efficient output.

Externality is a benefit or cost falling on people other than those involved in the activity's market. It can create a difference between private costs or values and social costs or values.

Corrective tax seeks to direct the market towards a more efficient output.

Greenhouse gases that accumulate excessively in the earth's atmosphere prevent heat from escaping and lead to global warming.

Marginal damage curve reflects the cost to society of an additional unit of pollution.

Marginal abatement curve reflects the cost to society of reducing the quantity of pollution by one unit.

Tradable permits are a market-based system aimed at reducing GHGs.

Carbon taxes are a market-based system aimed at reducing GHGs.

Exercises for Chapter 5

EXERCISE 5.1

Four teenagers live on your street. Each is willing to shovel snow from one driveway each day. Their "willingness to shovel" valuations (supply) are: Jean, \$10; Kevin, \$9; Liam, \$7; Margaret, \$5. Several households are interested in having their driveways shoveled, and their willingness to pay values (demand) are: Jones, \$8; Kirpinsky, \$4; Lafleur, \$7.50; Murray, \$6.

1. Draw the implied supply and demand curves as step functions.
2. How many driveways will be shoveled in equilibrium?
3. Compute the maximum possible sum for the consumer and supplier surpluses.
4. If a new (wealthy) family arrives on the block, that is willing to pay \$12 to have their driveway cleared, recompute the answers to parts (a), (b), and (c).

EXERCISE 5.2

Consider a market where supply curve is horizontal at $P=10$ and the demand curve has intercepts $\{34, 34\}$, and is defined by the relation $P=34-Q$.

1. Illustrate the market geometrically.
2. Impose a tax of \$2 per unit on the good so that the supply curve is now $P=12$. Illustrate the new equilibrium quantity.
3. Illustrate in your diagram the tax revenue generated.
4. Illustrate the deadweight loss of the tax.

EXERCISE 5.3

Next, consider an example of DWL in the labour market. Suppose the demand for labour is given by the fixed gross wage $W = \$16$. The supply is given by $W=0.8L$, indicating that the supply curve goes through the origin with a slope of 0.8.

1. Illustrate the market geometrically.
2. Calculate the supplier surplus, knowing that the equilibrium is $L=20$.
3. *Optional:* Suppose a wage tax is imposed that produces a net-of-tax wage equal to $W = \$12$. This can be seen as a downward shift in the demand curve. Illustrate the new quantity supplied and the new supplier's surplus.

EXERCISE 5.4

Governments are in the business of providing information to potential buyers. The first serious provision of information on the health consequences of tobacco use appeared in the United States Report of the Surgeon General in 1964.

1. How would you represent this intervention in a supply and demand for tobacco diagram?
2. Did this intervention "correct" the existing market demand?

EXERCISE 5.5

In deciding to drive a car in the rush hour, you think about the cost of gas and the time of the trip.

1. Do you slow down other people by driving?
2. Is this an externality, given that you yourself are suffering from slow traffic?

EXERCISE 5.6

Suppose that our local power station burns coal to generate electricity. The demand and supply functions for electricity are given by $P=12-0.5Q$ and $P=2+0.5Q$, respectively. The demand curve has intercepts $\{12, 24\}$ and the supply curve intercept is at \$2 with a

slope of one half. However, for each unit of electricity generated, there is an externality. When we factor this into the supply side of the market, the real social cost is increased by \$1 per unit. That is, the supply curve shifts upwards by \$1, and now takes the form $P=3+0.5Q$.

1. Illustrate the free-market equilibrium.
2. Illustrate the efficient (i.e. socially optimal) level of production.

EXERCISE 5.7

Your local dry cleaner, Bleached Brite, is willing to launder shirts at its cost of \$1.00 per shirt. The neighbourhood demand for this service is $P=5-0.005Q$, knowing that the demand intercepts are $\{ \$5, 1000 \}$.

1. Illustrate the market equilibrium.
2. Suppose that, for each shirt, Bleached Brite emits chemicals into the local environment that cause \$0.25 damage per shirt. This means the full cost of each shirt is \$1.25. Illustrate graphically the socially optimal number of shirts to be cleaned.
3. *Optional:* Calculate the socially optimal number of shirts to be cleaned.

EXERCISE 5.8

The supply curve for agricultural labour is given by $W=6+0.1L$, where W is the wage (price per unit) and L the quantity traded. Employers are willing to pay a wage of \$12 to all workers who are willing to work at that wage; hence the demand curve is $W=12$.

1. Illustrate the market equilibrium, if you are told that the equilibrium occurs where $L=60$.
2. Compute the supplier surplus at this equilibrium.

EXERCISE 5.9

Optional: The market demand for vaccine XYZ is given by $P=36-Q$ and the supply conditions are $P=20$; so \$20 represents the true cost of supplying a unit of vaccine. There is a positive externality associated with being vaccinated, and the real societal value is known and given by $P=36-(1/2)Q$. This new demand curve represents the true value to society of each vaccination. This is reflected in the private value demand curve rotating upward around the price intercept of \$36.

1. Illustrate the private and social demand curves on a diagram, with intercept values calculated.
 2. What is the market solution to this supply and demand problem?
 3. What is the socially optimal number of vaccinations?
1. The demand and supply functions behind these curves are $P=90-1Q$ and $P=30+(1/4)Q$. Equating supply and demand yields the solutions in the text.
 2. For example, if the demand curve shifts upwards, parallel, to become $P=120-1Q$ and the supply curve changes in slope to $P=30+(1/2)Q$, the new equilibrium solution is $\{ P = \$60, Q = 60 \}$.

This page titled [5: Welfare economics and externalities](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.1: Equity and efficiency

In modern mixed economies, markets and governments together determine the output produced and also who benefits from that output. In this chapter we explore a very broad question that forms the core of welfare economics: Even if market forces drive efficiency, are they a good way to allocate scarce resources in view of the fact that they not only give rise to inequality and poverty, but also fail to capture the impacts of productive activity on non-market participants? Mining impacts the environment, traffic results in road fatalities, alcohol, tobacco and opioids cause premature deaths. These products all generate secondary impacts beyond their stated objective. We frequently call these *external* effects.

The analysis of markets in this larger sense involves not just economic efficiency; public policy additionally has a normative content because policies can impact the various participants in different ways and to different degrees. Welfare economics, therefore, deals with both normative and positive issues.

Welfare economics assesses how well the economy allocates its scarce resources in accordance with the goals of efficiency and equity.

Political parties on the left and right disagree on how well a market economy works. Canada's New Democratic Party emphasizes the market's failings and the need for government intervention, while the Progressive Conservative Party believes, broadly, that the market fosters choice, incentives, and efficiency. What lies behind this disagreement? The two principal factors are efficiency and equity. Efficiency addresses the question of how well the economy's resources are used and allocated. In contrast, equity deals with how society's goods and rewards are, and should be, distributed among its different members, and how the associated costs should be apportioned.

Equity deals with how society's goods and rewards are, and should be, distributed among its different members, and how the associated costs should be apportioned.

Efficiency addresses the question of how well the economy's resources are used and allocated.

Equity is also concerned with how different generations share an economy's productive capabilities: More investment today makes for a more productive economy tomorrow, but more greenhouse gases today will reduce environmental quality tomorrow. These are inter-generational questions.

Climate change caused by global warming forms one of the biggest challenges for humankind at the present time. As we shall see in this chapter, economics has much to say about appropriate policies to combat warming. Whether pollution-abatement policies should be implemented today or down the road involves considerations of equity between generations. Our first task is to develop an analytical tool which will prove vital in assessing and computing welfare benefits and costs – economic surplus.

This page titled [5.1: Equity and efficiency](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.2: Consumer and producer surplus

An understanding of economic efficiency is greatly facilitated as a result of understanding two related measures: Consumer surplus and producer surplus. Consumer surplus relates to the demand side of the market, producer surplus to the supply side. Producer surplus is also termed supplier surplus. These measures can be understood with the help of a standard example, the market for city apartments.

The market for apartments

Table 5.1 and Figure 5.1 describe the hypothetical data. We imagine first a series of city-based students who are in the market for a standardized downtown apartment. These individuals are not identical; they value the apartment differently. For example, Alex enjoys comfort and therefore places a higher value on a unit than Brian. Brian, in turn, values it more highly than Cathy or Don. Evan and Frank would prefer to spend their money on entertainment, and so on. These valuations are represented in the middle column of the demand panel in Table 5.1, and also in Figure 5.1 with the highest valuations closest to the origin. The valuations reflect the willingness to pay of each consumer.

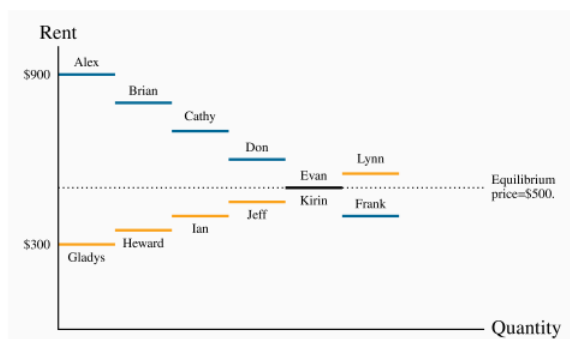
Table 5.1 Consumer and supplier surpluses

Demand		
Individual	Demand valuation	Surplus
Alex	900	400
Brian	800	300
Cathy	700	200
Don	600	100
Evan	500	0
Frank	400	0
Supply		
Individual	Reservation value	Surplus
Gladys	300	200
Heward	350	150
Ian	400	100
Jeff	450	50
Kirin	500	0
Lynn	550	0

On the supply side we imagine the market as being made up of different individuals or owners, who are willing to put their apartments on the market for different prices. Gladys will accept less rent than Heward, who in turn will accept less than Ian. The minimum prices that the suppliers are willing to accept are called *reservation* prices or values, and these are given in the lower part of Table 5.1. Unless the market price is greater than their reservation price, suppliers will hold back.

By definition, as stated in Chapter 3, the demand curve is made up of the valuations placed on the good by the various demanders. Likewise, the reservation values of the suppliers form the supply curve. If Alex is willing to pay \$900, then that is his demand price; if Heward is willing to put his apartment on the market for \$350, he is by definition willing to supply it for that price. Figure 5.1 therefore describes the demand and supply curves in this market. The steps reflect the willingness to pay of the buyers and the reservation valuations or prices of the suppliers.

Figure 5.1 The apartment market



Demanders and suppliers are ranked in order of the value they place on an apartment. The market equilibrium is where the marginal demand value of Evan equals the marginal supply value of Kirin at \$500. Five apartments are rented in equilibrium.

In this example, the equilibrium price for apartments will be \$500. Let us see why. At that price the value placed on the marginal unit supplied by Kirin equals Evan's willingness to pay. Five apartments will be rented. A sixth apartment will not be rented because Lynn will let her apartment only if the price reaches \$550. But the sixth potential demander is willing to pay only \$400. Note that, as usual, there is just a single price in the market. Each renter pays \$500, and therefore each supplier also receives \$500.

The consumer and supplier surpluses can now be computed. Note that, while Don is willing to pay \$600, he actually pays \$500. His consumer surplus is therefore \$100. In Figure 5.1, we can see that each consumer's surplus is the distance between the market price and the individual's valuation. These values are given in the final column of the top half of Table 5.1.

Consumer surplus is the excess of consumer willingness to pay over the market price.

Using the same reasoning, we can compute each supplier's surplus, which is the excess of the amount obtained for the rented apartment over the reservation price. For example, Heward obtains a surplus on the supply side of \$150, while Jeff gets \$50. Heward is willing to put his apartment on the market for \$350, but gets the equilibrium price/rent of \$500 for it. Hence his surplus is \$150.

Supplier or producer surplus is the excess of market price over the reservation price of the supplier.

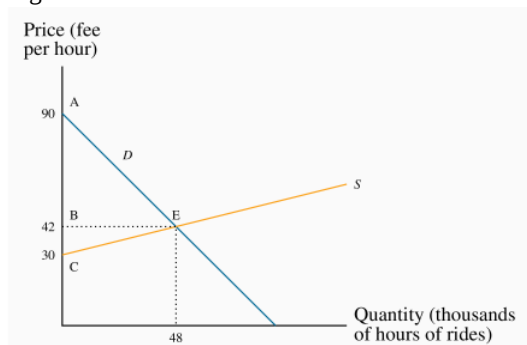
It should now be clear why these measures are called surpluses. *The suppliers and demanders are all willing to participate in this market because they earn this surplus.* It is a measure of their gain from being involved in the trading. The sum of each participant's surplus in the final column of Table 5.1 defines the total surplus in the market. Hence, on the demand side a total surplus arises of \$1,000 and on the supply side a value of \$500.

The taxi market

We do not normally think of demand and supply functions in terms of the steps illustrated in Figure 5.1. Usually there are so many participants in the market that the differences in reservation prices on the supply side and willingness to pay on the demand side are exceedingly small, and so the demand and supply curves are drawn as continuous lines. So our second example reflects this, and comes from the market for taxi rides. We might think of this as an *Uber-* or *Lyft-*type taxi operation.

Let us suppose that the demand and supply curves for taxi rides in a given city are given by the functions in Figure 5.2.

Figure 5.2 The taxi market



Consumer surplus is the area ABE, supplier surplus is the area BCE.

The demand curve represents the willingness to pay on the part of riders. The supply curve represents the willingness to supply on the part of drivers. The price per hour of rides defines the vertical axis; hours of rides (in thousands) are measured on the horizontal axis. The demand intercept of \$90 says that the person who values the ride most highly is willing to pay \$90 per hour. The downward slope of the demand curve states that other buyers are willing to pay less. On the supply side no driver is willing to supply his time and vehicle unless he obtains at least \$30 per hour. To induce additional suppliers a higher price must be paid, and this is represented by the upward sloping supply curve.

The intersection occurs at a price of \$42 per hour and the equilibrium number of ride-hours supplied is 48 thousand¹. Computing the surpluses is very straightforward. By definition the consumer surplus is the excess of the willingness to pay by each buyer above the uniform price. Buyers who value the ride most highly obtain the biggest surplus – the highest valuation rider gets a surplus of \$48 per hour – the difference between his willingness to pay of \$90 and the actual price of \$42. Each successive rider gets a slightly lower surplus until the final rider, who obtains zero. She pays \$42 and values the ride hours at \$42 also. On the supply side, the drivers who are willing to supply rides at the lowest reservation price (\$30 and above) obtain the biggest surplus. The 'marginal' supplier gets no surplus, because the price equals her reservation price.

From this discussion it follows that the consumer surplus is given by the area ABE and the supplier surplus by the area CBE. These are two triangular areas, and measured as half of the base by the perpendicular height. Therefore, in thousands of units:

Consumer Surplus	= (demand value - price) = area ABE
	= $(1/2) \times 48 \times \$48 = \$1,152$
Producer Surplus	= (price - reservation supply value) = area BEC
	= $(1/2) \times 48 \times \$12 = \288

The total surplus that arises in the market is the sum of producer and consumer surpluses, and since the units are in thousands of hours the total surplus here is $(\$1,152 + \$288) \times 1,000 = \$1,440,000$.

This page titled [5.2: Consumer and producer surplus](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.3: Efficient market outcomes

The definition and measurement of the surplus is straightforward provided the supply and demand functions are known. An important characteristic of the marketplace is that in certain circumstances it produces what we call an efficient outcome, or an efficient market. Such an outcome yields the highest possible sum of surpluses.

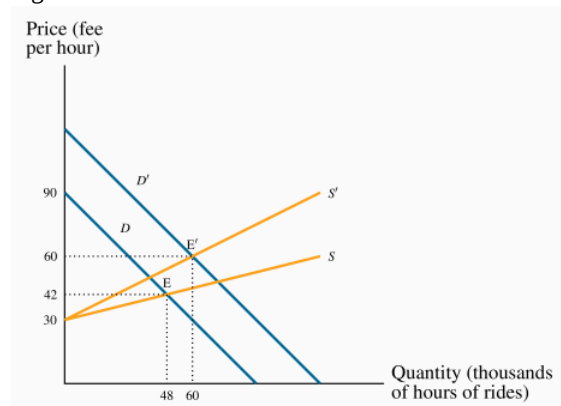
An **efficient market** maximizes the sum of producer and consumer surpluses.

To see that this outcome achieves the goal of maximizing the total surplus, consider what would happen if the quantity $Q=48$ in the taxi example were not supplied. Suppose that the city's taxi czar decreed that 50 units should be supplied, and the czar forced additional drivers on the road. If 2 additional units are to be traded in the market, consider the value of this at the margin. Suppliers value the supply more highly than the buyers are willing to pay. So on these additional 2 units negative surplus would accrue, thus reducing the total.

A second characteristic of the market equilibrium is that potential buyers who would like a cheaper ride and drivers who would like a higher hourly payment do not participate in the market. On the demand side those individuals who are unwilling to pay \$42/hour can take public transit, and on the supply side the those drivers who are unwilling to supply at \$42/hour can allocate their time to alternative activities. Obviously, only those who participate in the market benefit from a surplus.

One final characteristic of surplus measurement should be emphasized. That is, the surplus number is not unique, it depends upon the economic environment. We can illustrate this easily using the taxi example. A well recognized feature of *Uber* taxi rides is that the price varies with road and weather conditions. Poor weather conditions mean that there is an increased demand, and poor road or weather conditions mean that drivers are less willing to supply their services – their reservation payment increases. This situation is illustrated in Figure 5.3.

Figure 5.3 The taxi market



The curves represented by D' and S' represent the curves for bad weather: Taxi rides are more highly valued on the demand side, and drivers must be paid more to supply in less favourable work conditions.

The demand curve has shifted upwards and the supply curve has also changed in such a way that any quantity will now be supplied at a higher price. The new equilibrium is given by E' rather than E . There is a new equilibrium price-quantity combination that is *efficient in the new market conditions*. This illustrates that there is no such thing as a unique unchanging efficient outcome. When economic factors that influence the buyers' valuations (demand) or the suppliers' reservation prices (supply) change, then the efficient market outcome must be recomputed.

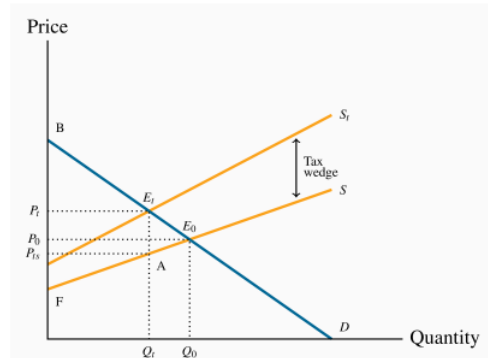
This page titled 5.3: Efficient market outcomes is shared under a CC BY-NC-SA 4.0 license and was authored, remixed, and/or curated by Douglas Curtis and Ian Irvine (Lyryx) via source content that was edited to the style and standards of the LibreTexts platform.

5.4: Taxation, surplus and efficiency

Despite enormous public interest in taxation and its impact on the economy, it is one of the least understood areas of public policy. In this section we will show how an understanding of two fundamental tools of analysis – elasticities and economic surplus – provides powerful insights into the field of taxation.

We begin with the simplest of cases: The federal government's goods and services tax (GST) or the provincial governments' sales taxes (PST). These taxes combined vary by province, but we suppose that a typical rate is 13 percent. In some provinces these two taxes are harmonized. Note that this is a *percentage*, or *ad valorem*, tax, not a *specific* tax of so many dollars per unit traded.

Figure 5.4 The efficiency cost of taxation



The tax shifts S to S_t and reduces the quantity traded from Q_0 to Q_t . At Q_t the demand value placed on an additional unit exceeds the supply valuation by $E_t A$. Since the tax keeps output at this lower level, the economy cannot take advantage of the additional potential surplus between Q_t and Q_0 . Excess burden = deadweight loss = $A E_t E_0$.

Figure 5.4 illustrates the supply and demand curves for some commodity. In the absence of taxes, the equilibrium E_0 is defined by the combination (P_0, Q_0) .

A 13-percent tax is now imposed, and the new supply curve S_t lies 13 percent above the no-tax supply S . A tax wedge is therefore imposed between the price the consumer must pay and the price that the supplier receives. The new equilibrium is E_t , and the new market price is at P_t . The price received by the supplier is lower than that paid by the buyer by the amount of the tax wedge. The post-tax supply price is denoted by P_{ts} .

There are two *burdens* associated with this tax. The first is the revenue burden, the amount of tax revenue paid by the market participants and received by the government. On each of the Q_t units sold, the government receives the amount $(P_t - P_{ts})$. Therefore, tax revenue is the amount $P_t E_t A P_{ts}$. As illustrated in Chapter 4, the degree to which the market price P_t rises above the no-tax price P_0 depends on the supply and demand elasticities.

A **tax wedge** is the difference between the consumer and producer prices.

The **revenue burden** is the amount of tax revenue raised by a tax.

The second burden of the tax is called the *excess burden*. The concepts of consumer and producer surpluses help us comprehend this. The effect of the tax has been to reduce consumer surplus by $P_t E_t E_0 P_0$. This is the reduction in the pre-tax surplus given by the triangle $P_0 B E_t$. By the same reasoning, supplier surplus is reduced by the amount $P_0 E_0 A P_{ts}$; prior to the tax it was $P_0 E_0 F$. Consumers and suppliers have therefore seen a reduction in their well-being that is measured by these dollar amounts. Nonetheless, the government has additional revenues amounting to $P_t E_t A P_{ts}$, and this tax imposition therefore represents a *transfer* from the consumers and suppliers in the marketplace to the government. Ultimately, the citizens should benefit from this revenue when it is used by the government, and it is therefore not considered to be a net loss of surplus.

However, there remains a part of the surplus loss that is not transferred, the triangular area $E_t E_0 A$. This component is called the excess burden, for the reason that it represents the component of the economic surplus that is not transferred to the government in the form of tax revenue. It is also called the deadweight loss, DWL.

The **excess burden**, or **deadweight loss**, of a tax is the component of consumer and producer surpluses forming a net loss to the whole economy.

The intuition behind this concept is not difficult. At the output Q_t , the value placed by consumers on the last unit supplied is $P_t (= E_t)$, while the production cost of that last unit is $P_{ts} (= A)$. But the potential surplus $(P_t - P_{ts})$ associated with producing an additional unit cannot be realized, because the tax dictates that the production equilibrium is at Q_t rather than any higher output. Thus, if output could be increased from Q_t to Q_0 , a surplus of value over cost would be realized on every additional unit equal to the vertical distance between the demand and supply functions D and S . Therefore, the loss associated with the tax is the area $E_t E_0 A$.

In public policy debates, this excess burden is rarely discussed. The reason is that notions of consumer and producer surpluses are not well understood by non-economists, despite the fact that the value of lost surpluses is frequently large. Numerous studies have estimated the excess burden associated with raising an additional dollar from the tax system. They rarely find that the excess burden is less than 25 percent of total expenditure. This is a sobering finding. It tells us that if the government wished to implement a new program by raising additional tax revenue, the benefits of the new program should be 25 percent greater than the amount expended on it!

The impact of taxes and other influences that result in an inefficient use of the economy's resources are frequently called distortions because they necessarily lead the economy away from the efficient output. The magnitude of the excess burden is determined by the elasticities of supply and demand in the markets where taxes are levied. To see this, return to Figure 5.4, and suppose that the demand curve through E_0 were more elastic (with the same supply curve, for simplicity). The post-tax equilibrium E_t would now yield a lower Q_t value and a price between P_t and P_0 . The resulting tax revenue raised and the magnitude of the excess burden would differ because of the new elasticity.

A **distortion** in resource allocation means that production is not at an efficient output.

This page titled [5.4: Taxation, surplus and efficiency](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

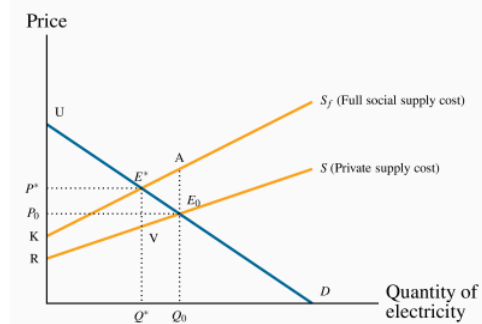
5.5: Market failures - externalities

The consumer and producer surplus concepts we have developed are extremely powerful tools of analysis, but the world is not always quite as straightforward as simple models indicate. For example, many suppliers generate pollutants that adversely affect the health of the population, or damage the environment, or both. The term **externality** is used to denote such impacts. Externalities impact individuals who are not participants in the market in question, and the effects of the externalities may not be captured in the market price. For example, electricity-generating plants that use coal reduce air quality, which, in turn, adversely impacts individuals who suffer from asthma or other lung ailments. While this is an example of a negative externality, externalities can also be positive.

An **externality** is a benefit or cost falling on people other than those involved in the activity's market. It can create a difference between private costs or values and social costs or values.

We will now show why markets characterized by externalities are not efficient, and also show how these externalities might be corrected or reduced. The essence of an externality is that it creates a divergence between private costs/benefits and social costs/benefits. If a steel producer pollutes the air, and the steel buyer pays only the costs incurred by the producer, then the buyer is not paying the full "social" cost of the product. The problem is illustrated in Figure 5.5.

Figure 5.5 Negative externalities and inefficiency



A negative externality is associated with this good. S reflects private costs, whereas S_f reflects the full social cost. The socially optimal output is Q^* , not the market outcome Q_0 . Beyond Q^* the real cost exceeds the demand value; therefore Q_0 is not an efficient output. A tax that increases P to P^* and reduces output is one solution to the externality.

Negative externalities

In Figure 5.5, the supply curve S represents the cost to the supplier, whereas S_f (the *full cost*) reflects, in addition, the cost of bad air to the population. Of course, we are assuming that this external cost is ascertainable, in order to be able to characterize S_f accurately. Note also that this illustration assumes that, as power output increases, the external cost *per unit* rises, because the difference between the two supply curves increases with output. This implies that low levels of pollution do less damage per unit: Perhaps the population has a natural tolerance for low levels, but higher levels cannot be tolerated easily and so the cost per unit is greater.

Despite the externality, *an efficient level of production can still be defined*. It is given by Q^* , not Q_0 . To see why, consider the impact of reducing output by one unit from Q_0 . At Q_0 the willingness of buyers to pay for the marginal unit supplied is E_0 . The (private) supply cost is also E_0 . But from a societal standpoint there is a pollution/health cost of AE_0 associated with that unit of production. The full cost, as represented by S_f , exceeds the buyer's valuation. Accordingly, if the last unit of output produced is cut, society gains by the amount AE_0 , because the cut in output reduces the excess of true cost over value.

Applying this logic to each unit of output between Q_0 and Q^* , it is evident that society can increase its well-being by the dollar amount equal to the area E^*AE_0 , as a result of reducing production.

Next, consider the consequences of reducing output further from Q^* . Note that some pollution is being created here, and environmentalists frequently advocate that pollution should be reduced to zero. However, an efficient outcome may not involve a zero level of pollution! If the production of power were reduced below Q^* , the loss in value to buyers, as a result of not being able to purchase the good, would exceed the full cost of its production.

If the government decreed that, instead of producing Q^* , no pollution would be tolerated, then society would forgo the possibility of earning the total real surplus equal to the area UE^*K . Economists do not advocate such a zero-pollution policy; rather, we

advocate a policy that permits a "tolerable" pollution level – one that still results in net benefits to society. In this particular example, the total cost of the tolerated pollution equals the area between the private and full supply functions, KE^*VR .

As a matter of policy, how is this market influenced to produce the amount Q^* rather than Q_0 ? One option would be for the government to intervene directly with production quotas for each firm. An alternative would be to impose a corrective tax on the good whose production causes the externality: With an appropriate increase in the price, consumers will demand a reduced quantity. In Figure 5.5 a tax equal to the dollar value VE^* would shift the supply curve upward by that amount and result in the quantity Q^* being traded.

A **corrective tax** seeks to direct the market towards a more efficient output.

We are now venturing into the field of environmental policy, where a corrective tax is usually called a carbon tax, and this is explored in the following section. The key conclusion of the foregoing analysis is that an efficient working of the market continues to have meaning in the presence of externalities. An efficient output level still maximizes economic surplus where surplus is correctly defined.

Positive externalities

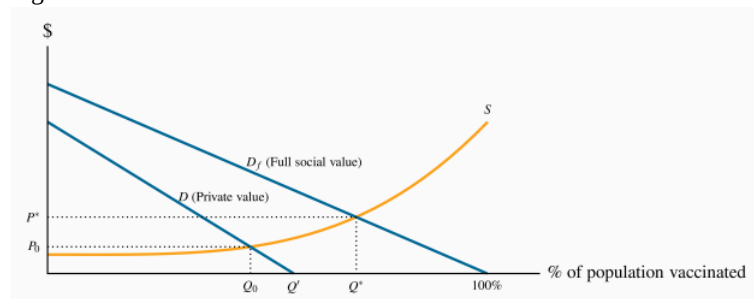
Externalities of the *positive* kind enable individuals or producers to get a type of 'free ride' on the efforts of others. Real world examples abound: When a large segment of the population is immunized against disease, the remaining individuals benefit on account of the reduced probability of transmission.

A less well recognized example is the benefit derived by many producers world-wide from research and development (R&D) undertaken in advanced economies and in universities and research institutes. The result is that society at large, including the corporate sector, gain from this enhanced understanding of science, the environment, or social behaviours.

The free market may not cope any better with these positive externalities than it does with negative externalities, and government intervention may be beneficial. Furthermore, firms that invest heavily in research and development would not undertake such investment if competitors could have a complete free ride and appropriate the fruits. This is why *patent laws* exist, as we shall see later in discussing Canada's competition policy. These laws prevent competitors from copying the product development of firms that invest in R&D. If such protection were not in place, firms would not allocate sufficient resources to R&D, which is a real engine of economic growth. In essence, the economy's research-directed resources would not be appropriately rewarded, and thus too little research would take place.

While patent protection is one form of corrective action, subsidies are another. We illustrated above that an appropriately formulated tax on a good that creates negative externalities can reduce demand for that good, and thereby reduce pollution. A subsidy can be thought of as a negative tax, and can stimulate the supply of goods and services that have positive externalities. Consider the example in Figure 5.6.

Figure 5.6 Positive externalities – the market for flu shots



The value to society of vaccinations exceeds the value to individuals: The greater the number of individuals vaccinated, the lower is the probability of others contracting the virus. D_f reflects this additional value. Consequently, the social optimum is Q^* which exceeds Q_0 .

Individuals have a demand for flu shots given by D . This reflects their private valuation – their personal willingness to pay. But the social value of flu shots is greater. When more individuals are vaccinated, the probability that others will be infected falls. Additionally, with higher rates of immunization, the health system will incur fewer costs in treating the infected. Therefore, the value to society of any quantity of flu shots is greater than the sum of the values that individuals place on them.

D_f reflects the full social value of any quantity of flu shots. In this instance the quantity axis measures the percentage of the population vaccinated, which has a maximum of 100%. If S is the supply curve, the socially optimal, efficient, market outcome is Q^* . The steeply upward-sloping section of S denotes that it may be very costly to vaccinate every last person – particularly those living in outlying communities. How can we influence the market to move from Q_0 towards Q^* ? One solution is a subsidy that would reduce the price to zero. In this case that gets us almost to the optimum, because the percentage of the population now choosing to be vaccinated is given by Q' . The zero price essentially makes the supply curve, as perceived by the population, to be running along the horizontal axis.

Note the social value of the improvement in moving from Q_0 to Q' ; the social value exceeds the social cost. But even at Q' further gains are available because at Q' the social value of additional vaccinations is greater than the social cost. Overall, at the point Q^* the social value is given by the area under the demand curve, and the social cost by the area under the supply curve.

This page titled [5.5: Market failures - externalities](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.6: Other market failures

There are other ways in which markets can fail to reflect accurately the social value or social cost of economic activity. Profit-seeking monopolies, which restrict output in order to increase profits, represent inefficient markets, and we will see why in the chapter on monopoly. Or the market may not deal very well with what are called public goods. These are goods, like radio and television service, national defence, or health information: With such goods and services many individuals can be supplied with the same good at the same total cost as one individual. We will address this problem in our chapter on government. And, of course, there are international externalities that cannot be corrected by national governments because the interests of adjoining states may differ: One economy may wish to see cheap coal-based electricity being supplied to its consumers, even if this means acid rain or reduced air quality in a neighbouring state. Markets may fail to supply an "efficient" amount of a good or service in all of these situations. Global warming is perhaps the best, and most extreme, example of international externalities and market failure.

This page titled [5.6: Other market failures](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.7: Environmental policy and climate change

Greenhouse gases

The greatest externality challenge in the modern world is to control our emissions of greenhouse gases. The emission of greenhouse gases (GHGs) is associated with a wide variety of economic activities such as coal-based power generation, oil-burning motors, wood-burning stoves, ruminant animals, *etc.* The most common GHG is carbon dioxide, methane is another. The gases, upon emission, circulate in the earth's atmosphere and, following an excessive build-up, prevent sufficient radiant heat from escaping. The result is a slow warming of the earth's surface and air temperatures. It is envisaged that such temperature increases will, in the long term, increase water temperatures and cause glacial melting, with the result that water levels worldwide will rise. In addition to the higher water levels, which the Intergovernmental Panel on Climate Change (IPCC) estimates will be between one foot and one metre by the end of the 21st century, oceans will become more acidic, weather patterns will change and weather events become more variable and severe. The changes will be latitude-specific and vary by economy and continent, and ultimately will impact the agricultural production abilities of certain economies.

Greenhouse gases that accumulate excessively in the earth's atmosphere prevent heat from escaping and lead to **global warming**.

While most scientific findings and predictions are subject to a degree of uncertainty, there is little disagreement in the scientific community on the long-term impact of increasing GHGs in the atmosphere. There is some skepticism as to whether the generally higher temperatures experienced in recent decades are completely attributable to anthropogenic activity since the industrial revolution, or whether they also reflect a natural cycle in the earth's temperature. But scientists agree that a continuance of the recent rate of GHG emissions is leading to serious climatic problems.

The major economic environmental challenge facing the world economy is this: Historically, GHG emissions have been strongly correlated with economic growth. The very high rate of economic growth in many large-population economies such as China and India that will be necessary to raise hundreds of millions out of poverty means that that historical pattern needs to be broken – GHG accumulation must be "decoupled" from economic growth.

GHGs as a common property

A critical characteristic of GHGs is that they are what we call in economics a 'common property': Every citizen in the world 'owns' them, every citizen has equal access to them, and it matters little where these GHGs originate. Consequently, if economy A reduces its GHG emissions, economy B may simply increase its emissions rather than incur the cost of reducing them. Hence, economy A's behaviour goes unrewarded. This is the crux of international agreements – or disagreements. Since GHGs are a common property, in order for A to have the incentive to reduce emissions, it needs to know that B will act correspondingly.

From the Kyoto Protocol to the Paris Accord

The world's first major response to climate concerns came in the form of the United Nations–sponsored Earth Summit in Rio de Janeiro in 1992. This was followed by the signing of the Kyoto Protocol in 1997, in which a group of countries committed themselves to reducing their GHG emissions relative to their 1990 emissions levels by the year 2012. Canada's Parliament subsequently ratified the Kyoto Protocol, and thereby agreed to meet Canada's target of a 6 percent reduction in GHGs relative to the amount emitted in 1990.

On a per-capita basis, Canada is one of the world's largest contributors to global warming, even though Canada's percentage of the total is just 2 percent. Many of the world's major economies refrained from signing the Protocol—most notably China, the United States, and India. Canada's emissions in 1990 amounted to approximately 600 giga tonnes (Gt) of carbon dioxide; but by the time we ratified the treaty in 2002, emissions were 25% above that level. Hence the signing was somewhat meaningless, in that Canada had virtually a zero possibility of attaining its target.

The target date of 2012 has come and gone and subsequent conferences in Copenhagen and Rio failed to yield an international agreement. But in Paris, December 2015, 195 economies committed to reduce their GHG emissions by specific amounts. Canada was a party to that agreement. Target reductions varied by country. Canada committed itself to reduce GHG emissions by 30% by the year 2030 relative to 2005 emissions levels. To this end the Liberal government of Prime Minister Justin Trudeau announced in late 2016 that if individual Canadian provinces failed to implement a carbon tax, or equivalent, the federal government would impose one unilaterally. The program involves a carbon tax of \$10 per tonne in 2018, that increases by \$10 per annum until it attains a value of \$50 in 2022. Some provinces already have GHG limitation systems in place (cap and trade systems - developed

below), and these provinces would not be subject to the federal carbon tax provided the province-level limitation is equivalent to the federal carbon tax.

Canada's GHG emissions

An excellent summary source of data on Canada's emissions and performance during the period 1990-2018 is available on Environment Canada's web site. See:

www.canada.ca/en/environment-climate-change/services/climate-change/greenhouse-gas-emissions/sources-sinks-executive-summary-2020.html#toc3

Canada, like many economies, has become more efficient in its use of energy (the main source of GHGs) in recent decades—its *use of energy per unit of total output* has declined steadily. Canada emitted 0.44 mega tonnes of CO₂ equivalent per billion dollars of GDP in 2005, and 0.36 mega tonnes in 2017. On a *per capita* basis Canada's emissions amounted to 22.9 tonnes in 2005, and dropped to 19.5 by 2017. This modest improvement in efficiency means that Canada's GDP is now *less energy intensive*. The critical challenge is to produce more output while using not just less energy per unit of output, but to use less energy in total.

While Canada's energy intensity (GHGs per unit of output) has dropped, overall emissions have increased by almost 20% since 1990. Furthermore, while developed economies have increased their efficiency, it is the *world's* efficiency that is ultimately critical. By outsourcing much of its manufacturing sector to China, Canada and the West have offloaded some of their most GHG-intensive activities. But GHGs are a common property resource.

Canada's GHG emissions also have a regional aspect: The *production* of oil and gas, which has created considerable wealth for all Canadians, is both energy intensive and concentrated in a limited number of provinces (Alberta, Saskatchewan and more recently Newfoundland and Labrador).

GHG measurement

GHG atmospheric concentrations are measured in parts per million (ppm). Current levels in the atmosphere are slightly above 400 ppm, and continued growth in concentration will lead to serious economic and social disruption. In the immediate pre-industrial revolution era concentrations were in the 280 ppm range. Hence, our world seems to be headed towards a doubling of GHG concentrations in the coming decades.

GHGs are augmented by the annual additions to the stock already in the atmosphere, and at the same time they decay—though very slowly. GHG-reduction strategies that propose an immediate reduction in emissions are more costly than those aimed at a more gradual reduction. For example, a slower investment strategy would permit in-place production and transportation equipment to reach the end of its economic life rather than be scrapped and replaced 'prematurely'. Policies that focus upon longer-term replacement are therefore less costly in this specific sense.

While not all economists and policy makers agree on the time scale for attacking the problem, the longer that GHG reduction is postponed, the greater the efforts will have to be in the long term—because GHGs will build up more rapidly in the near term.

A critical question in controlling GHG emissions relates to the cost of their control: How much of annual growth might need to be sacrificed in order to get emissions onto a sustainable path? Again estimates vary. The Stern Review (2006) proposed that, with an increase in technological capabilities, a strategy that focuses on the relative near-term implementation of GHG reduction measures might cost "only" a few percentage points of the value of world output. If correct, this is a low price to pay for risk avoidance in the longer term.

Nonetheless, such a reduction will require particular economic policies, and specific sectors will be impacted more than others.

Economic policies for climate change

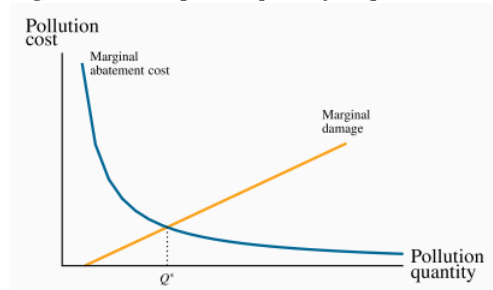
There are three main ways in which polluters can be controlled. One involves issuing direct controls; the other two involve incentives—in the form of pollution taxes, or on tradable "permits" to pollute.

To see how these different policies operate, consider first Figure 5.7. It is a standard diagram in environmental economics, and is somewhat similar to our supply and demand curves. On the horizontal axis is measured the quantity of environmental damage or pollution, and on the vertical axis its dollar value or cost. The upward-sloping damage curve represents the cost to society of each additional unit of pollution or gas, and it is therefore called a marginal damage curve. It is positively sloped to reflect the reality that, at low levels of emissions, the damage of one more unit is less than at higher levels. In terms of our earlier discussion, this

means that an increase in GHGs of 10 ppm when concentrations are at 300 ppm may be less damaging than a corresponding increase when concentrations are at 500 ppm.

The **marginal damage curve** reflects the cost to society of an additional unit of pollution.

Figure 5.7 The optimal quantity of pollution



Q^* represents the optimal amount of pollution. More than this would involve additional social costs because damages exceed abatement costs. Conversely, less than Q^* would require an abatement cost that exceeds the reduction in damage.

The second curve is the abatement curve. It reflects the cost of reducing emissions by one unit, and is therefore called a marginal abatement curve. This curve has a negative slope indicating that, as we reduce the total quantity of pollution produced (moving towards the origin on the horizontal axis), the cost of further unit reductions rises. This shape corresponds to reality. For example, halving the emissions of pollutants and gases from automobiles may be achieved by adding a catalytic converter and reducing the amount of lead in gasoline. But reducing those emissions all the way to zero requires the development of major new technologies such as electric cars—an enormously more costly undertaking.

The **marginal abatement curve** reflects the cost to society of reducing the quantity of pollution by one unit.

If producers are unconstrained in the amount of pollution they produce, they will produce more than what we will show is the optimal amount – corresponding to Q^* . This amount is optimal in the sense that at levels greater than Q^* the damage exceeds the cost of reducing the emissions. However, reducing emissions below Q^* would mean incurring a cost per unit reduction that exceeds the benefit of that reduction. Another way of illustrating this is to observe that at a level of pollution above Q^* the cost of reducing it is less than the damage it inflicts, and therefore a net gain accrues to society as a result of the reduction. But to reduce pollution below Q^* would involve an abatement cost greater than the reduction in pollution damage and therefore no net gain to society. This constitutes a first rule in optimal pollution policy.

An optimal quantity of pollution occurs when the marginal cost of abatement equals the marginal damage.

A second guiding principle emerges by considering a situation in which some firms are relatively 'clean' and others are 'dirty'. More specifically, a clean firm A may have already invested in new equipment that uses less energy per unit of output produced, or emits fewer pollutants per unit of output. In contrast, the dirty firm B uses older dirtier technology. Suppose furthermore that these two firms form a particular sector of the economy and that the government sets a limit on total pollution from this sector, and that this limit is less than what the two firms are currently producing. What is the least costly method to meet the target?

The intuitive answer to this question goes as follows: In order to reduce pollution at least cost to the sector, calculate what it would cost each firm to reduce pollution from its present level. Then implement a system so that the firm with the least cost of reduction is the first to act. In this case the 'dirty' firm will likely have a lower cost of abatement since it has not yet upgraded its physical plant. This leads to a second rule in pollution policy:

With many polluters, the least cost policy to society requires producers with the lowest abatement costs to act first.

This principle implies that policies which impose the same emission limits on firms may not be the least costly manner of achieving a target level of pollution. Let us now consider the use of tradable permits and corrective/carbon taxes as policy instruments. These are market-based systems aimed at reducing GHGs.

Tradable permits and **corrective/carbon taxes** are market-based systems aimed at reducing GHGs.

Incentive mechanism I: Tradable permits

A system of tradable permits is frequently called a 'cap and trade' system, because it limits or caps the total permissible emissions, while at the same time allows a market to develop in permits. For illustrative purposes, consider the hypothetical two-firm sector

we developed above, composed of firms A and B. Firm A has invested in clean technology, firm B has not. Thus it is less costly for B to reduce emissions than A if further reductions are required. Next suppose that each firm is allocated by the government a specific number of 'GHG emission permits'; and that the total of such permits is less than the amount of emissions at present, and that each firm is emitting more than its permits allow. How can these firms achieve the target set for this sector of the economy?

The answer is that they should be able to engage in mutually beneficial trade: If firm B has a lower cost of reducing emissions than A, then it may be in A's interest to pay B to reduce B's emissions heavily. Imagine that each firm is emitting 60 units of GHG, but they have permits to emit only 50 units each. And furthermore suppose it costs B \$20 to reduce GHGs by one unit, whereas it costs A \$30 to do this. In this situation A could pay B \$25 for several permits and this would benefit both firms. B can reduce GHGs at a cost of \$20 and is being paid \$25 to do this. In turn A would incur a cost of \$30 per unit to reduce his GHGs but he can buy permits from B for just \$25 and avoid the \$30 cost. Both firms gain, and the total cost to the economy is lower than if each firm had to reduce by the same amount.

The benefit of the cap 'n trade system is that it enables the marketplace to reduce GHGs at least cost.

The largest system of tradable permits currently operates in the European Union: The *EU Emissions Trading System*. It covers more than 10,000 large energy-using installations. Trading began in 2005. In North America a number of Western states and several Canadian provinces are joined, either as participants or observers, in the *Western Climate Initiative*, which is committed to reduce GHGs by means of tradable emissions permits. The longer-term goal of these systems is for the government to issue progressively fewer permits each year, and to include an ever larger share of GHG-emitting enterprises with the passage of time.

Policy in practice – international

In an ideal world, permits would be traded internationally, and such a system might be of benefit to developing economies: If the cost of reducing pollution is relatively low in developing economies because they have few controls in place, then developed economies, for whom the cost of GHG reduction is high, could induce firms in the developing world to undertake cost reductions. Such a trade would be mutually beneficial. For example, imagine in the above example that B is located in the developing world and A in the developed world. Both would obviously gain from such an arrangement, and because GHGs are a common property, the source of GHGs from a damage standpoint is immaterial.

Incentive mechanism II: Taxes

Corrective taxes are frequently called *Pigovian* taxes, after the economist Arthur Pigou. He advocated taxing activities that cause negative externalities. These taxes have been examined above in Section 5.4. Corrective taxes of this type can be implemented as part of a tax package *reform*. For example, taxpayers are frequently reluctant to see governments take 'yet more' of their money, in the form of new taxes. Such concerns can be addressed by reducing taxes in other sectors of the economy, in such a way that the package of tax changes maintains a 'revenue neutral' impact.

Revenues from taxes and permits

Taxes and tradable permits differ in that taxes generate revenue for the government from polluting producers, whereas permits may not generate revenue, or may generate less revenue. If the government simply *allocates* permits initially to all polluters, free of charge, and allows a market to develop, such a process generates no revenue to the government. While economists may advocate an *auction* of permits in the start-up phase of a tradable permits market, such a mechanism may run into political objections.

Setting taxes at the appropriate level requires knowledge of the cost and damage functions associated with GHGs. At the present time, economists and environmental scientists think that an appropriate price or tax on one tonne of GHG is in the \$40 – \$50 range. Such a tax would reduce emissions to a point where the longer-term impact of GHGs would not be so severe as otherwise.

British Columbia introduced a carbon tax of \$10 per tonne of GHG on fuels in 2008, and has increased that price regularly. This tax was designed to be revenue neutral in order to make it more acceptable. This means that British Columbia reduced its income tax rates by an amount such that income tax payments would fall by an amount equal to the revenue captured by the carbon tax.

GHG policy at the federal level in Canada is embodied in the *Greenhouse Gas Pollution Pricing Act of 2018*. As detailed earlier, the Act imposes a yearly increasing levy on emissions. The system is intended to be revenue neutral, in that the revenues will be returned to households in the form of a 'paycheck' by the federal government. Large emitters of GHGs are permitted a specific threshold number of tonnes of emission each year without being penalized. Beyond that threshold the above rates apply.

Will this amount of carbon taxation hurt consumers, and will it enable Canada to reach its 2030 GHG goal? As a specific example: the gasoline-pricing rule of thumb is that each \$10 in carbon taxation or pricing leads to an increase in the price of gasoline at the pump of about 2.5 cents. So a \$50 levy per tonne means gas at the pump should rise by 12.5 cents per litre. The proceeds are returned to households.

As for the goal of reaching the 2030 target announced at Paris: Environment Canada estimates that the pricing scheme will reduce GHG emissions by about 60 tonnes per annum. But Canada's goal stated in Paris is to reduce emissions in 2030 by approximately four times this amount. Under the Paris Accord, Canada stated that its 2030 goal would be to reduce emissions by 30% from their 2005 level of 725 MT, that is by an amount equal to approximately 220 tonnes.

Policy in practice – domestic large final emitters

Governments frequently focus upon quantities emitted by individual large firms, or *large final emitters (LFEs)*. In some economies, a relatively small number of producers are responsible for a disproportionate amount of an economy's total pollution, and limits are placed on those firms in the belief that significant economy-wide reductions can be achieved in this manner. One reason for concentrating on these LFEs is that the monitoring costs are relatively small compared to the costs associated with monitoring *all* firms in the economy. It must be kept in mind that pollution permits may be a legal requirement in some jurisdictions, but monitoring is still required, because firms could choose to risk polluting without owning a permit.

This page titled [5.7: Environmental policy and climate change](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.8: Conclusion

Welfare economics lies at the heart of public policy. Demand and supply curves can be interpreted as value curves and cost curves when there are no externalities involved. This is what enables us to define an efficient output of a product, and consequently an efficient use of the economy's resources. While efficiency is a central concept in economics, we must keep in mind that when the economic environment changes so too will the efficient use of resources, as we illustrated in Section 5.7.

In this chapter we have focused on equity issues through the lens of GHG emissions. The build-up of GHGs in our atmosphere invokes the concept of *intergenerational equity*: The current generation is damaging the environment and the costs of that damage will be borne by subsequent generations. Hence it is inequitable in the intergenerational sense for us to leave a negative legacy to succeeding generations. Equity arises within generations also. For example, how much more in taxes should the rich pay relative to the non-rich? We will explore this type of equity in our chapter on government.

This page titled [5.8: Conclusion](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.9: Key Terms

Welfare economics assesses how well the economy allocates its scarce resources in accordance with the goals of efficiency and equity.

Efficiency addresses the question of how well the economy's resources are used and allocated.

Equity deals with how society's goods and rewards are, and should be, distributed among its different members, and how the associated costs should be apportioned.

Consumer surplus is the excess of consumer willingness to pay over the market price.

Supplier or producer surplus is the excess of market price over the reservation price of the supplier.

Efficient market: maximizes the sum of producer and consumer surpluses.

Tax wedge is the difference between the consumer and producer prices.

Revenue burden is the amount of tax revenue raised by a tax.

Excess burden of a tax is the component of consumer and producer surpluses forming a net loss to the whole economy.

Deadweight loss of a tax is the component of consumer and producer surpluses forming a net loss to the whole economy.

Distortion in resource allocation means that production is not at an efficient output.

Externality is a benefit or cost falling on people other than those involved in the activity's market. It can create a difference between private costs or values and social costs or values.

Corrective tax seeks to direct the market towards a more efficient output.

Greenhouse gases that accumulate excessively in the earth's atmosphere prevent heat from escaping and lead to global warming.

Marginal damage curve reflects the cost to society of an additional unit of pollution.

Marginal abatement curve reflects the cost to society of reducing the quantity of pollution by one unit.

Tradable permits are a market-based system aimed at reducing GHGs.

Carbon taxes are a market-based system aimed at reducing GHGs.

This page titled [5.9: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine](#) ([Lyryx](#)) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

5.10: Exercises for Chapter 5

EXERCISE 5.1

Four teenagers live on your street. Each is willing to shovel snow from one driveway each day. Their "willingness to shovel" valuations (supply) are: Jean, \$10; Kevin, \$9; Liam, \$7; Margaret, \$5. Several households are interested in having their driveways shoveled, and their willingness to pay values (demand) are: Jones, \$8; Kirpinsky, \$4; Lafleur, \$7.50; Murray, \$6.

- Draw the implied supply and demand curves as step functions.
- How many driveways will be shoveled in equilibrium?
- Compute the maximum possible sum for the consumer and supplier surpluses.
- If a new (wealthy) family arrives on the block, that is willing to pay \$12 to have their driveway cleared, recompute the answers to parts (a), (b), and (c).

EXERCISE 5.2

Consider a market where supply curve is horizontal at $P=10$ and the demand curve has intercepts $\{ \$34, 34 \}$, and is defined by the relation $P=34-Q$.

- Illustrate the market geometrically.
- Impose a tax of \$2 per unit on the good so that the supply curve is now $P=12$. Illustrate the new equilibrium quantity.
- Illustrate in your diagram the tax revenue generated.
- Illustrate the deadweight loss of the tax.

EXERCISE 5.3

Next, consider an example of DWL in the labour market. Suppose the demand for labour is given by the fixed gross wage $W = \$16$. The supply is given by $W=0.8L$, indicating that the supply curve goes through the origin with a slope of 0.8.

- Illustrate the market geometrically.
- Calculate the supplier surplus, knowing that the equilibrium is $L=20$.
- Optional:* Suppose a wage tax is imposed that produces a net-of-tax wage equal to $W = \$12$. This can be seen as a downward shift in the demand curve. Illustrate the new quantity supplied and the new supplier's surplus.

EXERCISE 5.4

Governments are in the business of providing information to potential buyers. The first serious provision of information on the health consequences of tobacco use appeared in the United States Report of the Surgeon General in 1964.

- How would you represent this intervention in a supply and demand for tobacco diagram?
- Did this intervention "correct" the existing market demand?

EXERCISE 5.5

In deciding to drive a car in the rush hour, you think about the cost of gas and the time of the trip.

- Do you slow down other people by driving?
- Is this an externality, given that you yourself are suffering from slow traffic?

EXERCISE 5.6

Suppose that our local power station burns coal to generate electricity. The demand and supply functions for electricity are given by $P=12-0.5Q$ and $P=2+0.5Q$, respectively. The demand curve has intercepts $\{ \$12, 24 \}$ and the supply curve intercept is at \$2 with a slope of one half. However, for each unit of electricity generated, there is an externality. When we factor this into the supply side of the market, the real social cost is increased by \$1 per unit. That is, the supply curve shifts upwards by \$1, and now takes the form $P=3+0.5Q$.

- Illustrate the free-market equilibrium.
- Illustrate the efficient (i.e. socially optimal) level of production.

EXERCISE 5.7

Your local dry cleaner, Bleached Brite, is willing to launder shirts at its cost of \$1.00 per shirt. The neighbourhood demand for this service is $P=5-0.005Q$, knowing that the demand intercepts are $\{\$5, 1000\}$.

- Illustrate the market equilibrium.
- Suppose that, for each shirt, Bleached Brite emits chemicals into the local environment that cause \$0.25 damage per shirt. This means the full cost of each shirt is \$1.25. Illustrate graphically the socially optimal number of shirts to be cleaned.
- Optional:* Calculate the socially optimal number of shirts to be cleaned.

EXERCISE 5.8

The supply curve for agricultural labour is given by $W=6+0.1L$, where W is the wage (price per unit) and L the quantity traded. Employers are willing to pay a wage of \$12 to all workers who are willing to work at that wage; hence the demand curve is $W=12$.

- Illustrate the market equilibrium, if you are told that the equilibrium occurs where $L=60$.
- Compute the supplier surplus at this equilibrium.

EXERCISE 5.9

Optional: The market demand for vaccine XYZ is given by $P=36-Q$ and the supply conditions are $P=20$; so \$20 represents the true cost of supplying a unit of vaccine. There is a positive externality associated with being vaccinated, and the real societal value is known and given by $P=36-(1/2)Q$. This new demand curve represents the true value to society of each vaccination. This is reflected in the private value demand curve rotating upward around the price intercept of \$36.

- Illustrate the private and social demand curves on a diagram, with intercept values calculated.
 - What is the market solution to this supply and demand problem?
 - What is the socially optimal number of vaccinations?
- The demand and supply functions behind these curves are $P=90-1Q$ and $P=30+(1/4)Q$. Equating supply and demand yields the solutions in the text.
 - For example, if the demand curve shifts upwards, parallel, to become $P=120-1Q$ and the supply curve changes in slope to $P=30+(1/2)Q$, the new equilibrium solution is $\{P = \$60, Q = 60\}$.

This page titled [5.10: Exercises for Chapter 5](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

SECTION OVERVIEW

Unit 3: Decision Making by Consumer and Producers

Individuals and producers make choices when they interact with one another. We begin by supposing that individuals act in a manner that is at least in their self-interest – they transact in markets because they derive satisfaction or utility from doing so. Individuals may also be altruistic. But in each case we suppose they plan their actions in a way that is consistent with a goal. Utility may be measurable or only comparable, and we derive the tools of choice optimization using these two alternative notions of utility in Chapter 6.

Chapter 7 furnishes a link between buyers and sellers, between savers and investors, between households and firms. We explore the workings of corporations, and illustrate how investment in many different firms can reduce risk for an individual investor.

Like individual buyers, sellers and producers also act with a plan to reach a goal. They are interested in making a profit from their enterprise and produce goods and services at the least cost that is consistent with available production technologies. Cost minimization and profit maximization in the short run and the long run are explored in Chapter 8.

6: Individual choice

- 6.1: Rationality
- 6.2: Choice with measurable utility
- 6.3: Choice with ordinal utility
- 6.4: Applications of indifference analysis
- 6.5: Key Terms
- 6.6: Exercises for Chapter 6

7: Firms, investors and capital markets

- 7.1: Business organization
- 7.2: Profit
- 7.3: Risk and the investor
- 7.4: Risk pooling and diversification
- 7.5: Conclusion
- 7.6: Key Terms
- 7.7: Exercises for Chapter 7

8: Production and cost

- 8.1: Efficient production
- 8.2: The time frame
- 8.3: Production in the short run
- 8.4: Costs in the short run
- 8.5: Fixed costs and sunk costs
- 8.6: Long-run production and costs
- 8.7: Technological change- globalization and localization
- 8.8: Clusters, learning by doing, scope economics
- 8.9: Conclusion
- 8.10: Key Terms
- 8.11: Exercises for Chapter 8

This page titled [Unit 3: Decision Making by Consumer and Producers](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

6: Individual choice

Chapter 6: Individual choice

In this chapter we will explore:

6.1	Rationality
6.2	Consumer choice with measurable utility
6.3	Consumer choice with ordinal utility
6.4	Applications of indifference analysis

6.1 Rationality

A critical behavioural assumption in economics is that agents operate in a way that is oriented towards achieving a goal. This can be interpreted to mean that individuals and firms maximize their personal well-being and/or their profits. These players may have other goals in addition: Philanthropy and the well-being of others are consistent with individual optimization.

If individuals are to achieve their goals then they must act in a manner that will get them to their objective; broadly, they must act in a rational manner. The theory of individual maximization that we will develop in this chapter is based on that premise or assumption. In assuming individuals are rational we need not assume that they have every piece of information available to them that might be relevant for a specific decision or choice. Nor need we assume that they have super computers in their brain when they evaluate alternative possible strategies.

What we do need to assume, however, is that individuals act in a manner that is consistent with obtaining a given objective. The modern theory of behavioural economics and behavioural psychology examines decision making in a wide range of circumstances and has uncovered many fascinating behaviours – some of which are developed in Application Box 6.1 below.

We indicated in Chapter 1 that as social scientists, we require a *reliable model* of behaviour, that is, *a way of describing the essentials of choice that is consistent with everyday observations on individual behaviour patterns*. In this chapter, our aim is to understand more fully the behavioural forces that drive the demand side of the economy.

Economists analyze individual decision making using two different, yet complementary, approaches – utility analysis and indifference analysis. We begin by portraying individuals as maximizing their *measurable utility* (sometimes called *cardinal utility*); then progress to indifference analysis, where a weaker assumption is made on the ability of individuals to measure their satisfaction. In this second instance we do not assume that individuals can measure their utility numerically, only that they can say if one collection of goods and services yields them greater satisfaction than another group. This ranking of choices corresponds to what is sometimes called *ordinal utility* – because individuals can *order* groups of goods and services in ascending order of satisfaction. In each case individuals are perceived as rational maximizers or optimizers: They allocate their income so as to choose the outcome that will make them as well off as possible.

The second approach to consumer behaviour is frequently omitted in introductory texts. It can be omitted here without interpreting the flow of ideas, although it does yield additional insights into consumer choice and government policy. As in preceding chapters, we begin the analysis with a motivating numerical example.

Application Box 6.1 Rationality and impulse

A number of informative and popular books on decision making have appeared recently. Their central theme is that our decision processes should not be viewed solely as a rational computer – operating in one single mode only, and unmoved by our emotions or history. Psychologists now know that our brains have at least two decision modes, and these are developed by economics Nobel Prize winner Daniel Kahneman in his book "Thinking, Fast and Slow". One part of our brain operates in a rational goal-oriented forward-looking manner (the 'slow' part), another is motivated by immediate gratification (the 'fast' part). Decisions that we observe in the world about us reflect these different mechanisms.

Richard Thaler, a Chicago economist and his law professor colleague Cass Sunstein, have developed a role for public policy in their book entitled "Nudge". They too argue that individuals do not inevitably operate in their own best long-term interests, and as a consequence individuals frequently require a *nudge* by government to make the long-term choice rather than the short-term choice. For example, when individuals begin a new job, they might be automatically enrolled in the company pension plan and be given

the freedom to opt out, rather than not be enrolled and given the choice to opt in. Such policies are deemed to be 'soft paternalism'. They are paternalistic for the obvious reason – another organism is directing, but they are also soft in that they are not binding.

6.2 Choice with measurable utility

Neal loves to pump his way through the high-altitude powder at the Whistler ski and snowboard resort. His student-rate lift-ticket cost is \$30 per visit. He also loves to frequent the jazz bars in downtown Vancouver, and each such visit costs him \$20. With expensive passions, Neal must allocate his monthly entertainment budget carefully. He has evaluated how much satisfaction, measured in utils, he obtains from each snowboard outing and each jazz club visit. We assume that these utils are measurable, and use the term cardinal utility to denote this. These measurable utility values are listed in columns 2 and 3 of Table 6.1. They define the total utility he gets from various amounts of the two activities.

Table 6.1 Utils from snowboarding and jazz

1	2	3	4	5	6	7
Visit	Total	Total	Marginal	Marginal	Marginal	Marginal
#	snowboard	jazz	snowboard	jazz utils	snowboard	jazz utils
	utils	utils	utils		utils per \$	per \$
1	72	52	72	52	2.4	2.6
2	132	94	60	42	2.0	2.1
3	182	128	50	34	1.67	1.7
4	224	156	42	28	1.4	1.4
5	260	180	36	24	1.2	1.2
6	292	201	32	21	1.07	1.05
7	321	220	29	19	0.97	0.95

Price of snowboard visit=\$30. Price of jazz club visit=\$20.

Cardinal utility is a measurable concept of satisfaction.

Total utility is a measure of the total satisfaction derived from consuming a given amount of goods and services.

Neal's total utility from each activity in this example is independent of the amount of the other activity he engages in. These total utilities are plotted in Figures 6.1 and 6.2. Clearly, more of each activity yields more utility, so the additional or marginal utility (*MU*) of each activity is positive. This positive marginal utility for any amount of the good consumed, no matter how much, reflects the assumption of *non-satiation*—more is always better. Note, however, that the decreasing slopes of the total utility curves show that *total utility is increasing at a diminishing rate*. While more is certainly better, each additional visit to Whistler or a jazz club augments Neal's utility by a smaller amount. At the margin, his additional utility declines: He has diminishing marginal utility. The marginal utilities associated with snowboarding and jazz are entered in columns 4 and 5 of Table 6.1. They are the differences in total utility values when consumption increases by one unit. For example, when Neal makes a sixth visit to Whistler his total utility increases from 260 utils to 292 utils. His marginal utility for the sixth unit is therefore 32 utils, as defined in column 4. In light of this example, it should be clear that we can define marginal utility as:

$$\text{Marginal Utility} = \frac{\text{additional utility}}{\text{additional consumption}} \text{ or, } MU = \frac{\Delta U}{\Delta C}, \quad (6.1)$$

where ΔC denotes the change in the quantity consumed of the good or service in question.

Marginal utility is the addition to total utility created when one more unit of a good or service is consumed.

Diminishing marginal utility implies that the addition to total utility from each extra unit of a good or service consumed is declining.

Figure 6.1 TU from snowboarding

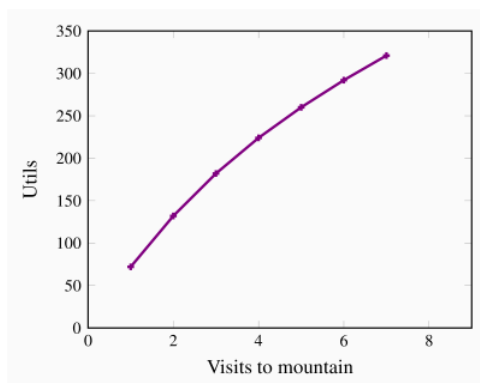


Figure 6.2 TU from jazz

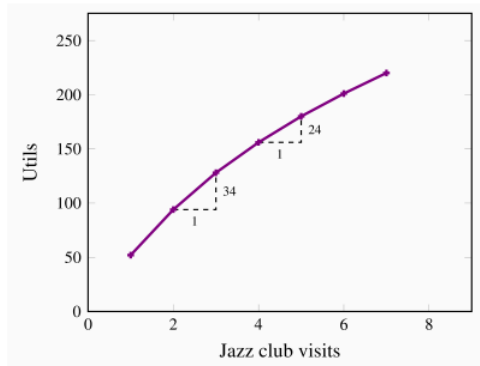


Figure 6.3 MU from snowboarding

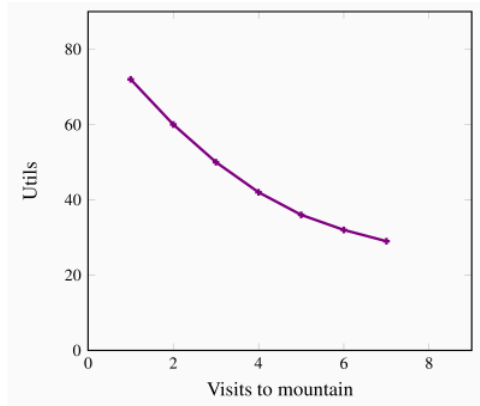
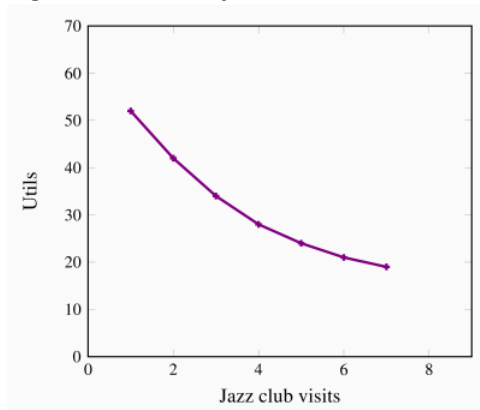


Figure 6.4 MU from jazz



The marginal utilities associated with consuming different amounts of the two goods are plotted in Figures 6.3 and 6.4, using the data from columns 4 and 5 in Table 6.1. These functions are declining, as indicated by their negative slope. It should also be clear that the *MU* curves can be derived from the *TU* curves. For example, in figure 6.2, when going from 2 units to 3 units of Jazz, *TU* increases by 34 units. But 34/1 is the slope of the *TU* function in this range of consumption – the vertical distance divided by the horizontal distance. Similarly, if jazz consumption increases from 4 units to 5 units the corresponding change in *TU* is 24 units, again the vertical distance divided by the horizontal distance, and so the slope of the function. In short, the *MU* is the slope of the *TU* function.

Now that Neal has defined his utility schedules, he must consider the price of each activity. Ultimately, when deciding how to allocate his monthly entertainment budget, he must evaluate how much utility he gets from each dollar spent on snowboarding and jazz: What "bang for his buck" does he get? Let us see how he might go about allocating his budget. *When he has fully spent his budget in the manner that will yield him greatest utility, we say that he has attained equilibrium*, because he will have no incentive to change his expenditure patterns.

If he boards once, at a cost of \$30, he gets 72 utils of satisfaction, which is 2.4 utils per dollar spent ($=72/30$). One visit to a jazz club would yield him 2.6 utils per dollar ($=52/20$). Initially, therefore, his dollars give him more *utility per dollar* when spent on jazz. His *MU* per dollar spent on each activity is given in the final two columns of the table. These values are obtained by dividing the *MU* associated with each additional unit by the good's price.

We will assume that Neal has a budget of \$200. He realizes that his initial expenditure should be on a jazz club visit, because he gets more utility per dollar spent there. Having made one such expenditure, he sees that a second jazz outing would yield him 2.1 utils per dollar expended, while a first visit to Whistler would yield him 2.4 utils per dollar. Accordingly, his second activity is a snowboard outing.

Having made one jazz and one snowboarding visit, he then decides upon a second jazz club visit for the same reason as before—utility value for his money. He continues to allocate his budget in this way until his budget is exhausted. In our example, this occurs when he spends \$120 on four snowboarding outings and \$80 on four jazz club visits. At this consumer equilibrium, he gets the same utility value per dollar for the last unit of each activity consumed. This is a necessary condition for him to be maximizing his utility, that is, to be in equilibrium.

Consumer equilibrium occurs when marginal utility per dollar spent on the last unit of each good is equal.

To be absolutely convinced of this, imagine that Neal had chosen instead to board twice and to visit the jazz clubs seven times; this combination would also exhaust his \$200 budget exactly. With such an allocation, he would get 2.0 utils per dollar spent on his marginal (second) snowboard outing, but just 0.95 utils per dollar spent on his marginal (seventh) jazz club visit.¹ If, instead, he were to reallocate his budget in favour of snowboarding, he would get 1.67 utils per dollar spent on a third visit to the hills. By reducing the number of jazz visits by one, he would lose 0.95 utils per dollar reallocated. Consequently, the utility gain from a reallocation of his budget towards snowboarding would outweigh the utility loss from allocating fewer dollars to jazz. His initial allocation, therefore, was not an optimum, or equilibrium.

Only when the utility per dollar expended on each activity is equal at the margin will Neal be optimizing. When that condition holds, a reallocation would be of no benefit to him, because the gains from one more dollar on boarding would be exactly offset by the loss from one dollar less spent on jazz. Therefore, we can write the equilibrium condition as

Equilibrium requires: $\frac{MU_s}{P_s} = \frac{MU_j}{P_j}$ or $\frac{MU_s}{MU_j} = \frac{P_s}{P_j}$.	(6.2)
--	-------

While this example has just two goods, in the more general case of many goods, this same *condition must hold for all pairs of goods* on which the consumer allocates his or her budget.

From utility to demand

Utility theory is a useful way of analyzing how a consumer makes choices. But in the real world we do not observe a consumer's utility, either total or marginal. Instead, his or her behaviour in the marketplace is observed through the demand curve. How are utility and demand related?

Demand functions relate the quantity of a good consumed to the price of that good, other things being equal. So let us trace out the effects of a price change on demand, with the help of this utility framework. We will introduce a simplification here: Goods are divisible, or that they come in small packages relative to income. Think, for example, of kilometres driven per year, or liters of

gasoline purchased. Conceptualizing things in this way enables us to imagine more easily experiments in which small amounts of a budget are allocated one dollar at a time. In contrast, in the snowboard/jazz example, we had to reallocate the budget in lumps of \$30 or \$20 at a time because we could not "fractionalize" these goods.

The effects of a price change on a consumer's demand can be seen through the condition that describes his or her equilibrium. If income is allocated to, say, three goods $\{a, b, c\}$, such that $MU_a/P_a = MU_b/P_b = MU_c/P_c$, and the price of, say, good b falls, the consumer must reallocate the budget so that once again the MU s per dollar spent are all equated. How does he do this? Clearly, if he purchases more or less of any one good, the MU changes. If the price of good b falls, then the consumer initially gets more utility from good b for the last dollar he spends on it (the denominator in the expression MU_b/P_b falls, and consequently the value of the ratio rises to a value greater than the values for goods a and c).

The consumer responds to this, in the first instance, by buying more of the cheaper good. He obtains more total utility as a consequence, and in the process will get *less utility at the margin* from that good. In essence, the numerator in the expression then falls, in order to realign it with the lower price. This equality also provides an underpinning for what is called the law of demand: More of a good is demanded at a lower price. If the price of any good falls, then, in order for the equilibrium condition to be re-established, the MU of that good must be driven down also. Since MU declines when more is purchased, this establishes that demand curves must slope downwards.

The **law of demand** states that, other things being equal, more of a good is demanded the lower is its price.

However, the effects of a price decline are normally more widespread than this, because the quantities of other goods consumed may also change. As explained in earlier chapters, the decline in the price of good b will lead the consumer to purchase more units of *complementary goods* and fewer units of goods that are *substitutes*. So the whole budget allocation process must be redetermined in response to any price change. But at the end of the day, a new equilibrium must be one where the marginal utility per dollar spent on each good is equal.

Applying the theory

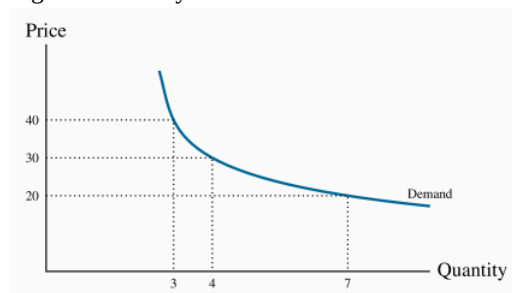
The demand curves developed in Chapter 3 can be related to the foregoing utility analysis. In our example, Neal purchased four lift tickets at Whistler when the price was \$30. We can think of this combination as one point on his demand curve, where the "other things kept constant" are the price of jazz, his income, his tastes, etc.

Suppose now that the price of a lift ticket increased to \$40. How could we find another point on his demand curve corresponding to this price, using the information in Table 6.1? The marginal utility per dollar associated with each visit to Whistler could be recomputed by dividing the values in column 4 by 40 rather than 30, yielding a new column 6. We would then determine a new allocation of his budget between the two goods that would maximize utility. After such a calculation we would find that he makes three visits to Whistler and four jazz-club visits. Thus, the combination ($P_s = \$40, Q_s = 3$) is another point on his demand curve. Note that this allocation exactly exhausts his \$200 budget.

By setting the price equal to \$20, this exercise could be performed again, and the outcome will be a quantity demanded of lift tickets equal to seven (plus three jazz club visits). Thus, the combination ($P_s = \$20, Q_s = 7$) is another point on his demand curve. Figure 6.5 plots a demand curve going through these three points.

By repeating this exercise for many different prices, the demand curve is established. We have now linked the demand curve to utility theory.

Figure 6.5 Utility to demand



When $P = \$30$, the consumer finds the quantity such that MU/P is equal for all purchases. The corresponding quantity purchased is 4 tickets. At prices of \$40 and \$20 the equilibrium condition implies quantities of 3 and 7 respectively.

Application Box 6.2 Individual and Collective Utility

The example developed in the text is not far removed from what economists do in practice. From a philosophical standpoint, economists are supposed to be interested in the well-being of the citizens who make up an economy or a country. To determine how 'well-off' citizens may be, social scientists frequently carry out surveys on how 'content' or 'happy' people are in their every-day lives. For example, the *Earth Institute at Columbia University* regularly produces a 'World Happiness Report'. The report is based upon responses to survey questions in numerous economies. One of the measures it uses to compare utility levels is the *Cantril ladder*. This is an 11-point scale running from 0 to 10, with the lowest value signifying the worst possible life, and 10 the highest possible quality of life. In reporting their findings, the researchers are essentially claiming that some economies have, on average, more contented or happier, people than others. Utility can be considered in exactly this way: A higher reported value on the Cantril ladder suggests higher utility.

A slightly different measure of well-being across economies is given by the *United Nations Human Development Index*. In this case, countries score high by having a high level of income, good health (as measured by life expectancy), and high levels of education, as measured by the number of years of education completed or envisaged.

In practice, social scientists are very comfortable using utility-based concepts to describe the economic circumstances of individuals in different economies.

6.3 Choice with ordinal utility

The budget constraint

In the preceding section, we assumed that utility is measurable in order to better understand how consumers allocate their budgets, and how this process is reflected in the market demands that are observed. The belief that utility might be measurable is not too extreme in the modern era. Neuroscientists are mapping more and more of the human brain and understanding how it responds to positive and negative stimuli. At the same time, numerous sociological surveys throughout the world ask individuals to rank their happiness on a scale of one to ten, or something similar, with a view to making comparisons between individual-level and group-level happiness – see Application Box 6.2. Nonetheless, not every scientist may be convinced that we should formulate behavioural rules on this basis. Accordingly we now examine the economics of consumer behaviour without this strong assumption. We assume instead that individuals are able to identify (a) different combinations of goods and services that yield equal satisfaction, and (b) combinations of goods and services that yield more satisfaction than other combinations. In contrast to measurable (or cardinal) utility, this concept is called ordinal utility, because it assumes only that consumers can order utility bundles rather than quantify the utility.

Ordinal utility assumes that individuals can rank commodity bundles in accordance with the level of satisfaction associated with each bundle.

The budget constraint

Neal's monthly expenditure limit, or budget constraint, is \$200. In addition, he faces a price of \$30 for lift tickets and \$20 per visit to jazz clubs. Therefore, using S to denote the number of snowboard outings and J the number of jazz club visits, if he spends his entire budget it must be true that the sum of expenditures on each activity exhausts his budget or income (I):

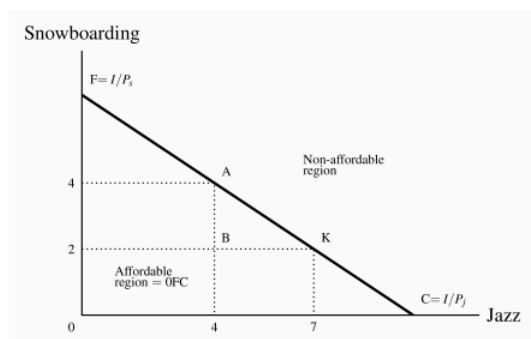
Expenditure on snowboarding + expenditure on Jazz	= Income
(Price of $S \times$ quantity of S) + (price of $J \times$ quantity of J)	= Income
$P_S S + P_J J = I$ or $\$30S + \$20J$	= \$200

Since many different combinations of the two goods are affordable, it follows that the budget constraint defines all bundles of goods that the consumer can afford with a given budget.

The **budget constraint** defines all bundles of goods that the consumer can afford with a given budget.

The budget constraint, then, is just what it claims to be—a limit on behaviour. Neal's budget constraint is illustrated in Figure 6.6, where the amount of each good consumed is given on the axes. If he spends all of his \$200 income on jazz, he can make exactly ten jazz club visits ($\$200/\$20 = 10$). The calculation also applies to visits to Whistler. The intercept value is always obtained by dividing income by the price of the good or activity in question.

Figure 6.6 The budget line



FC is the budget constraint and defines the affordable combinations of snowboarding and jazz. F represents all income spent on snowboarding. Thus $F = I/P_s$. Similarly $C = I/P_j$. Points above FC are not attainable. The slope = $OF/OC = (I/P_s)/(I/P_j) = P_j/P_s = 20/30 = 2/3$. The affordable set is OFC.

In addition to these affordable extremes, Neal can also afford many other bundles, e.g., $(S=2, J=7)$, or $(S=4, J=4)$, or $(S=6, J=1)$. The set of feasible, or affordable, combinations is bounded by the budget line, and this is illustrated in Figure 6.6.

The **affordable set** of goods and services for the consumer is bounded by the budget line from above; the **non-affordable set** lies strictly above the budget line.

The slope of the budget line is informative. As illustrated in Chapter 1, it indicates how many snowboard visits must be sacrificed for one additional jazz visit; it defines the consumer's *trade-offs*. To illustrate: Suppose Neal is initially at point A $(J=4, S=4)$, and moves to point K $(J=7, S=2)$. Clearly, both points are affordable. In making the move, he trades two snowboard outings in order to get three additional jazz club visits, a trade-off of $2/3$. This trade-off is the slope of the budget line, which, in Figure 6.6, is $AB/BK = -2/3$, where the negative sign reflects the downward slope.

Could it be that this ratio reflects the two prices $(\$20/\$30)$? The answer is yes: The slope of the budget line is given by the vertical distance divided by the horizontal distance, OF/OC . The points F and C were obtained by dividing income by the respective price—remember that the jazz intercept is $\$200/\$20 = 10$. Formally, that is I/P_j . The intercept on the snowboard axis is likewise I/P_s . Accordingly, the slope of the budget constraint is:

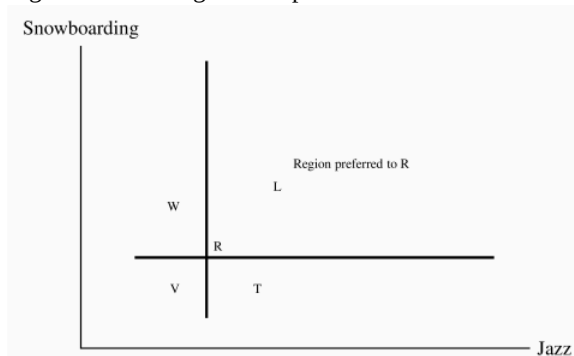
$$\text{Slope} = OF/OC = \frac{(I/P_s)}{(I/P_j)} = \frac{I}{P_s} \times \frac{P_j}{I} = \frac{P_j}{P_s}.$$

Since the budget line has a negative slope, it is technically correct to define it with a negative sign. But, as with elasticities, the sign is frequently omitted.

Tastes and indifference

We now consider how to represent a consumer's tastes in two dimensions, given that he can order, or rank, different consumption bundles, and that he can define a series of different bundles that all yield the same satisfaction. We limit ourselves initially to considering just "goods," and not "bads" such as pollution.

Figure 6.7 Ranking consumption bundles



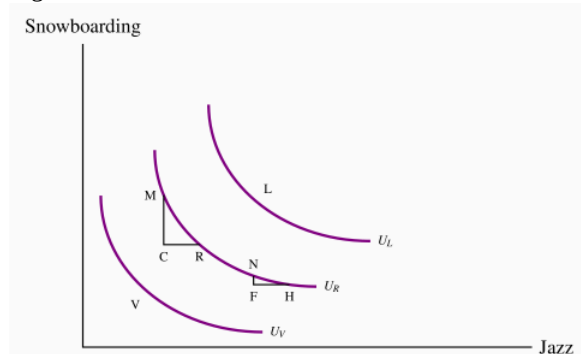
L is preferred to R since more of each good is consumed at L, while points such as V are less preferred than R. Points W and T contain more of one good and less of the other than R. Consequently, we cannot say if they are preferred to R without knowing how the consumer trades the goods off—that is, his preferences.

Figure 6.7 examines the implications of these assumptions about tastes. Each point shows a consumption bundle of snowboarding and jazz. Let us begin at bundle R. Since more of a good is preferred to less, any point such as L, which lies to the northeast of R, is preferred to R, since L offers more of both goods than R. Conversely, points to the southwest of R offer *less of each good* than R, and therefore R is preferred to a point such as V.

Without knowing the consumer's tastes, we cannot be sure at this stage how points in the northwest and southeast regions compare with R. At W or T, the consumer has more of one good and less of the other than at R. Someone who really likes snowboarding might prefer W to R, but a jazz buff might prefer T to R.

Let us now ask Neal to disclose his tastes, by asking him to define several combinations of snowboarding and jazz that *yield him exactly the same degree of satisfaction as the combination at R*. Suppose further, for reasons we shall understand shortly, that his answers define a series of points that lie on the beautifully smooth contour U_R in Figure 6.8. Since he is indifferent between all points on U_R by construction, this contour is an indifference curve.

Figure 6.8 Indifference curves



An indifference curve defines a series of consumption bundles, all of which yield the same satisfaction. The slope of an indifference curve is the marginal rate of substitution (*MRS*) and defines the number of units of the good on the vertical axis that the individual will trade for one unit of the good on the horizontal axis. The *MRS* declines as we move south-easterly, because the consumer values the good more highly when he has less of it.

An **indifference curve** defines combinations of goods and services that yield the same level of satisfaction to the consumer.

Pursuing this experiment, we could take other points in Figure 6.8, such as L and V, and ask the consumer to define bundles that would yield the same level of satisfaction, or indifference. These combinations would yield additional contours, such as U_L and U_V in Figure 6.8. This process yields a series of indifference curves that together form an indifference map.

An **indifference map** is a set of indifference curves, where curves further from the origin denote a higher level of satisfaction.

Let us now explore the properties of this map, and thereby understand why the contours have their smooth convex shape. They have four properties. The first three follow from our preceding discussion, and the fourth requires investigation.

1. Indifference curves *further from the origin reflect higher levels of satisfaction*.
2. Indifference curves are *negatively sloped*. This reflects the fact that if a consumer gets more of one good she should have less of the other in order to remain indifferent between the two combinations.
3. Indifference curves *cannot intersect*. If two curves were to intersect at a given point, then we would have two different levels of satisfaction being associated with the same commodity bundle—an impossibility.
4. Indifference curves are convex when viewed from the origin, reflecting a *diminishing marginal rate of substitution*.

The convex shape reflects an important characteristic of preferences: When consumers have a lot of some good, they value a marginal unit of it less than when they have a small amount of that good. More formally, they have a *higher marginal valuation at low consumption levels*—that first cup of coffee in the morning provides greater satisfaction than the second or third cup.

Consider the various points on U_R , starting at M in Figure 6.8. At M Neal snowboards a lot; at N he boards much less. The convex shape of his indifference map shows that he values a marginal snowboard trip more at N than at M. To see this, consider what happens as he moves along his indifference curve, starting at M. We have chosen the coordinates on U_R so that, in moving from M to R, and again from N to H, the additional amount of jazz is the same: $CR = FH$. From M, if Neal moves to R, he consumes an additional amount of jazz, CR. By definition of the indifference curve, he is willing to give up MC snowboard outings. The ratio

MC/CR defines his willingness to substitute one good for the other. This ratio, being a vertical distance divided by a horizontal distance, is the slope of the indifference curve and is called the marginal rate of substitution, *MRS*.

The **marginal rate of substitution** is the slope of the indifference curve. It defines the amount of one good the consumer is willing to sacrifice in order to obtain a given increment of the other, while maintaining utility unchanged.

At N, the consumer is willing to sacrifice the amount NF of boarding to get the same additional amount of jazz. Note that, when he boards *less*, as at N, he is willing to give up less boarding than when he has a lot of it, as at M, in order to get the same additional amount of jazz. His willingness to substitute *diminishes* as he moves from M to N: The quantity NF is less than the quantity MC. In order to reflect this taste characteristic, the indifference curve has a diminishing marginal rate of substitution: A flatter slope as we move down along its surface.

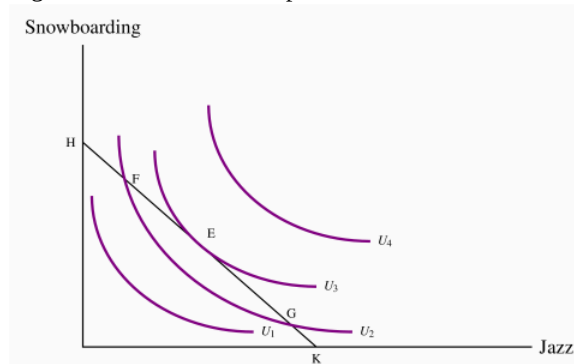
A **diminishing marginal rate of substitution** reflects a higher marginal value being associated with smaller quantities of any good consumed.

Optimization

We are now in a position to examine how the consumer optimizes—how he gets to the highest level of satisfaction possible. The constraint on his behaviour is the affordable set defined in Figure 6.6, the budget line.

Figure 6.9 displays several of Neal's indifference curves in conjunction with his budget constraint. We propose that he maximizes his utility, or satisfaction, at the point E, on the indifference curve denoted by U_3 . While points such as F and G are also on the boundary of the affordable set, they do not yield as much satisfaction as E, because E lies on a higher indifference curve. *The highest possible level of satisfaction is attained, therefore, when the budget line touches an indifference curve at just a single point—that is, where the constraint is tangent to the indifference curve.* E is such a point.

Figure 6.9 The consumer optimum



The budget constraint constrains the individual to points on or below HK. The highest level of satisfaction attainable is U_3 , where the budget constraint just touches, or is just tangent to, it. At this optimum the slope of the budget constraint ($-P_j/P_s$) equals the *MRS*.

This tangency between the budget constraint and an indifference curve requires that the slopes of each be the same at the point of tangency. We have already established that the slope of the budget constraint is the negative of the price ratio ($= -P_x/P_y$). The slope of the indifference curve is the marginal rate of substitution *MRS*. It follows, therefore, that the consumer optimizes where the marginal rate of substitution equals the slope of the price line.

Optimization requires:

Slope of Indifference curve = marginal rate of substitution = $-\frac{P_j}{P_s}$.	(6.3)
--	-------

A **consumer optimum** occurs where the chosen consumption bundle is a point such that the price ratio equals the marginal rate of substitution.

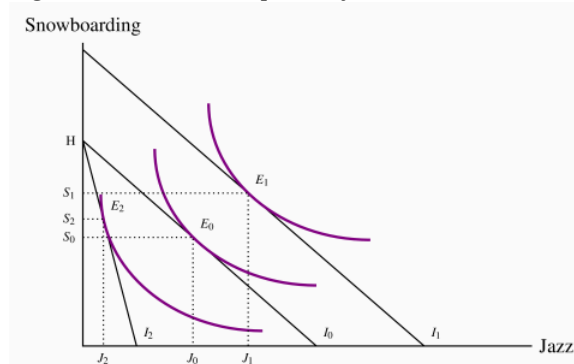
Notice the resemblance between this condition and the one derived in the first section as Equation 6.2. There we argued that equilibrium requires the ratio of the marginal utilities be same as the ratio of prices. Here we show that the *MRS* must equal the ratio of prices. In fact, with a little mathematics it can be shown that the *MRS* is indeed the same as the (negative of the) ratio of the marginal utilities: $MRS = -MU_j/MU_s$. Therefore the two conditions are in essence the same! However, it was not necessary to assume

that an individual can actually measure his utility in obtaining the result that the *MRS* should equal the price ratio in equilibrium. The concept of ordinal utility is sufficient.

Adjusting to income changes

Suppose now that Neal's income changes from \$200 to \$300. How will this affect his consumption decisions? In Figure 6.10, this change is reflected in a *parallel* outward shift of the budget constraint. Since no price change occurs, the slope remains constant. By recomputing the ratio of income to price for each activity, we find that the new snowboard and jazz intercepts are 10 ($= \$300/\30) and 15 ($= \$300/\20), respectively. Clearly, the consumer can attain a higher level of satisfaction—at a new tangency to a higher indifference curve—as a result of the size of the affordable set being expanded. In Figure 6.10, the new equilibrium is at E_1 .

Figure 6.10 Income and price adjustments



An income increase shifts the budget constraint from I_0 to I_1 . This enables the consumer to attain a higher indifference curve. A price rise in jazz tickets rotates the budget line I_0 inwards around the snowboard intercept to I_2 . The price rise reflects a lower real value of income and results in a lower equilibrium level of satisfaction.

Adjusting to price changes

Next, consider the impact of a price change from the initial equilibrium E_0 in Figure 6.10. Suppose that jazz now costs more. This reduces the purchasing power of the given budget of \$200. The new jazz intercept is therefore reduced. The budget constraint becomes steeper and rotates around the snowboard intercept H , which is unchanged because its price is constant. The new equilibrium is at E_2 , which reflects a lower level of satisfaction because the affordable set has been reduced by the price increase. As explained in Section 6.2, E_0 and E_2 define points on the demand curve for jazz (J_0 and J_2): They reflect the consumer response to a change in the price of jazz with all other things held constant. In contrast, the price increase for jazz *shifts* the demand curve for snowboarding: As far as the demand curve for snowboarding is concerned, a change in the price of jazz is one of those things other than own-price that determine its position.

Philanthropy

Individuals in the foregoing analysis aim to maximize their utility, given that they have a fixed budget. Note that this behavioural assumption does not rule out the possibility that these same individuals may be philanthropic – that is, they get utility from the act of giving to their favourite charity or the United Way or Centre-aide. To see this suppose that donations give utility to the individual in question – she gets a 'warm glow' feeling as a result of giving, which is to say she gets utility from the activity. There is no reason why we cannot put charitable donations on one axis and some other good or combination of goods on the remaining axis. At equilibrium, the marginal utility per dollar of contributions to charity should equal the marginal utility per dollar of expenditure on other goods; or, stated in terms of ordinal utility, the marginal rate of substitution between philanthropy and any other good should equal the ratio of their prices. Evidently the price of a dollar of charitable donations is one dollar.

6.4 Applications of indifference analysis

Price impacts: Complements and substitutes

The nature of complements and substitutes, defined in Chapter 4, can be further understood with the help of Figure 6.10. The new equilibrium E_2 has been drawn so that the increase in the price of jazz results in more snowboarding—the quantity of S increases to S_2 from S_0 . These goods are substitutes in this picture, because snowboarding *increases* in response to an *increase* in the price of jazz. If the new equilibrium E_2 were at a point yielding a lower level of S than S_0 , we would conclude that they were complements.

Cross-price elasticities

Continuing with the same price increase in jazz, we could compute the *percentage* change in the quantity of snowboarding demanded as a result of the *percentage* change in the jazz price. In this example, the result would be a positive elasticity value, because the quantity change in snowboarding and the price change in jazz are both in the same direction, each being positive.

Income impacts: Normal and inferior goods

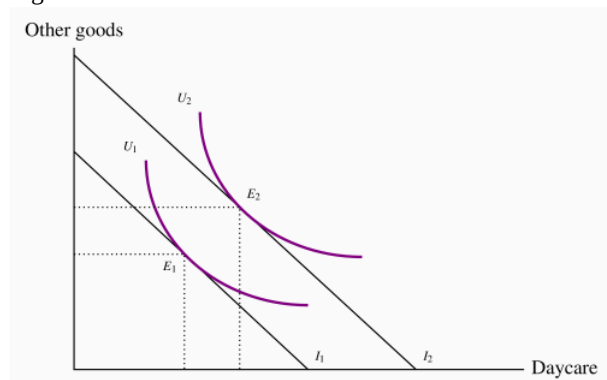
We know from Chapter 4 that the quantity demanded of a *normal good* increases in response to an income increase, whereas the quantity demanded of an *inferior good* declines. Clearly, both jazz and boarding are normal goods, as illustrated in Figure 6.10, because more of each one is demanded in response to the income increase from I_0 to I_1 . It would challenge the imagination to think that either of these goods might be inferior. But if J were to denote junky (inferior) goods and S super goods, we could envisage an equilibrium E_1 to the northwest of E_0 in response to an income increase, along the constraint I_1 ; less J and more S would be consumed in response to the income increase.

Policy: Income transfers and price subsidies

Government policies that improve the purchasing power of low-income households come in two main forms: Pure income transfers and price subsidies. *Social Assistance* payments ("welfare") or *Employment Insurance* benefits, for example, provide an increase in income to the needy. Subsidies, on the other hand, enable individuals to purchase particular goods or services at a lower price—for example, rent or daycare subsidies.

In contrast to taxes, which *reduce* the purchasing power of the consumer, subsidies and income transfers *increase* purchasing power. The impact of an income transfer, compared with a pure price subsidy, can be analyzed using Figures 6.11 and 6.12.

Figure 6.11 Income transfer



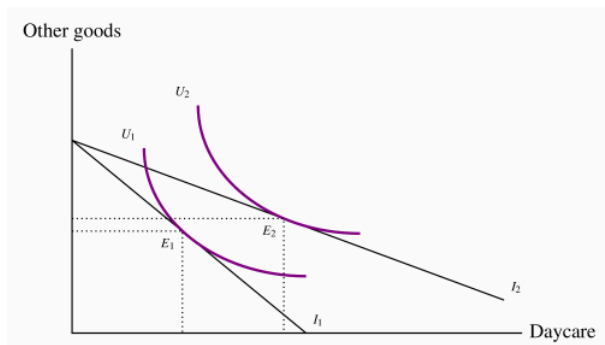
An increase in income due to a government transfer shifts the budget constraint from I_1 to I_2 . This parallel shift increases the quantity consumed of the target good (daycare) *and* other goods, unless one is inferior.

In Figure 6.11, an *income transfer* increases income from I_1 to I_2 . The new equilibrium at E_2 reflects an increase in utility, and an increase in the consumption of *both* daycare and other goods.

Suppose now that a government program administrator decides that, while helping this individual to purchase more daycare accords with the intent of the transfer, she does not intend that government money should be used to purchase other goods. She therefore decides that a daycare *subsidy* program might better meet this objective than a pure income transfer.

A daycare subsidy reduces the price of daycare and therefore *rotates the budget constraint outwards around the intercept on the vertical axis*. At the equilibrium in Figure 6.12, purchases of other goods change very little, and therefore most of the additional purchasing power is allocated to daycare.

Figure 6.12 Price subsidy



A subsidy to the targeted good, by reducing its price, rotates the budget constraint from I_1 to I_2 . This induces the consumer to direct expenditure more towards daycare and less towards other goods than an income transfer that does not change the relative prices.

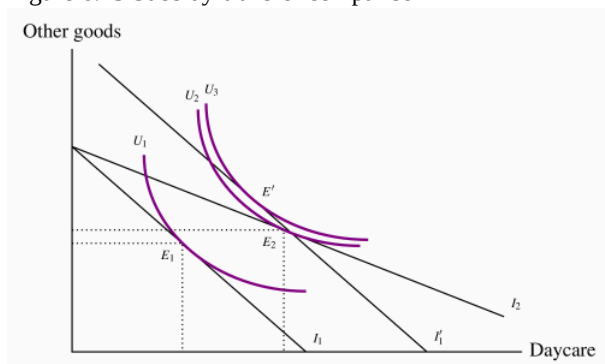
Let us take the example one stage further. From the initial equilibrium E_1 in Figure 6.12, suppose that, instead of a subsidy that took the individual to E_2 , we gave an income transfer *that enabled the consumer to purchase the combination E_2* . Such a transfer is represented in Figure 6.13 by a parallel outward shift of the budget constraint from I_1 to I'_1 , going through the point E_2 . We now have a subsidy policy and an alternative income transfer policy, each permitting the same consumption bundle (E_2). The interesting aspect of this pair of possibilities is that the income transfer will enable the consumer to attain a higher level of satisfaction—for example, at point E' —and will also induce her to consume more of the good on the vertical axis. The higher level of satisfaction comes about because the consumer has more latitude in allocating the additional real income.

Application Box 6.3 Daycare subsidies in Quebec

The Quebec provincial government subsidizes daycare heavily. In the public-sector network called the "Centres de la petite enfance", families can place their children in daycare for less than \$10 per day, while families that use the private sector are permitted a generous tax allowance for their daycare costs. This policy is designed to enable households to limit the share of their income expended on daycare. It is described in Figure 6.13.

The consequences of strong subsidization are not negligible: Excess demand, to such an extent that children are frequently placed on waiting lists for daycare places long before their parents intend to use the service. Annual subsidy costs amount to almost \$2 billion per year. At the same time, it has been estimated that the policy has enabled many more parents to enter the workforce than otherwise would have.

Figure 6.13 Subsidy-transfer comparison



A price subsidy to the targeted good induces the individual to move from E_1 to E_2 , facing a budget constraint I_2 . An income transfer that permits him to consume E_2 is given by I'_1 ; but it also permits him to attain a higher level of satisfaction, denoted by E' on the indifference curve U_3 .

The price of giving

Imagine now that the good on the horizontal axis is charitable donations, rather than daycare, and the government decides that for every dollar given the individual will see a reduction in their income tax of 50 cents. This is equivalent to cutting the 'price' of donations in half, because a donation of one dollar now costs the individual half of that amount. Graphically the budget constraint rotates outward with the vertical intercept unchanged. Since donations now cost less the individual has increased spending power as a result of the price reduction for donations. The price reduction is designed to increase the attractiveness of donations to the utility maximizing consumer.

Key Terms

Cardinal utility is a measurable concept of satisfaction.

Total utility is a measure of the total satisfaction derived from consuming a given amount of goods and services.

Marginal utility is the addition to total utility created when one more unit of a good or service is consumed.

Diminishing marginal utility implies that the addition to total utility from each extra unit of a good or service consumed is declining.

Consumer equilibrium occurs when marginal utility per dollar spent on the last unit of each good is equal.

Law of demand states that, other things being equal, more of a good is demanded the lower is its price.

Ordinal utility assumes that individuals can rank commodity bundles in accordance with the level of satisfaction associated with each bundle.

Budget constraint defines all bundles of goods that the consumer can afford with a given budget.

Affordable set of goods and services for the consumer is bounded by the budget line from above; the **non-affordable set** lies strictly above the budget line.

Indifference curve defines combinations of goods and services that yield the same level of satisfaction to the consumer.

Indifference map is a set of indifference curves, where curves further from the origin denote a higher level of satisfaction.

Marginal rate of substitution is the slope of the indifference curve. It defines the amount of one good the consumer is willing to sacrifice in order to obtain a given increment of the other, while maintaining utility unchanged.

Diminishing marginal rate of substitution reflects a higher marginal value being associated with smaller quantities of any good consumed.

Consumer optimum occurs where the chosen consumption bundle is a point such that the price ratio equals the marginal rate of substitution.

Exercises for Chapter 6

EXERCISE 6.1

In the example given in Table 6.1, suppose Neal experiences a small increase in income. Will he allocate it to snowboarding or jazz? [*Hint: At the existing equilibrium, which activity will yield the higher MU for an additional dollar spent on it?*]

EXERCISE 6.2

Suppose that utility depends on the square root of the amount of good X consumed: $U = \sqrt{X}$.

1. In a spreadsheet enter the values 1... 16 as the X column (col A), and in the adjoining column (B) compute the value of utility corresponding to each quantity of X . To do this use the 'SQRT' command. For example, the entry in cell B3 will be of the form '=SQRT(A3)'.
2. In the third column enter the marginal utility (MU) associated with each value of X – the change in utility in going from one value of X to the next.
3. Use the 'graph' tool to map the relationship between U and X .
4. Use the graph tool to map the relationship between MU and X .

EXERCISE 6.3

Instead of the square-root utility function in Exercise 6.2, suppose that utility takes the form $U=x^2$.

1. Follow the same procedure as in the previous question – graph the utility function.
2. Why is this utility function not consistent with our beliefs on utility?

EXERCISE 6.4

1. Plot the utility function $U=2X$, following the same procedure as in the previous questions.
2. Next plot the marginal utility values in a graph. What do we notice about the behaviour of the MU ?

EXERCISE 6.5

Let us see if we can draw a utility function for beer. In this instance the individual may reach a point where he takes too much.

1. If the utility function is of the form $U=6X-X^2$, plot the utility values for X values in the range 1...8, using either a spreadsheet or manual calculations.
2. At how many units of X (beer) is the individual's utility maximized?
3. At how many beers does the utility become negative?

EXERCISE 6.6

Cappuccinos, C , cost \$3 each, and music downloads of your favourite artist, M , cost \$1 each from your iTunes store. Income is \$24.

1. Draw the budget line, with cappuccinos on the vertical axis, and music on the horizontal axis, and compute the values of the intercepts.
2. What is the slope of the budget constraint, and what is the opportunity cost of 1 cappuccino?
3. Are the following combinations of goods in the affordable set: (4C and 9M), (6C and 2M), (3C and 15M)?
4. Which combination(s) above lie inside the affordable set, and which lie on the boundary?

EXERCISE 6.7

George spends his income on gasoline and "other goods."

1. First, draw a budget constraint, with gasoline on the horizontal axis.
2. Suppose now that, in response to a gasoline shortage in the economy, the government imposes a *ration* on each individual that limits the purchase of gasoline to an amount less than the gasoline intercept of the budget constraint. Draw the new effective budget constraint.

EXERCISE 6.8

Suppose that you are told that the indifference curves defining the trade-off for two goods took the form of straight lines. Which of the four properties outlines in Section 6.3 would such indifference curves violate?

EXERCISE 6.9

Draw an indifference map with several indifference curves and several budget constraints corresponding to different possible levels of income. Note that these budget constraints should all be parallel because only income changes, not prices. Now find some optimizing (tangency) points. Join all of these points. You have just constructed what is called an income-consumption curve. Can you understand why it is called an income-consumption curve?

EXERCISE 6.10

Draw an indifference map again, in conjunction with a set of budget constraints. This time the budget constraints should each have a different price of good X and the same price for good Y .

1. Draw in the resulting equilibria or tangencies and join up all of these points. You have just constructed a price-consumption curve for good X . Can you understand why the curve is so called?
2. Now repeat part (a), but keep the price of X constant and permit the price of Y to vary. The resulting set of equilibrium points will form a price consumption curve for good Y .

EXERCISE 6.11

Suppose that movies are a normal good, but public transport is inferior. Draw an indifference map with a budget constraint and initial equilibrium. Now let income increase and draw a plausible new equilibrium, noting that one of the goods is inferior.

1. Note that, with two snowboard outings and seven jazz club visits, total utility is 352 ($=132+220$), while the optimal combination of four of each yields a total utility of 380 ($=224+156$).

This page titled [6: Individual choice](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis](#) and [Ian Irvine](#) (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

6.1: Rationality

A critical behavioural assumption in economics is that agents operate in a way that is oriented towards achieving a goal. This can be interpreted to mean that individuals and firms maximize their personal well-being and/or their profits. These players may have other goals in addition: Philanthropy and the well-being of others are consistent with individual optimization.

If individuals are to achieve their goals then they must act in a manner that will get them to their objective; broadly, they must act in a rational manner. The theory of individual maximization that we will develop in this chapter is based on that premise or assumption. In assuming individuals are rational we need not assume that they have every piece of information available to them that might be relevant for a specific decision or choice. Nor need we assume that they have super computers in their brain when they evaluate alternative possible strategies.

What we do need to assume, however, is that individuals act in a manner that is consistent with obtaining a given objective. The modern theory of behavioural economics and behavioural psychology examines decision making in a wide range of circumstances and has uncovered many fascinating behaviours – some of which are developed in Application Box 6.1 below.

We indicated in Chapter 1 that as social scientists, we require a *reliable model* of behaviour, that is, *a way of describing the essentials of choice that is consistent with everyday observations on individual behaviour patterns*. In this chapter, our aim is to understand more fully the behavioural forces that drive the demand side of the economy.

Economists analyze individual decision making using two different, yet complementary, approaches – utility analysis and indifference analysis. We begin by portraying individuals as maximizing their *measurable utility* (sometimes called *cardinal utility*); then progress to indifference analysis, where a weaker assumption is made on the ability of individuals to measure their satisfaction. In this second instance we do not assume that individuals can measure their utility numerically, only that they can say if one collection of goods and services yields them greater satisfaction than another group. This ranking of choices corresponds to what is sometimes called *ordinal utility* – because individuals can *order* groups of goods and services in ascending order of satisfaction. In each case individuals are perceived as rational maximizers or optimizers: They allocate their income so as to choose the outcome that will make them as well off as possible.

The second approach to consumer behaviour is frequently omitted in introductory texts. It can be omitted here without interpreting the flow of ideas, although it does yield additional insights into consumer choice and government policy. As in preceding chapters, we begin the analysis with a motivating numerical example.

Application Box 6.1 Rationality and impulse

A number of informative and popular books on decision making have appeared recently. Their central theme is that our decision processes should not be viewed solely as a rational computer – operating in one single mode only, and unmoved by our emotions or history. Psychologists now know that our brains have at least two decision modes, and these are developed by economics Nobel Prize winner Daniel Kahneman in his book "Thinking, Fast and Slow". One part of our brain operates in a rational goal-oriented forward-looking manner (the 'slow' part), another is motivated by immediate gratification (the 'fast' part). Decisions that we observe in the world about us reflect these different mechanisms.

Richard Thaler, a Chicago economist and his law professor colleague Cass Sunstein, have developed a role for public policy in their book entitled "Nudge". They too argue that individuals do not inevitably operate in their own best long-term interests, and as a consequence individuals frequently require a *nudge* by government to make the long-term choice rather than the short-term choice. For example, when individuals begin a new job, they might be automatically enrolled in the company pension plan and be given the freedom to opt out, rather than not be enrolled and given the choice to opt in. Such policies are deemed to be 'soft paternalism'. They are paternalistic for the obvious reason – another organism is directing, but they are also soft in that they are not binding.

This page titled [6.1: Rationality](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

6.2: Choice with measurable utility

Neal loves to pump his way through the high-altitude powder at the Whistler ski and snowboard resort. His student-rate lift-ticket cost is \$30 per visit. He also loves to frequent the jazz bars in downtown Vancouver, and each such visit costs him \$20. With expensive passions, Neal must allocate his monthly entertainment budget carefully. He has evaluated how much satisfaction, measured in utils, he obtains from each snowboard outing and each jazz club visit. We assume that these utils are measurable, and use the term cardinal utility to denote this. These measurable utility values are listed in columns 2 and 3 of Table 6.1. They define the total utility he gets from various amounts of the two activities.

Table 6.1 Utils from snowboarding and jazz

1	2	3	4	5	6	7
Visit	Total	Total	Marginal	Marginal	Marginal	Marginal
#	snowboard	jazz	snowboard	jazz utils	snowboard	jazz utils
	utils	utils	utils		utils per \$	per \$
1	72	52	72	52	2.4	2.6
2	132	94	60	42	2.0	2.1
3	182	128	50	34	1.67	1.7
4	224	156	42	28	1.4	1.4
5	260	180	36	24	1.2	1.2
6	292	201	32	21	1.07	1.05
7	321	220	29	19	0.97	0.95

Price of snowboard visit=\$30. Price of jazz club visit=\$20.

Cardinal utility is a measurable concept of satisfaction.

Total utility is a measure of the total satisfaction derived from consuming a given amount of goods and services.

Neal's total utility from each activity in this example is independent of the amount of the other activity he engages in. These total utilities are plotted in Figures 6.1 and 6.2. Clearly, more of each activity yields more utility, so the additional or marginal utility (*MU*) of each activity is positive. This positive marginal utility for any amount of the good consumed, no matter how much, reflects the assumption of *non-satiation*—more is always better. Note, however, that the decreasing slopes of the total utility curves show that *total utility is increasing at a diminishing rate*. While more is certainly better, each additional visit to Whistler or a jazz club augments Neal's utility by a smaller amount. At the margin, his additional utility declines: He has diminishing marginal utility. The marginal utilities associated with snowboarding and jazz are entered in columns 4 and 5 of Table 6.1. They are the differences in total utility values when consumption increases by one unit. For example, when Neal makes a sixth visit to Whistler his total utility increases from 260 utils to 292 utils. His marginal utility for the sixth unit is therefore 32 utils, as defined in column 4. In light of this example, it should be clear that we can define marginal utility as:

$\text{Marginal Utility} = \frac{\text{additional utility}}{\text{additional consumption}} \text{ or, } MU = \frac{\Delta U}{\Delta C},$	(6.1)
--	-------

where ΔC denotes the change in the quantity consumed of the good or service in question.

Marginal utility is the addition to total utility created when one more unit of a good or service is consumed.

Diminishing marginal utility implies that the addition to total utility from each extra unit of a good or service consumed is declining.

Figure 6.1 TU from snowboarding

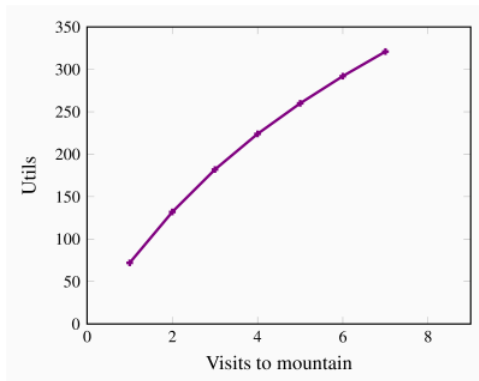


Figure 6.2 TU from jazz

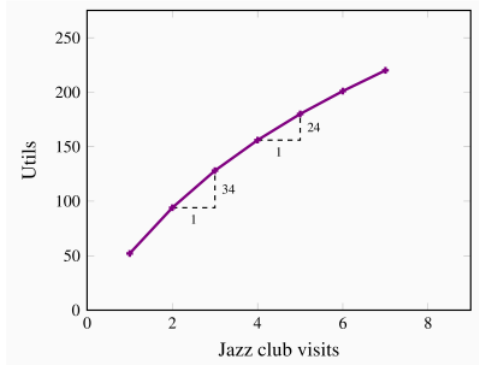


Figure 6.3 MU from snowboarding

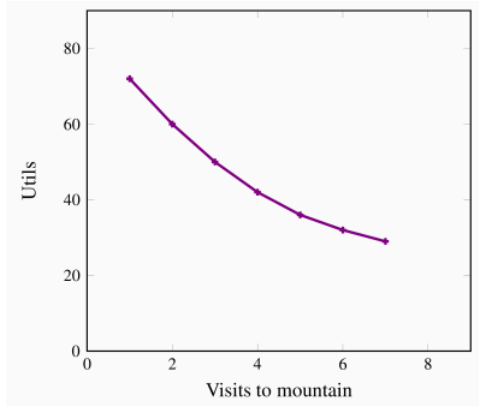
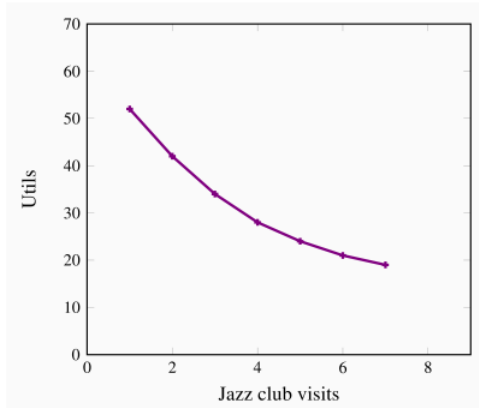


Figure 6.4 MU from jazz



The marginal utilities associated with consuming different amounts of the two goods are plotted in Figures 6.3 and 6.4, using the data from columns 4 and 5 in Table 6.1. These functions are declining, as indicated by their negative slope. It should also be clear

that the MU curves can be derived from the TU curves. For example, in figure 6.2, when going from 2 units to 3 units of Jazz, TU increases by 34 units. But $34/1$ is the slope of the TU function in this range of consumption – the vertical distance divided by the horizontal distance. Similarly, if jazz consumption increases from 4 units to 5 units the corresponding change in TU is 24 units, again the vertical distance divided by the horizontal distance, and so the slope of the function. In short, the MU is the slope of the TU function.

Now that Neal has defined his utility schedules, he must consider the price of each activity. Ultimately, when deciding how to allocate his monthly entertainment budget, he must evaluate how much utility he gets from each dollar spent on snowboarding and jazz: What "bang for his buck" does he get? Let us see how he might go about allocating his budget. *When he has fully spent his budget in the manner that will yield him greatest utility, we say that he has attained equilibrium*, because he will have no incentive to change his expenditure patterns.

If he boards once, at a cost of \$30, he gets 72 utils of satisfaction, which is 2.4 utils per dollar spent ($=72/30$). One visit to a jazz club would yield him 2.6 utils per dollar ($=52/20$). Initially, therefore, his dollars give him more *utility per dollar* when spent on jazz. His MU per dollar spent on each activity is given in the final two columns of the table. These values are obtained by dividing the MU associated with each additional unit by the good's price.

We will assume that Neal has a budget of \$200. He realizes that his initial expenditure should be on a jazz club visit, because he gets more utility per dollar spent there. Having made one such expenditure, he sees that a second jazz outing would yield him 2.1 utils per dollar expended, while a first visit to Whistler would yield him 2.4 utils per dollar. Accordingly, his second activity is a snowboard outing.

Having made one jazz and one snowboarding visit, he then decides upon a second jazz club visit for the same reason—as utility value for his money. He continues to allocate his budget in this way until his budget is exhausted. In our example, this occurs when he spends \$120 on four snowboarding outings and \$80 on four jazz club visits. At this consumer equilibrium, he gets the same utility value per dollar for the last unit of each activity consumed. This is a necessary condition for him to be maximizing his utility, that is, to be in equilibrium.

Consumer equilibrium occurs when marginal utility per dollar spent on the last unit of each good is equal.

To be absolutely convinced of this, imagine that Neal had chosen instead to board twice and to visit the jazz clubs seven times; this combination would also exhaust his \$200 budget exactly. With such an allocation, he would get 2.0 utils per dollar spent on his marginal (second) snowboard outing, but just 0.95 utils per dollar spent on his marginal (seventh) jazz club visit.¹ If, instead, he were to reallocate his budget in favour of snowboarding, he would get 1.67 utils per dollar spent on a third visit to the hills. By reducing the number of jazz visits by one, he would lose 0.95 utils per dollar reallocated. Consequently, the utility gain from a reallocation of his budget towards snowboarding would outweigh the utility loss from allocating fewer dollars to jazz. His initial allocation, therefore, was not an optimum, or equilibrium.

Only when the utility per dollar expended on each activity is equal at the margin will Neal be optimizing. When that condition holds, a reallocation would be of no benefit to him, because the gains from one more dollar on boarding would be exactly offset by the loss from one dollar less spent on jazz. Therefore, we can write the equilibrium condition as

Equilibrium requires: $\frac{MU_s}{P_s} = \frac{MU_j}{P_j}$ or $\frac{MU_s}{MU_j} = \frac{P_s}{P_j}$	(6.2)
--	-------

While this example has just two goods, in the more general case of many goods, this same *condition must hold for all pairs of goods* on which the consumer allocates his or her budget.

From utility to demand

Utility theory is a useful way of analyzing how a consumer makes choices. But in the real world we do not observe a consumer's utility, either total or marginal. Instead, his or her behaviour in the marketplace is observed through the demand curve. How are utility and demand related?

Demand functions relate the quantity of a good consumed to the price of that good, other things being equal. So let us trace out the effects of a price change on demand, with the help of this utility framework. We will introduce a simplification here: Goods are divisible, or that they come in small packages relative to income. Think, for example, of kilometres driven per year, or liters of gasoline purchased. Conceptualizing things in this way enables us to imagine more easily experiments in which small amounts of a

budget are allocated one dollar at a time. In contrast, in the snowboard/jazz example, we had to reallocate the budget in lumps of \$30 or \$20 at a time because we could not "fractionalize" these goods.

The effects of a price change on a consumer's demand can be seen through the condition that describes his or her equilibrium. If income is allocated to, say, three goods $\{a, b, c\}$, such that $MU_a/P_a = MU_b/P_b = MU_c/P_c$, and the price of, say, good b falls, the consumer must reallocate the budget so that once again the MU s per dollar spent are all equated. How does he do this? Clearly, if he purchases more or less of any one good, the MU changes. If the price of good b falls, then the consumer initially gets more utility from good b for the last dollar he spends on it (the denominator in the expression MU_b/P_b falls, and consequently the value of the ratio rises to a value greater than the values for goods a and c).

The consumer responds to this, in the first instance, by buying more of the cheaper good. He obtains more total utility as a consequence, and in the process will get *less utility at the margin* from that good. In essence, the numerator in the expression then falls, in order to realign it with the lower price. This equality also provides an underpinning for what is called the law of demand: More of a good is demanded at a lower price. If the price of any good falls, then, in order for the equilibrium condition to be re-established, the MU of that good must be driven down also. Since MU declines when more is purchased, this establishes that demand curves must slope downwards.

The **law of demand** states that, other things being equal, more of a good is demanded the lower is its price.

However, the effects of a price decline are normally more widespread than this, because the quantities of other goods consumed may also change. As explained in earlier chapters, the decline in the price of good b will lead the consumer to purchase more units of *complementary goods* and fewer units of goods that are *substitutes*. So the whole budget allocation process must be redetermined in response to any price change. But at the end of the day, a new equilibrium must be one where the marginal utility per dollar spent on each good is equal.

Applying the theory

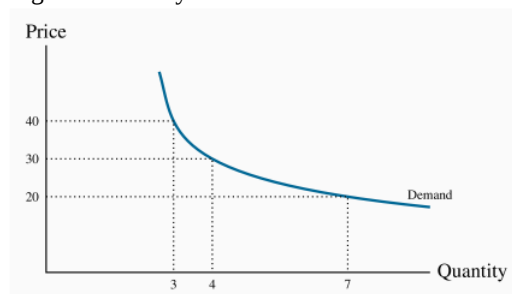
The demand curves developed in Chapter 3 can be related to the foregoing utility analysis. In our example, Neal purchased four lift tickets at Whistler when the price was \$30. We can think of this combination as one point on his demand curve, where the "other things kept constant" are the price of jazz, his income, his tastes, etc.

Suppose now that the price of a lift ticket increased to \$40. How could we find another point on his demand curve corresponding to this price, using the information in Table 6.1? The marginal utility per dollar associated with each visit to Whistler could be recomputed by dividing the values in column 4 by 40 rather than 30, yielding a new column 6. We would then determine a new allocation of his budget between the two goods that would maximize utility. After such a calculation we would find that he makes three visits to Whistler and four jazz-club visits. Thus, the combination $(P_s = \$40, Q_s = 3)$ is another point on his demand curve. Note that this allocation exactly exhausts his \$200 budget.

By setting the price equal to \$20, this exercise could be performed again, and the outcome will be a quantity demanded of lift tickets equal to seven (plus three jazz club visits). Thus, the combination $(P_s = \$20, Q_s = 7)$ is another point on his demand curve. Figure 6.5 plots a demand curve going through these three points.

By repeating this exercise for many different prices, the demand curve is established. We have now linked the demand curve to utility theory.

Figure 6.5 Utility to demand



When $P = \$30$, the consumer finds the quantity such that MU/P is equal for all purchases. The corresponding quantity purchased is 4 tickets. At prices of \$40 and \$20 the equilibrium condition implies quantities of 3 and 7 respectively.

Application Box 6.2 Individual and Collective Utility

The example developed in the text is not far removed from what economists do in practice. From a philosophical standpoint, economists are supposed to be interested in the well-being of the citizens who make up an economy or a country. To determine how 'well-off' citizens may be, social scientists frequently carry out surveys on how 'content' or 'happy' people are in their every-day lives. For example, the *Earth Institute at Columbia University* regularly produces a 'World Happiness Report'. The report is based upon responses to survey questions in numerous economies. One of the measures it uses to compare utility levels is the *Cantril ladder*. This is an 11-point scale running from 0 to 10, with the lowest value signifying the worst possible life, and 10 the highest possible quality of life. In reporting their findings, the researchers are essentially claiming that some economies have, on average, more contented or happier, people than others. Utility can be considered in exactly this way: A higher reported value on the Cantril ladder suggests higher utility.

A slightly different measure of well-being across economies is given by the *United Nations Human Development Index*. In this case, countries score high by having a high level of income, good health (as measured by life expectancy), and high levels of education, as measured by the number of years of education completed or envisaged.

In practice, social scientists are very comfortable using utility-based concepts to describe the economic circumstances of individuals in different economies.

This page titled [6.2: Choice with measurable utility](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

6.3: Choice with ordinal utility

The budget constraint

In the preceding section, we assumed that utility is measurable in order to better understand how consumers allocate their budgets, and how this process is reflected in the market demands that are observed. The belief that utility might be measurable is not too extreme in the modern era. Neuroscientists are mapping more and more of the human brain and understanding how it responds to positive and negative stimuli. At the same time, numerous sociological surveys throughout the world ask individuals to rank their happiness on a scale of one to ten, or something similar, with a view to making comparisons between individual-level and group-level happiness – see Application Box 6.2. Nonetheless, not every scientist may be convinced that we should formulate behavioural rules on this basis. Accordingly we now examine the economics of consumer behaviour without this strong assumption. We assume instead that individuals are able to identify (a) different combinations of goods and services that yield equal satisfaction, and (b) combinations of goods and services that yield more satisfaction than other combinations. In contrast to measurable (or cardinal) utility, this concept is called ordinal utility, because it assumes only that consumers can order utility bundles rather than quantify the utility.

Ordinal utility assumes that individuals can rank commodity bundles in accordance with the level of satisfaction associated with each bundle.

The budget constraint

Neal's monthly expenditure limit, or budget constraint, is \$200. In addition, he faces a price of \$30 for lift tickets and \$20 per visit to jazz clubs. Therefore, using S to denote the number of snowboard outings and J the number of jazz club visits, if he spends his entire budget it must be true that the sum of expenditures on each activity exhausts his budget or income (I):

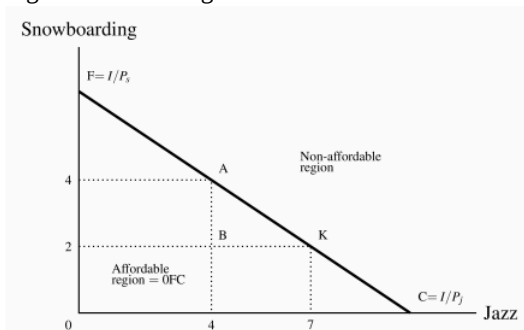
Expenditure on snowboarding + expenditure on Jazz	= Income
(Price of $S \times$ quantity of S) + (price of $J \times$ quantity of J)	= Income
$P_S S + P_J J = I$ or $\$30S + \$20J$	= \$200

Since many different combinations of the two goods are affordable, it follows that the budget constraint defines all bundles of goods that the consumer can afford with a given budget.

The **budget constraint** defines all bundles of goods that the consumer can afford with a given budget.

The budget constraint, then, is just what it claims to be—a limit on behaviour. Neal's budget constraint is illustrated in Figure 6.6, where the amount of each good consumed is given on the axes. If he spends all of his \$200 income on jazz, he can make exactly ten jazz club visits ($\$200/\$20 = 10$). The calculation also applies to visits to Whistler. The intercept value is always obtained by dividing income by the price of the good or activity in question.

Figure 6.6 The budget line



FC is the budget constraint and defines the affordable combinations of snowboarding and jazz. F represents all income spent on snowboarding. Thus $F = I/P_S$. Similarly $C = I/P_J$. Points above FC are not attainable. The slope = $OF/OC = (I/P_S)/(I/P_J) = P_J/P_S = 20/30 = 2/3$. The affordable set is OFC.

In addition to these affordable extremes, Neal can also afford many other bundles, e.g., ($S=2, J=7$), or ($S=4, J=4$), or ($S=6, J=1$). The set of feasible, or affordable, combinations is bounded by the budget line, and this is illustrated in Figure 6.6.

The **affordable set** of goods and services for the consumer is bounded by the budget line from above; the **non-affordable set** lies strictly above the budget line.

The slope of the budget line is informative. As illustrated in Chapter 1, it indicates how many snowboard visits must be sacrificed for one additional jazz visit; it defines the consumer's *trade-offs*. To illustrate: Suppose Neal is initially at point A ($J=4, S=4$), and moves to point K ($J=7, S=2$). Clearly, both points are affordable. In making the move, he trades two snowboard outings in order to get three additional jazz club visits, a trade-off of $2/3$. This trade-off is the slope of the budget line, which, in Figure 6.6, is $AB/BK=-2/3$, where the negative sign reflects the downward slope.

Could it be that this ratio reflects the two prices ($\$20/\30)? The answer is yes: The slope of the budget line is given by the vertical distance divided by the horizontal distance, OF/OC . The points F and C were obtained by dividing income by the respective price—remember that the jazz intercept is $\$200/\$20 = 10$. Formally, that is I/P_J . The intercept on the snowboard axis is likewise I/P_S . Accordingly, the slope of the budget constraint is:

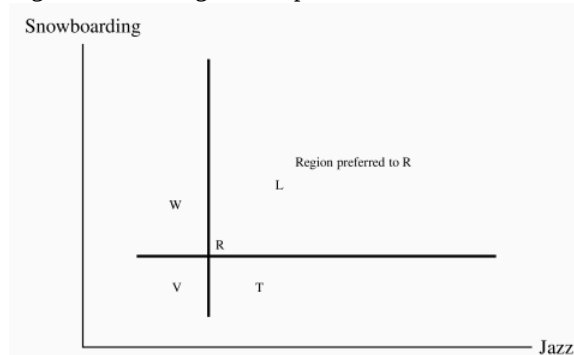
$$\text{Slope} = OF/OC = \frac{(I/P_S)}{(I/P_J)} = \frac{I}{P_S} \times \frac{P_J}{I} = \frac{P_J}{P_S}.$$

Since the budget line has a negative slope, it is technically correct to define it with a negative sign. But, as with elasticities, the sign is frequently omitted.

Tastes and indifference

We now consider how to represent a consumer's tastes in two dimensions, given that he can order, or rank, different consumption bundles, and that he can define a series of different bundles that all yield the same satisfaction. We limit ourselves initially to considering just "goods," and not "bads" such as pollution.

Figure 6.7 Ranking consumption bundles



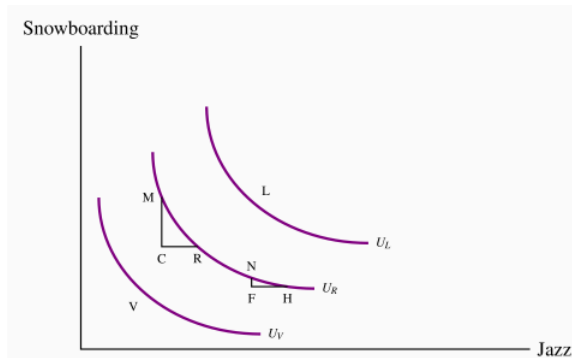
L is preferred to R since more of each good is consumed at L, while points such as V are less preferred than R. Points W and T contain more of one good and less of the other than R. Consequently, we cannot say if they are preferred to R without knowing how the consumer trades the goods off—that is, his preferences.

Figure 6.7 examines the implications of these assumptions about tastes. Each point shows a consumption bundle of snowboarding and jazz. Let us begin at bundle R. Since more of a good is preferred to less, any point such as L, which lies to the northeast of R, is preferred to R, since L offers more of both goods than R. Conversely, points to the southwest of R offer *less of each good* than R, and therefore R is preferred to a point such as V.

Without knowing the consumer's tastes, we cannot be sure at this stage how points in the northwest and southeast regions compare with R. At W or T, the consumer has more of one good and less of the other than at R. Someone who really likes snowboarding might prefer W to R, but a jazz buff might prefer T to R.

Let us now ask Neal to disclose his tastes, by asking him to define several combinations of snowboarding and jazz that *yield him exactly the same degree of satisfaction as the combination at R*. Suppose further, for reasons we shall understand shortly, that his answers define a series of points that lie on the beautifully smooth contour U_R in Figure 6.8. Since he is indifferent between all points on U_R by construction, this contour is an indifference curve.

Figure 6.8 Indifference curves



An indifference curve defines a series of consumption bundles, all of which yield the same satisfaction. The slope of an indifference curve is the marginal rate of substitution (MRS) and defines the number of units of the good on the vertical axis that the individual will trade for one unit of the good on the horizontal axis. The MRS declines as we move south-easterly, because the consumer values the good more highly when he has less of it.

An **indifference curve** defines combinations of goods and services that yield the same level of satisfaction to the consumer.

Pursuing this experiment, we could take other points in Figure 6.8, such as L and V, and ask the consumer to define bundles that would yield the same level of satisfaction, or indifference. These combinations would yield additional contours, such as U_L and U_V in Figure 6.8. This process yields a series of indifference curves that together form an indifference map.

An **indifference map** is a set of indifference curves, where curves further from the origin denote a higher level of satisfaction.

Let us now explore the properties of this map, and thereby understand why the contours have their smooth convex shape. They have four properties. The first three follow from our preceding discussion, and the fourth requires investigation.

1. Indifference curves *further from the origin reflect higher levels of satisfaction*.
2. Indifference curves are *negatively sloped*. This reflects the fact that if a consumer gets more of one good she should have less of the other in order to remain indifferent between the two combinations.
3. Indifference curves *cannot intersect*. If two curves were to intersect at a given point, then we would have two different levels of satisfaction being associated with the same commodity bundle—an impossibility.
4. Indifference curves are convex when viewed from the origin, reflecting a *diminishing marginal rate of substitution*.

The convex shape reflects an important characteristic of preferences: When consumers have a lot of some good, they value a marginal unit of it less than when they have a small amount of that good. More formally, they have a *higher marginal valuation at low consumption levels*—that first cup of coffee in the morning provides greater satisfaction than the second or third cup.

Consider the various points on U_R , starting at M in Figure 6.8. At M Neal snowboards a lot; at N he boards much less. The convex shape of his indifference map shows that he values a marginal snowboard trip more at N than at M. To see this, consider what happens as he moves along his indifference curve, starting at M. We have chosen the coordinates on U_R so that, in moving from M to R, and again from N to H, the additional amount of jazz is the same: $CR = FH$. From M, if Neal moves to R, he consumes an additional amount of jazz, CR. By definition of the indifference curve, he is willing to give up MC snowboard outings. The ratio MC/CR defines his willingness to substitute one good for the other. This ratio, being a vertical distance divided by a horizontal distance, is the slope of the indifference curve and is called the marginal rate of substitution, MRS .

The **marginal rate of substitution** is the slope of the indifference curve. It defines the amount of one good the consumer is willing to sacrifice in order to obtain a given increment of the other, while maintaining utility unchanged.

At N, the consumer is willing to sacrifice the amount NF of boarding to get the same additional amount of jazz. Note that, when he boards *less*, as at N, he is willing to give up less boarding than when he has a lot of it, as at M, in order to get the same additional amount of jazz. His willingness to substitute *diminishes* as he moves from M to N: The quantity NF is less than the quantity MC. In order to reflect this taste characteristic, the indifference curve has a diminishing marginal rate of substitution: A flatter slope as we move down along its surface.

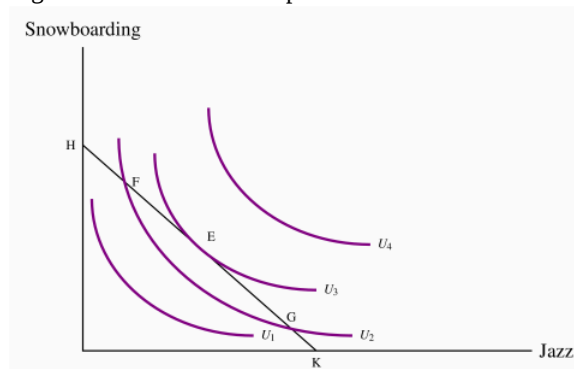
A **diminishing marginal rate of substitution** reflects a higher marginal value being associated with smaller quantities of any good consumed.

Optimization

We are now in a position to examine how the consumer optimizes—how he gets to the highest level of satisfaction possible. The constraint on his behaviour is the affordable set defined in Figure 6.6, the budget line.

Figure 6.9 displays several of Neal's indifference curves in conjunction with his budget constraint. We propose that he maximizes his utility, or satisfaction, at the point E, on the indifference curve denoted by U_3 . While points such as F and G are also on the boundary of the affordable set, they do not yield as much satisfaction as E, because E lies on a higher indifference curve. *The highest possible level of satisfaction is attained, therefore, when the budget line touches an indifference curve at just a single point—that is, where the constraint is tangent to the indifference curve.* E is such a point.

Figure 6.9 The consumer optimum



The budget constraint constrains the individual to points on or below HK. The highest level of satisfaction attainable is U_3 , where the budget constraint just touches, or is just tangent to, it. At this optimum the slope of the budget constraint ($-P_J/P_S$) equals the MRS .

This tangency between the budget constraint and an indifference curve requires that the slopes of each be the same at the point of tangency. We have already established that the slope of the budget constraint is the negative of the price ratio ($= -P_J/P_S$). The slope of the indifference curve is the marginal rate of substitution MRS . It follows, therefore, that the consumer optimizes where the marginal rate of substitution equals the slope of the price line.

Optimization requires:

Slope of Indifference curve = marginal rate of substitution = $-\frac{P_J}{P_S}$	(6.3)
--	-------

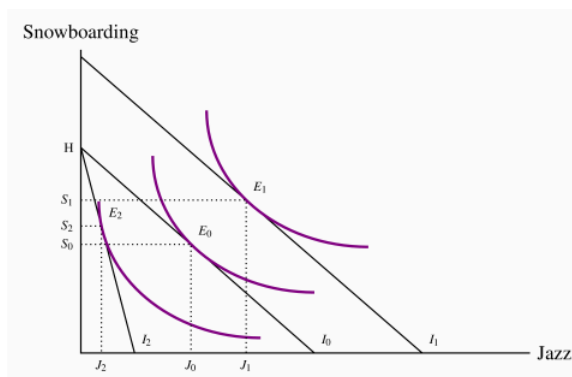
A **consumer optimum** occurs where the chosen consumption bundle is a point such that the price ratio equals the marginal rate of substitution.

Notice the resemblance between this condition and the one derived in the first section as Equation 6.2. There we argued that equilibrium requires the ratio of the marginal utilities be same as the ratio of prices. Here we show that the MRS must equal the ratio of prices. In fact, with a little mathematics it can be shown that the MRS is indeed the same as the (negative of the) ratio of the marginal utilities: $MRS = -MU_J/MU_S$. Therefore the two conditions are in essence the same! However, it was not necessary to assume that an individual can actually measure his utility in obtaining the result that the MRS should equal the price ratio in equilibrium. The concept of ordinal utility is sufficient.

Adjusting to income changes

Suppose now that Neal's income changes from \$200 to \$300. How will this affect his consumption decisions? In Figure 6.10, this change is reflected in a *parallel* outward shift of the budget constraint. Since no price change occurs, the slope remains constant. By recomputing the ratio of income to price for each activity, we find that the new snowboard and jazz intercepts are 10 ($= \$300/\30) and 15 ($= \$300/\20), respectively. Clearly, the consumer can attain a higher level of satisfaction—at a new tangency to a higher indifference curve—as a result of the size of the affordable set being expanded. In Figure 6.10, the new equilibrium is at E_1 .

Figure 6.10 Income and price adjustments



An income increase shifts the budget constraint from I_0 to I_1 . This enables the consumer to attain a higher indifference curve. A price rise in jazz tickets rotates the budget line I_0 inwards around the snowboard intercept to I_2 . The price rise reflects a lower real value of income and results in a lower equilibrium level of satisfaction.

Adjusting to price changes

Next, consider the impact of a price change from the initial equilibrium E_0 in Figure 6.10. Suppose that jazz now costs more. This reduces the purchasing power of the given budget of \$200. The new jazz intercept is therefore reduced. The budget constraint becomes steeper and rotates around the snowboard intercept H , which is unchanged because its price is constant. The new equilibrium is at E_2 , which reflects a lower level of satisfaction because the affordable set has been reduced by the price increase. As explained in Section 6.2, E_0 and E_2 define points on the demand curve for jazz (J_0 and J_2): They reflect the consumer response to a change in the price of jazz with all other things held constant. In contrast, the price increase for jazz *shifts* the demand curve for snowboarding: As far as the demand curve for snowboarding is concerned, a change in the price of jazz is one of those things other than own-price that determine its position.

Philanthropy

Individuals in the foregoing analysis aim to maximize their utility, given that they have a fixed budget. Note that this behavioural assumption does not rule out the possibility that these same individuals may be philanthropic – that is, they get utility from the act of giving to their favourite charity or the United Way or Centre-aide. To see this suppose that donations give utility to the individual in question – she gets a 'warm glow' feeling as a result of giving, which is to say she gets utility from the activity. There is no reason why we cannot put charitable donations on one axis and some other good or combination of goods on the remaining axis. At equilibrium, the marginal utility per dollar of contributions to charity should equal the marginal utility per dollar of expenditure on other goods; or, stated in terms of ordinal utility, the marginal rate of substitution between philanthropy and any other good should equal the ratio of their prices. Evidently the price of a dollar of charitable donations is one dollar.

This page titled 6.3: Choice with ordinal utility is shared under a CC BY-NC-SA 4.0 license and was authored, remixed, and/or curated by Douglas Curtis and Ian Irvine (Lyryx) via source content that was edited to the style and standards of the LibreTexts platform.

6.4: Applications of indifference analysis

Price impacts: Complements and substitutes

The nature of complements and substitutes, defined in Chapter 4, can be further understood with the help of Figure 6.10. The new equilibrium E_2 has been drawn so that the increase in the price of jazz results in more snowboarding—the quantity of S increases to S_2 from S_0 . These goods are substitutes in this picture, because snowboarding *increases* in response to an *increase* in the price of jazz. If the new equilibrium E_2 were at a point yielding a lower level of S than S_0 , we would conclude that they were complements.

Cross-price elasticities

Continuing with the same price increase in jazz, we could compute the *percentage* change in the quantity of snowboarding demanded as a result of the *percentage* change in the jazz price. In this example, the result would be a positive elasticity value, because the quantity change in snowboarding and the price change in jazz are both in the same direction, each being positive.

Income impacts: Normal and inferior goods

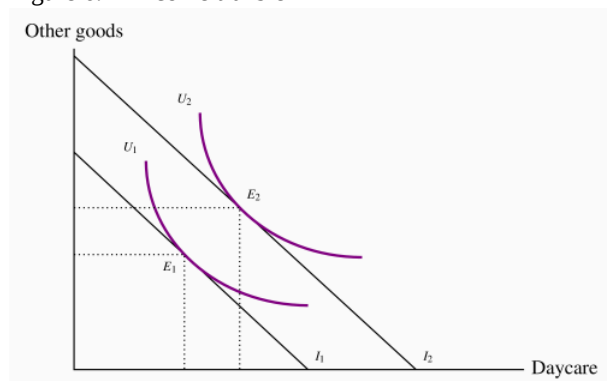
We know from Chapter 4 that the quantity demanded of a *normal good* increases in response to an income increase, whereas the quantity demanded of an *inferior good* declines. Clearly, both jazz and boarding are normal goods, as illustrated in Figure 6.10, because more of each one is demanded in response to the income increase from I_0 to I_1 . It would challenge the imagination to think that either of these goods might be inferior. But if J were to denote junky (inferior) goods and S super goods, we could envisage an equilibrium E_1 to the northwest of E_0 in response to an income increase, along the constraint I_1 ; less J and more S would be consumed in response to the income increase.

Policy: Income transfers and price subsidies

Government policies that improve the purchasing power of low-income households come in two main forms: Pure income transfers and price subsidies. *Social Assistance* payments ("welfare") or *Employment Insurance* benefits, for example, provide an increase in income to the needy. Subsidies, on the other hand, enable individuals to purchase particular goods or services at a lower price—for example, rent or daycare subsidies.

In contrast to taxes, which *reduce* the purchasing power of the consumer, subsidies and income transfers *increase* purchasing power. The impact of an income transfer, compared with a pure price subsidy, can be analyzed using Figures 6.11 and 6.12.

Figure 6.11 Income transfer



An increase in income due to a government transfer shifts the budget constraint from I_1 to I_2 . This parallel shift increases the quantity consumed of the target good (daycare) *and* other goods, unless one is inferior.

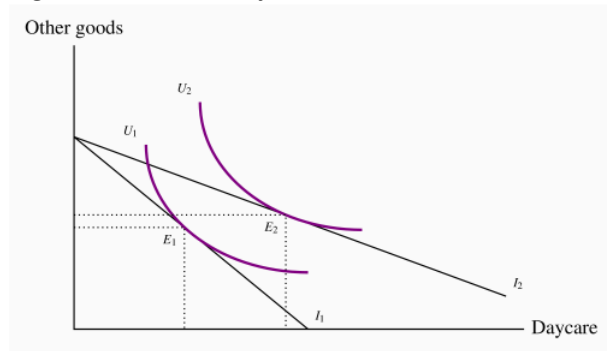
In Figure 6.11, an *income transfer* increases income from I_1 to I_2 . The new equilibrium at E_2 reflects an increase in utility, and an increase in the consumption of *both* daycare and other goods.

Suppose now that a government program administrator decides that, while helping this individual to purchase more daycare accords with the intent of the transfer, she does not intend that government money should be used to purchase other goods. She therefore decides that a daycare *subsidy* program might better meet this objective than a pure income transfer.

A daycare subsidy reduces the price of daycare and therefore *rotates the budget constraint outwards around the intercept on the vertical axis*. At the equilibrium in Figure 6.12, purchases of other goods change very little, and therefore most of the additional

purchasing power is allocated to daycare.

Figure 6.12 Price subsidy



A subsidy to the targeted good, by reducing its price, rotates the budget constraint from I_1 to I_2 . This induces the consumer to direct expenditure more towards daycare and less towards other goods than an income transfer that does not change the relative prices.

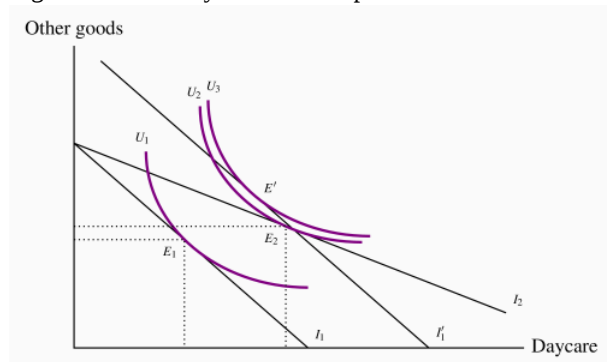
Let us take the example one stage further. From the initial equilibrium E_1 in Figure 6.12, suppose that, instead of a subsidy that took the individual to E_2 , we gave an income transfer *that enabled the consumer to purchase the combination E_2* . Such a transfer is represented in Figure 6.13 by a parallel outward shift of the budget constraint from I_1 to I'_1 , going through the point E_2 . We now have a subsidy policy and an alternative income transfer policy, each permitting the same consumption bundle (E_2). The interesting aspect of this pair of possibilities is that the income transfer will enable the consumer to attain a higher level of satisfaction—for example, at point E' —and will also induce her to consume more of the good on the vertical axis. The higher level of satisfaction comes about because the consumer has more latitude in allocating the additional real income.

Application Box 6.3 Daycare subsidies in Quebec

The Quebec provincial government subsidizes daycare heavily. In the public-sector network called the "Centres de la petite enfance", families can place their children in daycare for less than \$10 per day, while families that use the private sector are permitted a generous tax allowance for their daycare costs. This policy is designed to enable households to limit the share of their income expended on daycare. It is described in Figure 6.13.

The consequences of strong subsidization are not negligible: Excess demand, to such an extent that children are frequently placed on waiting lists for daycare places long before their parents intend to use the service. Annual subsidy costs amount to almost \$2 billion per year. At the same time, it has been estimated that the policy has enabled many more parents to enter the workforce than otherwise would have.

Figure 6.13 Subsidy-transfer comparison



A price subsidy to the targeted good induces the individual to move from E_1 to E_2 , facing a budget constraint I_2 . An income transfer that permits him to consume E_2 is given by I'_1 ; but it also permits him to attain a higher level of satisfaction, denoted by E' on the indifference curve U_3 .

The price of giving

Imagine now that the good on the horizontal axis is charitable donations, rather than daycare, and the government decides that for every dollar given the individual will see a reduction in their income tax of 50 cents. This is equivalent to cutting the 'price' of donations in half, because a donation of one dollar now costs the individual half of that amount. Graphically the budget constraint

rotates outward with the vertical intercept unchanged. Since donations now cost less the individual has increased spending power as a result of the price reduction for donations. The price reduction is designed to increase the attractiveness of donations to the utility maximizing consumer.

This page titled [6.4: Applications of indifference analysis](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

6.5: Key Terms

Cardinal utility is a measurable concept of satisfaction.

Total utility is a measure of the total satisfaction derived from consuming a given amount of goods and services.

Marginal utility is the addition to total utility created when one more unit of a good or service is consumed.

Diminishing marginal utility implies that the addition to total utility from each extra unit of a good or service consumed is declining.

Consumer equilibrium occurs when marginal utility per dollar spent on the last unit of each good is equal.

Law of demand states that, other things being equal, more of a good is demanded the lower is its price.

Ordinal utility assumes that individuals can rank commodity bundles in accordance with the level of satisfaction associated with each bundle.

Budget constraint defines all bundles of goods that the consumer can afford with a given budget.

Affordable set of goods and services for the consumer is bounded by the budget line from above; the **non-affordable set** lies strictly above the budget line.

Indifference curve defines combinations of goods and services that yield the same level of satisfaction to the consumer.

Indifference map is a set of indifference curves, where curves further from the origin denote a higher level of satisfaction.

Marginal rate of substitution is the slope of the indifference curve. It defines the amount of one good the consumer is willing to sacrifice in order to obtain a given increment of the other, while maintaining utility unchanged.

Diminishing marginal rate of substitution reflects a higher marginal value being associated with smaller quantities of any good consumed.

Consumer optimum occurs where the chosen consumption bundle is a point such that the price ratio equals the marginal rate of substitution.

This page titled [6.5: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

6.6: Exercises for Chapter 6

EXERCISE 6.1

In the example given in Table 6.1, suppose Neal experiences a small increase in income. Will he allocate it to snowboarding or jazz? [Hint: At the existing equilibrium, which activity will yield the higher MU for an additional dollar spent on it?]

EXERCISE 6.2

Suppose that utility depends on the square root of the amount of good X consumed: $U = \sqrt{X}$.

- In a spreadsheet enter the values 1... 16 as the X column (col A), and in the adjoining column (B) compute the value of utility corresponding to each quantity of X . To do this use the 'SQRT' command. For example, the entry in cell B3 will be of the form '=SQRT(A3)'.
- In the third column enter the marginal utility (MU) associated with each value of X – the change in utility in going from one value of X to the next.
- Use the 'graph' tool to map the relationship between U and X .
- Use the graph tool to map the relationship between MU and X .

EXERCISE 6.3

Instead of the square-root utility function in Exercise 6.2, suppose that utility takes the form $U=x^2$.

- Follow the same procedure as in the previous question – graph the utility function.
- Why is this utility function not consistent with our beliefs on utility?

EXERCISE 6.4

- Plot the utility function $U=2X$, following the same procedure as in the previous questions.
- Next plot the marginal utility values in a graph. What do we notice about the behaviour of the MU ?

EXERCISE 6.5

Let us see if we can draw a utility function for beer. In this instance the individual may reach a point where he takes too much.

- If the utility function is of the form $U=6X-X^2$, plot the utility values for X values in the range 1...8, using either a spreadsheet or manual calculations.
- At how many units of X (beer) is the individual's utility maximized?
- At how many beers does the utility become negative?

EXERCISE 6.6

Cappuccinos, C , cost \$3 each, and music downloads of your favourite artist, M , cost \$1 each from your iTunes store. Income is \$24.

- Draw the budget line, with cappuccinos on the vertical axis, and music on the horizontal axis, and compute the values of the intercepts.
- What is the slope of the budget constraint, and what is the opportunity cost of 1 cappuccino?
- Are the following combinations of goods in the affordable set: (4C and 9M), (6C and 2M), (3C and 15M)?
- Which combination(s) above lie inside the affordable set, and which lie on the boundary?

EXERCISE 6.7

George spends his income on gasoline and "other goods."

- First, draw a budget constraint, with gasoline on the horizontal axis.
- Suppose now that, in response to a gasoline shortage in the economy, the government imposes a *ration* on each individual that limits the purchase of gasoline to an amount less than the gasoline intercept of the budget constraint. Draw the new effective budget constraint.

EXERCISE 6.8

Suppose that you are told that the indifference curves defining the trade-off for two goods took the form of straight lines. Which of the four properties outlines in Section 6.3 would such indifference curves violate?

EXERCISE 6.9

Draw an indifference map with several indifference curves and several budget constraints corresponding to different possible levels of income. Note that these budget constraints should all be parallel because only income changes, not prices. Now find some optimizing (tangency) points. Join all of these points. You have just constructed what is called an income-consumption curve. Can you understand why it is called an income-consumption curve?

EXERCISE 6.10

Draw an indifference map again, in conjunction with a set of budget constraints. This time the budget constraints should each have a different price of good X and the same price for good Y .

- Draw in the resulting equilibria or tangencies and join up all of these points. You have just constructed a price-consumption curve for good X . Can you understand why the curve is so called?
- Now repeat part (a), but keep the price of X constant and permit the price of Y to vary. The resulting set of equilibrium points will form a price consumption curve for good Y .

EXERCISE 6.11

Suppose that movies are a normal good, but public transport is inferior. Draw an indifference map with a budget constraint and initial equilibrium. Now let income increase and draw a plausible new equilibrium, noting that one of the goods is inferior.

- Note that, with two snowboard outings and seven jazz club visits, total utility is 352 ($=132+220$), while the optimal combination of four of each yields a total utility of 380 ($=224+156$).

This page titled [6.6: Exercises for Chapter 6](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7: Firms, investors and capital markets

Chapter 7: Firms, investors and capital markets

In this chapter we will explore:

7.1	Business organization
7.2	Corporate goals – profit
7.3	Risk and the investor
7.4	Pooling risks

7.1 Business organization

Suppliers of goods and services to the marketplace come in a variety of forms; some are small, some are large. But, whatever their size, suppliers choose an organizational structure that is appropriate for their business: Aircraft, oil rigs, social media and information services are produced by large corporations; dental services and family health are provided by individual professionals or private partnerships.

The initial material of this chapter addresses organizational forms, their goals and their operation. We then examine why individuals choose to invest in firms, and illustrate that such investment provides individual investors with a means both to earning a return on their savings and to managing the risk associated with investing. Uncertainty regarding the future is a central consideration.

Understanding the way firms and capital markets function is crucial to understanding our economic history and how different forms of social and economic institutions interact. For example, seventeenth-century Amsterdam had a thriving bourgeoisie, well-developed financial markets, and investors with savings. This environment facilitated the channeling of investors' funds to firms specializing in trade and nautical conquest. This tiny state was then the source of some of the world's leading explorers and traders, and it had colonies stretching to Indonesia. The result was economic growth and prosperity.

In contrast, for much of the twentieth century, the Soviet Union dominated a huge territory covering much of Asia and Europe. But capital markets were non-existent, independent firms were stifled, and economic decline ultimately ensued. Much of the enormous difference in the respective patterns of economic development can be explained by the fact that one state fostered firms, capital markets, and legal institutions, while the other did not. In terms of our production possibility frontier: One set of institutional arrangements was conducive to expanding the possibilities; the other was not. Sustainable new businesses invariably require investors at an early point in the lifecycle of the business. Accordingly, financial and legal institutions that facilitate the flow of savings and financial investment into new enterprises perform a vital function in the economy.

Businesses, or firms, have several different forms. At the smallest scale, a business takes the form of a sole proprietor or sole trader who is the exclusive owner. A sole trader gets all of the revenues from the firm and incurs all of the costs. Hence he may make profits or be personally liable for the losses. In the latter case his business or even personal assets may be confiscated to cover debts. Personal bankruptcy may result.

Sole proprietor is the single owner of a business.

If a business is to grow, partners may be required. Such partners can inject money in exchange for a share of future profits. Firms where trust is involved, such as legal or accounting firms, typically adopt this structure. A firm is given credibility when customers see that partners invest their own wealth in it.

Partnership: a business owned jointly by two or more individuals, who share in the profits and are jointly responsible for losses.

In order to expand and grow, a firm will need cash, perhaps partners, and investors. Providers of family health and dental services rely primarily on human expertise, and therefore they need relatively little physical capital. Hence their cash start-up needs are limited. But firms that produce aircraft, or develop software and organizational systems, need vast amounts of money for capital investment; pharmaceuticals may need a billion dollars worth of research and development to bring a new drug to the marketplace; ride-sharing companies need billions in order to establish their business globally. Such businesses must form corporations – also known as companies. Not all corporations are public; some are privately held, but relatively few large corporations are not publicly traded.

Large organizations have several inherent advantages over small organizations when a high output level is required. Specialization in particular tasks leads to increased efficiency for production workers. At the same time, non-production workers can perform a multitude of different tasks. If a large corporation decided to contract out every task involved in bringing its product to market, the costs of such agreements would be prohibitively high. In addition, synergies can arise from teamwork. New ideas and better work flow are more likely to materialize when individuals work in close proximity than when working as isolated units, no matter how efficient they may be individually. A key aspect of such large organizations is that *they have a legal identity separate from the managers and owners*.

Corporation or company is an organization with a legal identity separate from its owners that produces and trades.

The owners of a corporation are known as its shareholders, and their object is usually to make profits. There also exist non-profit corporations whose objective may be philanthropic. Since our focus is upon markets, we will generally assume that profits form the objective of a typical corporation. The profits that accrue to a corporation may be paid to the shareholders in the form of a dividend, or retained in the corporation for future use. When large profits (or losses) accrue the value of the corporation increases (or decreases), and this is reflected in the value of each share of the company. If the value of each share in the company increases (decreases) there is a capital gain (loss) to the owners of the shares – the shareholders. In any given year shareholders may receive a dividend and also obtain a capital gain (or loss). The sum of the dividend and capital gain represents the return to owning corporate stock in that year. When this sum is adjusted for inflation it is termed the real return on corporate stock

Shareholders invest in corporations and therefore are the owners.

Dividends are payments made from after-tax profits to company shareholders.

Capital gains (losses) arise from the ownership of a corporation when an individual sells a share at a price higher (lower) than when the share was purchased.

Real return to corporate stock is the inflation-adjusted sum of dividends and capital gain (or loss).

A key difference between a company and a partnership is that a company involves limited liability, whereas a partnership does not. Limited liability means that the liability of the company is limited to the value of the company's assets. Shareholders cannot be further liable for any wrongdoing on the part of the company. Accordingly, partnerships and sole traders normally insure themselves and their operations. For example, all specialist doctors carry malpractice insurance, and engineers insure themselves against error.

Limited liability means that the liability of the company is limited to the value of the company's assets.

Corporations use capital, labour, and human expertise to produce a good, to supply a service, or to act as an intermediary. Corporations are required to produce an annual income statement that accurately describes the operation of the firm. An example is given in Table 7.1.

Table 7.1 The Regal Bank of Toronto, 2025

Total Revenue	\$	32.0b
Net income post tax	\$	4.80b
Shares outstanding		640m
Net income/share	\$	7.50
Dividends/share	\$	2.50
Share price	\$	72.0
Market capitalization	\$	46.08b

The data in Table 7.1 define the main financial characteristics of an imaginary bank: the Regal Bank of Toronto in the year 2025. "Net income post-tax" represents after-tax profits. There are 640 million shares outstanding, and thus each share could be attributed a profit of \$7.50 ($= \$4.80b \div 640m$). Of this amount, \$2.50 is distributed to shareholders in the form of dividends per share. The remainder is held by the Corporation in the form of retained earnings - to be used for future investment primarily. Each share traded at a price of \$72.00. Given that there were 640 million shares, the total market valuation of the corporation at that time stood at \$46.08 billion ($= 640m \times \72.00).

Such information is publicly available for a vast number of corporations at the 'finance' section of major search engines such as *Google* or *Yahoo*.

Retained earnings are the profits retained by a company for reinvestment and not distributed as dividends.

In Canada, the corporate sector as a whole tends to hold on to more than half of after-tax profits in the form of retained earnings. However there exists considerable variety in the behaviour of corporations, and most firms establish a pattern of how profits are allocated between dividends and retained earnings. In the Table 7.1 example, one third of profits are distributed; yet some corporations have a no-dividend policy. In these latter cases the benefit to investing in a firm must come in the form of capital gain to the owners of the shares.

7.2 Profit

Ownership and corporate goals

As economists, we believe that profit maximization accurately describes a typical firm's objective. However, since large firms are not run by their owners but by their executives or agents, it is frequently hard for the shareholders to know exactly what happens within a company. Even the board of directors—the guiding managerial group—may not be fully aware of the decisions, strategies, and practices of their executives and managers. Occasionally things go wrong, sometimes as a result of managers deciding to follow their own interests rather than the interests of the company. In technical terms, the interests of the corporation and its shareholders might not be aligned with the interests of its managers. For example, managers might have a short horizon and take steps to increase their own income in the short term, knowing that they will move to another job before the long-term effects of their decisions impact the firm.

At the same time, the marketplace for the ownership of corporations exerts a certain discipline: If firms are not as productive or profitable as possible, they may become *subject to takeover* by other firms. Fear of such takeover can induce executives and boards to maximize profits.

The shareholder-manager relationship is sometimes called a principal-agent relationship, and it can give rise to a principal-agent problem. If it is costly or difficult to monitor the behaviour of an agent because the agent has additional information about his own performance, the principal may not know if the agent is working to achieve the firm's goals. This is the principal-agent problem.

Principal or owner: delegates decisions to an agent, or manager.

Agent: usually a manager who works in a corporation and is directed to follow the corporation's interests.

Principal-agent problem: arises when the principal cannot easily monitor the actions of the agent, who therefore may not act in the best interests of the principal.

In an effort to deal with such a challenge, corporate executives frequently get bonuses or stock options that are related to the overall profitability of their firm. Stock options usually take the form of an executive being allowed to purchase the company's stock in the future – but at a price that is predetermined. If the company's profits do increase, then the price of the company's stock will reflect this and increase likewise. Hence the executive has an incentive to work with the objective of increasing profits because that will enable him to buy the company stock in the future at a lower price than it will be worth.

Stock option: an option to buy the stock of the company at a future date for a fixed, predetermined price.

The threat of takeover and the structure of rewards, together, imply that the assumption of profit maximization is a reasonable one.

Application Box 7.1 The 'Sub-Prime' mortgage crisis: A principal-agent problem

With a decline in interest and mortgage rates in the early part of the twenty first century, many individuals believed they could afford to buy a house because the borrowing costs were lower than before. Employees and managers of lending companies believed likewise, and they structured loans in such a way as to provide an incentive to low-income individuals to borrow. These mortgage loans frequently enabled purchasers to buy a house with only a 5% down payment, in some cases even less, coupled with a repayment schedule that saw low repayments initially but higher repayments subsequently. The initial interest cost was so low in many of these mortgages that it was even lower than the 'prime' rate – the rate banks charge to their most prized customers.

The crisis that resulted became known as the 'sub-prime' mortgage crisis. In many cases loan officers got bonuses based on the total value of loans they oversaw, regardless of the quality or risk associated with the loan. The consequence was that they had the incentive to make loans to customers to whom they would not have lent, had these employees and managers been lending their own money, or had they been remunerated differently. The outcomes were disastrous for numerous lending institutions. When interest

rates climbed, borrowers could not repay their loans. The construction industry produced a flood of houses that, combined with the sale of houses that buyers could no longer afford, sent housing prices through the floor. This in turn meant that recent house purchasers were left with negative value in their homes – the value of their property was less than what they paid for it. Many such 'owners' simply returned the keys to their bank, declared bankruptcy and walked away. Some lenders went bankrupt; some were bailed out by the government, others bought by surviving firms. This is a perfect example of the principal agent problem – the managers of the lending institutions and their loan officers did not have the incentive to act in the interest of the owners of those institutions.

The broader consequence of this lending practice was a financial collapse greater than any since the Depression of the nineteen thirties. Assets of the world's commercial and investment banks plummeted in value. Their assets included massive loans and investments both directly and indirectly to the real estate market, and when real estate values fell, so inevitably did the value of the assets based on this sector. Governments around the world had to buy up bad financial assets from financial institutions, or invest massive amounts of taxpayer money in these same institutions. Otherwise the world's financial system might have collapsed, with unknowable consequences.

Taxpayers and shareholders together bore the burden of this disastrous investment policy. Shareholders in many banks saw their shares drop in value to just a few percent of what they had been worth a year or two prior to the collapse.

Economic and accounting profit

Economists and accountants frequently differ in how they measure profits. An accountant stresses the financial flows of corporate activity; the economist is, in addition, concerned with opportunity cost. Imagine that Felicity has just inherited \$250,000 and decides to pursue her dream by opening a clothing boutique. She quits her job that pays her \$55,000 per annum, invests her inheritance in the purchase of a small retail space on the high street and launches her business. At the end of her first year she records \$110,000 in clothing sales, which she purchased from the wholesaler for \$50,000. She pays herself a salary of \$35,000 and has no other accounting costs because she owns her physical capital – the store. Her accounting profit for the year is given by the margin returned between the buying and selling price of her clothing (\$60,000) minus her incurred costs (\$35,000) in salary. Her accounting profit is thus \$25,000. Should she be content with this sum?

Felicity's economist friend, Prudence, informs Felicity that her enterprise is not returning a profit by economic standards. Prudence points out that Felicity could earn \$55,000 as an alternative to working in her own store, hence there is an additional implicit cost of \$20,000 to be considered, because Felicity only draws a salary of \$35,000. Furthermore, Felicity has invested \$250,000 in her business to avoid rent. But that sum, invested at the going interest rate of 4%, could earn her \$10,000 per annum. That too is a foregone income stream so it is an implicit cost. Altogether, the additional implicit costs, not included in the accounting flows, amount to \$30,000, and these implicit costs exceed the 'accounting profits'. Thus no economic profits are being made, because the economist includes implicit costs in her profit calculation. In economic terms Felicity would be better off by returning to her job and investing her inheritance. That strategy would generate an income of \$65,000, as opposed to the income of \$60,000 that she generates from the boutique – a salary of \$35,000 plus an accounting profit of \$25,000.

We can summarize this: Accounting profit is the difference between revenues and explicit costs. Economic profit is the difference between revenue and the sum of explicit and implicit costs. Explicit costs are the measured financial costs; Implicit costs represent the opportunity cost of the resources used in production.

Accounting profit: is the difference between revenues and explicit costs.

Economic profit: is the difference between revenue and the sum of explicit and implicit costs.

Explicit costs: are the measured financial costs.

Implicit costs: represent the opportunity cost of the resources used in production.

We will return to these concepts in the following chapters. Opportunity cost, or implicit costs, are critical in determining the long-run structure of certain sectors in the economy.

7.3 Risk and the investor

Firms cannot grow without investors. A successful firm's founder always arrives at a point where more investment is required if her enterprise is to expand. Frequently, she will not be able to secure a sufficiently large loan for such growth, and therefore must induce outsiders to buy shares in her firm. She may also realize that expansion carries risk, and she may want others to share in this

risk. Risk plays a central role in the life of the firm and the investor. Most investors prefer to avoid risk, but are prepared to assume a limited amount of it if the anticipated rewards are sufficiently attractive.

An illustration of risk-avoidance is to be seen in the purchase of home insurance. Most home owners who run even a small risk of seeing their house burn down, being flooded, or damaged by a gas leak purchase insurance. By doing so they are avoiding risk. But how much are they willing to pay for such insurance? If the house is worth \$500,000 and the probability of its being destroyed is one in one thousand in a given year then, using an averaging perspective, individuals should be willing to pay an insurance premium of \$500 per annum. That insurance premium represents what actuaries call a 'fair' gamble: If the probability of disaster is one in one thousand, then the 'fair' premium should be one thousandth the value of the home that is being insured. If the insurance company insures millions of homes, then on average it will have to pay for the replacement of one house for every one thousand houses it insures each year. So by charging homeowners a price that exceeds \$500 the insurer will cover not only the replacement cost of homes, but in addition cover her administrative costs and perhaps make a profit. Insurers operate on the basis of what we sometimes call the 'law of large numbers'.

In fact however, most individuals are willing to pay more than this 'fair' amount, and actually do pay more. If the insurance premium is \$750 or \$1,000 the home-owner is paying more than is actuarially 'fair', but a person who dislikes risk may be willing to pay such an amount in order to avoid the risk of being uninsured.

Our challenge now is to explain why individuals who purchase home insurance on terms that are less than actuarially 'fair' in order to avoid risk are simultaneously willing to invest their retirement savings into risky companies. Companies, like homes, are risky; while they may not collapse or implode in any given year, they can have good or bad returns in any given year. Corporate returns are inherently unpredictable and therefore risky. The key to understanding the willingness of risk-averse individuals to invest in risky firms is to be found in the pooling of risks.

7.4 Risk pooling and diversification

Silicon Valley: from angel investors to public corporation

Risky firms frequently succeed in attracting investment through the capital market in the modern economy. A typical start-up firm in the modern economy originates in the form of an idea. The developers of *Uber* got the idea of simplifying and streamlining the ride sharing business (sometimes called the taxi business). In its simplest form the inventors developed an App that would link users to drivers. The developers of *Airbnb* got the idea that spare accommodations in individual homes could be used to satisfy the needs of travellers. The founders developed an efficient means of putting potential renters/guests in communication with suppliers of rooms, houses and condominiums. *WeWork* was founded on the belief that workers and small corporations, particularly those in the technology sector, have immediate and changing space needs. The result is that *WeWork* provides flexible work space on a 'just in time' basis, frequently on a shared overhead basis. Each of these corporations acts as an intermediary, and sells intermediation services.

Typically the initial funding for new ideas comes from a couple of founding partners who develop a model or prototype of their software or their business. Following trials, the founders may approach potential investors for 'small' amounts that will fund expansion. Such investors are frequently 'angel' investors because they are a source of funding that makes the difference between expansion and death for the venture in question. If the venture shows promise the founders seek 'round A' funding, and this funding may come from venture capitalists who specialize in new ventures. Further evidence of possible success may result in 'round B' funding, and this frequently amounts to hundreds of millions of dollars.

Venture funds are managed by partners, or capitalists, with reputations for being better able than most to predict which start-ups will ultimately see profitability. These venture capitalists invest both their own funds and the funds of individuals who entrust their accumulated savings to the investing partnership.

But extremely high risk is associated with most start-ups. Funding frequently takes place in an environment where the venture has no revenue; merely a product in the course of development. Winners in the new economy are recognized and celebrated. Bill Gates' *Microsoft*, Jeff Bezos' *Amazon*, Steve Jobs' *Apple*, and Larry Page and Sergey Brin's *Google* are corporate giants with valuations approach one trillion dollars. But Elizabeth Holmes' *Theranos*, once with an implicit valuation of several billion dollars has expired. *Theranos* hired several hundred employees with the aim of developing blood tests for scores of purposes using just a pin-prick of blood. But it was a failure, despite attracting hundreds of millions of dollars in investment. In the year 2019 Canada had dozens of cannabis-based firms listed on the Canadian Securities Exchange. None of these is earning a profit in 2020, some have negligible earnings, yet several have a value in excess of one billion dollars. It is highly improbable that all will survive.

Capital market: a set of financial institutions that funnels financing from investors into bonds and stocks.

Venture capital: investment in a business venture, where the ultimate outcome is highly unpredictable.

How can we reconcile the fact that, while these firms carry extraordinary uncertainty, investors are still willing to part with large sums of money to fund development? And the investors are not only billionaires with a good sense of the marketplace; private individuals who save for their retirement also invest in risky firms on the advice of their financial manager. Let us explore how and why.

Dealing with risk

Most (sensible) investors hold a portfolio of investments, which is a combination of different stocks and bonds. By investing in different stocks and bonds rather than concentrating in one single investment or type of investment, an individual diversifies her portfolio, which is to say she engages in risk pooling. Venture capital partnerships act in the same way. They recognize that their investments in start-ups will yield both complete failures and some roaring successes; the variation in their outcomes will exceed the variation in the outcomes of an investor who invests in 'mature' corporations.

Portfolio: a combination of assets that is designed to secure an income from investing and to reduce risk.

Risk pooling: Combining individual risks in such a way that the aggregate risk is reduced.

A rigorous theory underlies this "don't put all of your eggs in the one basket" philosophy. The essentials of diversification or pooling are illustrated in the example given in Table 7.2 below.¹ There are two risky stocks here: Natural Gas (NG) and technology (Tech). Each stock is priced at \$100, and over time it is observed that each yields a \$10 return in good times and \$0 in bad times. The investor has \$200 to invest, and each sector independently has a 50% probability ($p=0.5$) of good or bad times. This means that each stock should yield a \$5 return *on average*: half of all outcomes will yield \$10 and half will yield zero. The challenge here is to develop an investment strategy that minimizes the risk for the investor.

At this point we need a specific working definition of risk. We define it in terms of how much variation a stock might experience in its returns from year to year. Each of NG and Tech have returns of either \$0 or \$10, with equal probability. But what if the Tech returns were either +\$20 or -\$10 with equal probability; or +\$40 or -\$30 with equal probability? In each of these alternative scenarios the average outcome remains the same: A positive average return of \$5. If the returns profile to NG remains unchanged we would say that Tech is a riskier stock (than NG) if its returns were defined by one of the alternatives here. Note that the average return is unchanged, and we are defining risk in terms of the greater spread in the possible returns around an unchanged average. The key to minimizing risk in the investor's portfolio lies in exploring how the variation in returns can be minimized by pooling risks.

Risk measurement: A higher degree of risk is associated with increased variation in the possible returns around an unchanged mean return.

Table 7.2 Investment strategies with risky assets

Strategy	Expected returns with probabilities		
\$200 in NG	220 ($p=0.5$)		200 ($p=0.5$)
\$200 in Tech	220 ($p=0.5$)		200 ($p=0.5$)
\$100 in each	220 ($p=0.25$)	210 ($p=0.5$)	200 ($p=0.25$)

The outcomes from three different investment strategies are illustrated in Table 7.2. By investing all of her \$200 in either NG or Tech, she will obtain \$220 half of the time and \$200 half of the time, as indicated in the first two outcome rows. But by diversifying through buying one of each stock, as illustrated in the final row, she reduces the variability of her portfolio. To see why note that, since the performance of each stock is independent, there is now only a one chance in four that both stocks do well, and therefore there is a 25 percent probability of earning \$220. By the same reasoning, there is a 25 percent probability of earning \$200. But there is a 50 percent chance that one stock will perform well and the other poorly. When that happens, she gets a return of \$210. In contrast to the outcomes defined in rows 1 and 2, the diversification strategy in row 3 *yields fewer extreme potential outcomes and more potential outcomes that lie closer to the mean outcome*.

Diversification reduces the total risk of a portfolio by pooling risks across several different assets whose individual returns behave independently.

Further diversification could reduce the variation in possible returns even further. To see this, imagine that, rather than having a choice between investing in one or two stocks, we could invest in four different stocks with the same returns profile as the two given in the table above. In such a case, the likelihood of getting extreme returns would be even lower than when investing in two stocks. This is because, if the returns to each stock are independent of the returns on the remaining stocks, it becomes increasingly improbable that all, or almost all, of the stocks will experience favorable (or unfavorable) returns in the same year. Now, imagine that we had 8 stocks, or 16, or 32, or 64, etc. The "magic" of diversification is that the same average return can be attained, yet variability can be reduced. If it can be reduced sufficiently by adding ever more stocks to the portfolio, then even a highly risk-averse individual can build a portfolio that is compatible with buying into risky firms.

We can conclude from this simple example that there need be no surprise over the fact that risk-averse individuals are willing, at the same time, to pay a high home-insurance premium to avoid risk, and simultaneously invest in risky ventures.

Application Box 7.2 The value of a financial advisor

The modern economy has thousands of highly-trained financial advisors. The successful ones earn huge salaries. But there is a puzzle: Why do such advisors exist? Can they predict the behaviour of the market any better than an uninformed advisor? Two insights help us answer the question.

First, Burton Malkiel wrote a best seller called *A Random Walk down Wall Street*. He provided ample evidence that a portfolio chosen on the basis of a monkey throwing darts at a list of stocks would do just as well as the average portfolio constructed by your friendly financial advisor.

Second, there are costs of transacting: An investor who builds a portfolio must devote time to the undertaking, and incur the associated financial trading cost. In recognizing this, investors may choose to invest in what they call mutual funds – a diversified collection of stocks – or may choose to employ a financial advisor who will essentially perform the same task of building a diversified portfolio. But, on average, financial advisors cannot beat the market, even though many individual investors would like to believe otherwise.

At this point we may reasonably ask why individuals choose to invest any of their funds in a "safe" asset – perhaps cash or Canadian Government bonds. After all, if their return to bonds is lower on average than the return to stocks, and they can diversify away much of the risk associated with stocks, why not get the higher average returns associated with stocks and put little or nothing in the safer asset? The reason is that it is impossible to fully diversify. When a recession hits, for example, the *whole stock market* may take a dive, because profits fall across the whole economy. On account of this possibility, we cannot ever arrive at a portfolio where the returns to the different stocks are completely independent. As a consequence, the rational investor will decide to put some funds in bonds in order to reduce this systematic risk component that is associated with the whole market. This is not a completely risk-free strategy because such assets can depreciate in value with inflation.

To see how the whole market can change dramatically students can go to any publicly accessible financial data site – such as *Yahoo Finance* and attempt to plot the TSX for the period 2005 – present, or the NASDAQ index from the mid-nineties to the present. The year 2020 is a particularly appropriate year to examine. In that year the coronavirus pandemic struck with disastrous impacts on stock markets worldwide. Initially almost all stocks declined in value. But within a matter of days investors realized that some firms would perform better in this particular downturn: those specializing in home delivery and those specializing in home exercise equipment for example. The stock valuations of firms such as Shopify, Amazon and Peloton shot up, while the valuations of traditional auto makers languished.

Efficiency and Allocation

We have now come full circle. We started this chapter by describing the key role in economic development and growth played by firms and capital markets. Capital markets channel the funds of individual investors to risk-taking firms. Such firms—whether they are Dutch spice importers in the seventeenth century, the Hudson's Bay Company in nineteenth-century Canada, communications corporations such as *Airbnb* or *Expedia*, or some high tech start-ups in Silicon Valley—are engines of growth and play a pivotal role in an economy's development. Capital markets are what make it possible for these firms to attract the savings of risk-averse individuals. By enabling individuals to diversify their portfolios, capital markets form the link between individuals and firms.

But capital markets fulfill another function, or at least they frequently do. They are a means of funnelling financial capital into ventures that appear to have a future return. It is not possible for each individual saver to perform the research necessary on a series of existing or new corporations, or 'ventures', in order to be able to invest in a knowledgeable manner. That is one reason we have financial intermediaries. When individuals deposit their savings with a bank, or with a financial manager, these individuals are anticipating that their savings will be protected and that a return will be forthcoming. A bank may promise a return of a fixed

percent if an individual deposits her money in a guaranteed investment certificate. Alternatively, if the individual saver wishes to take on some risk she can place her savings with her financial manager, an equity fund, or even a venture capitalist. These intermediaries are better at assessing risks and returns than most private individuals. This in turn means that the return to the individual from entrusting their savings to one of them should on average exceed the returns that the individual would earn herself by following some investment strategy.

If the professional investor indeed invests in more profitable ventures than an amateur investor, then that intermediary is performing an efficiency function for the whole economy: He does better at directing the economy's savings to where it is more productive on a macro level. This in turn means that the economy should have a higher growth rate than if savings are allocated towards ventures that are less likely to grow and satisfy a need or a demand in the economy.

Consider the example of *Airbnb* that we cited earlier. The original intent of this corporation was to provide the owners of unused (home) space the opportunity to earn a return on that space. *Airbnb* was thus a transformative mechanism, in that it enabled unused resources to be more fully utilized - by linking potential buyers who were willing to pay for the product, with potential sellers who were willing to supply at a price buyers were willing to pay. Unused resources became utilized, created a surplus and contributed to growth in the macro economy.

While financial intermediaries perform a valuable service to both individual savers, and the economy at large, we should not expect that intermediaries always make optimal decisions. However, these analysts have research resources available, and thus they have a comparative advantage over individuals for whom investing is a part-time activity. By being more efficient than individuals, financial intermediaries perform their broader economic allocation function, even if that is an unintended by-product of their professional activity.

At times professional investors suffer from what has been called 'irrational exuberance'. Crowd psychology creeps into the investment world from time to time, sometimes with devastating consequences. In the late 1990s tech stocks were all the rage and the stock market that specialized in trading such stocks saw the capital value of these stocks rise to stratospheric heights. The NASDAQ index stood at about 5,000 in March 2,000 but crashed to 1,300 by January of 2003. The run-up in NASDAQ valuations in the late nineties turned out to be a bubble.

Conclusion

We next turn to examine decision making *within* the firm. Firms must make the right decisions if they are to grow and provide investors with a satisfactory return. Firms that survive the growth process and ultimately bring a product to market are the survivors of the uncertainty surrounding product development.

Key Terms

Sole proprietor is the single owner of a business and is responsible for all profits and losses.

Partnership: a business owned jointly by two or more individuals, who share in the profits and are jointly responsible for losses.

Corporation or company is an organization with a legal identity separate from its owners that produces and trades.

Shareholders invest in corporations and therefore are the owners. They have limited liability personally if the firm incurs losses.

Dividends are payments made from after-tax profits to company shareholders.

Capital gains (losses) arise from the ownership of a corporation when an individual sells a share at a price higher (lower) than when the share was purchased.

Real return on corporate stock: the sum of dividend plus capital gain, adjusted for inflation.

Real return: the nominal return minus the rate of inflation.

Limited liability means that the liability of the company is limited to the value of the company's assets.

Retained earnings are the profits retained by a company for reinvestment and not distributed as dividends.

Principal or owner: delegates decisions to an agent, or manager.

Agent: usually a manager who works in a corporation and is directed to follow the corporation's interests.

Principal-agent problem: arises when the principal cannot easily monitor the actions of the agent, who therefore may not act in the best interests of the principal.

Stock option: an option to buy the stock of the company at a future date for a fixed, predetermined price.

Accounting profit: is the difference between revenues and explicit costs.

Economic profit: is the difference between revenue and the sum of explicit and implicit costs.

Explicit costs: are the measured financial costs.

Implicit costs: represent the opportunity cost of the resources used in production.

Capital market: a set of financial institutions that funnels financing from investors into bonds and stocks.

Portfolio: a combination of assets that is designed to secure an income from investing and to reduce risk.

Risk pooling: a means of reducing risk and increasing utility by aggregating or pooling multiple independent risks.

Risk: the risk associated with an investment can be measured by the dispersion in possible outcomes. A greater dispersion in outcomes implies more risk.

Diversification reduces the total risk of a portfolio by pooling risks across several different assets whose individual returns behave independently.

Exercises for Chapter 7

EXERCISE 7.1

Henry is contemplating opening a microbrewery and investing his savings of \$100,000 in it. He will quit his current job as a quality controller at Megawaiser where he is paid an annual salary of \$50,000. He plans on paying himself a salary of \$40,000 at the microbrewery. He also anticipates that his beer sales minus all costs other than his salary will yield him a surplus of \$55,000 per annum. The rate of return on savings is 7%.

1. Calculate the accounting profits envisaged by Henry.
2. Calculate the economic profits.
3. Should Henry open the microbrewery?
4. If all values except the return on savings remain the same, what rate of return would leave him indifferent between opening the brewery and not?

EXERCISE 7.2

You see an advertisement for life insurance for everyone 55 years of age and older. The advertisement says that no medical examination is required prior to purchasing insurance. If you are a very healthy 57-year old, do you think you will get a good deal from purchasing this insurance?

EXERCISE 7.3

In which of the following are risks being pooled, and in which would risks likely be spread by insurance companies?

1. Insurance against Alberta's Bow River Valley flooding.
2. Life insurance.
3. Insurance for the voice of Avril Lavigne or Celine Dion.
4. Insuring the voices of the lead vocalists in Metallica, Black Eyed Peas, Incubus, Evanescence, Green Day, and Jurassic Five.

EXERCISE 7.4

Your house has a one in five hundred probability chance of burning down in any given year. It is valued at \$350,000.

1. What insurance premium would be actuarially fair for this situation?
2. If the owner is willing to pay a premium of \$900, does she dislike risk or is she indifferent to risk?

EXERCISE 7.5

If individuals experience diminishing marginal utility from income it means that their utility function will resemble the total utility functions developed graphically in Section 6.2. Let us imagine specifically that if Y is income and U is utility, the individual gets utility from income according to the relation $U = \sqrt{Y}$.

1. In a spreadsheet or using a calculator, calculate the amount of utility the individual gets for all income values running from \$1 to \$25.
2. Graph the result with utility on the vertical axis and income on the horizontal axis, and verify from its shape that the marginal utility of income is declining.
3. Using your calculations, how much utility will the individual get from \$4, \$9 and \$16?
4. Suppose now that income results from a lottery and half of the time the individual gets \$4 and half of the time he gets \$16. How much utility will he get on average?
5. Now suppose he gets \$10 each time with certainty. How much utility will he get from this?
6. Since \$10 is exactly an average of \$4 and \$16, can you explain why \$10 with certainty gives him more utility than getting \$4 and \$16 each half of the time?

EXERCISE 7.6

In Question 7.5, suppose that the individual gets utility according to the relation $U = \frac{1}{2}Y$. Repeat the calculations for each part of the question and see if you can understand why the answers are different.

1. A different form of risk management is defined by the idea of risk spreading. Imagine that an oil supertanker has to be insured and the owner approaches one insurer. In the event of the tanker being ship-wrecked the damage caused by the resulting oil spill would be catastrophic – both to the environment and the insurance company. In this instance the insurance company may not benefit from the law of large numbers – it may not be insuring thousands of tankers and therefore would find it difficult to balance the potential claims with the annual insurance premiums. As a consequence, an insurer will spread the potential cost among other insurers – this is called risk spreading. The world's major insurers, such as *Lloyd's of London*, have hundreds of syndicates who each take on a small proportion of a big risk. These syndicates may again choose to subdivide their share among others, until the big risk becomes widely spread. In this way it is possible to insure against almost any event or possibility, no matter how large.

This page titled [7: Firms, investors and capital markets](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7.1: Business organization

Suppliers of goods and services to the marketplace come in a variety of forms; some are small, some are large. But, whatever their size, suppliers choose an organizational structure that is appropriate for their business: Aircraft, oil rigs, social media and information services are produced by large corporations; dental services and family health are provided by individual professionals or private partnerships.

The initial material of this chapter addresses organizational forms, their goals and their operation. We then examine why individuals choose to invest in firms, and illustrate that such investment provides individual investors with a means both to earning a return on their savings and to managing the risk associated with investing. Uncertainty regarding the future is a central consideration.

Understanding the way firms and capital markets function is crucial to understanding our economic history and how different forms of social and economic institutions interact. For example, seventeenth-century Amsterdam had a thriving bourgeoisie, well-developed financial markets, and investors with savings. This environment facilitated the channeling of investors' funds to firms specializing in trade and nautical conquest. This tiny state was then the source of some of the world's leading explorers and traders, and it had colonies stretching to Indonesia. The result was economic growth and prosperity.

In contrast, for much of the twentieth century, the Soviet Union dominated a huge territory covering much of Asia and Europe. But capital markets were non-existent, independent firms were stifled, and economic decline ultimately ensued. Much of the enormous difference in the respective patterns of economic development can be explained by the fact that one state fostered firms, capital markets, and legal institutions, while the other did not. In terms of our production possibility frontier: One set of institutional arrangements was conducive to expanding the possibilities; the other was not. Sustainable new businesses invariably require investors at an early point in the lifecycle of the business. Accordingly, financial and legal institutions that facilitate the flow of savings and financial investment into new enterprises perform a vital function in the economy.

Businesses, or firms, have several different forms. At the smallest scale, a business takes the form of a sole proprietor or sole trader who is the exclusive owner. A sole trader gets all of the revenues from the firm and incurs all of the costs. Hence he may make profits or be personally liable for the losses. In the latter case his business or even personal assets may be confiscated to cover debts. Personal bankruptcy may result.

Sole proprietor is the single owner of a business.

If a business is to grow, partners may be required. Such partners can inject money in exchange for a share of future profits. Firms where trust is involved, such as legal or accounting firms, typically adopt this structure. A firm is given credibility when customers see that partners invest their own wealth in it.

Partnership: a business owned jointly by two or more individuals, who share in the profits and are jointly responsible for losses.

In order to expand and grow, a firm will need cash, perhaps partners, and investors. Providers of family health and dental services rely primarily on human expertise, and therefore they need relatively little physical capital. Hence their cash start-up needs are limited. But firms that produce aircraft, or develop software and organizational systems, need vast amounts of money for capital investment; pharmaceuticals may need a billion dollars worth of research and development to bring a new drug to the marketplace; ride-sharing companies need billions in order to establish their business globally. Such businesses must form corporations – also known as companies. Not all corporations are public; some are privately held, but relatively few large corporations are not publicly traded.

Large organizations have several inherent advantages over small organizations when a high output level is required. Specialization in particular tasks leads to increased efficiency for production workers. At the same time, non-production workers can perform a multitude of different tasks. If a large corporation decided to contract out every task involved in bringing its product to market, the costs of such agreements would be prohibitively high. In addition, synergies can arise from teamwork. New ideas and better work flow are more likely to materialize when individuals work in close proximity than when working as isolated units, no matter how efficient they may be individually. A key aspect of such large organizations is that *they have a legal identity separate from the managers and owners*.

Corporation or company is an organization with a legal identity separate from its owners that produces and trades.

The owners of a corporation are known as its shareholders, and their object is usually to make profits. There also exist non-profit corporations whose objective may be philanthropic. Since our focus is upon markets, we will generally assume that profits form the

objective of a typical corporation. The profits that accrue to a corporation may be paid to the shareholders in the form of a dividend, or retained in the corporation for future use. When large profits (or losses) accrue the value of the corporation increases (or decreases), and this is reflected in the value of each share of the company. If the value of each share in the company increases (decreases) there is a capital gain (loss) to the owners of the shares – the shareholders. In any given year shareholders may receive a dividend and also obtain a capital gain (or loss). The sum of the dividend and capital gain represents the return to owning corporate stock in that year. When this sum is adjusted for inflation it is termed the real return on corporate stock

Shareholders invest in corporations and therefore are the owners.

Dividends are payments made from after-tax profits to company shareholders.

Capital gains (losses) arise from the ownership of a corporation when an individual sells a share at a price higher (lower) than when the share was purchased.

Real return to corporate stock is the inflation-adjusted sum of dividends and capital gain (or loss).

A key difference between a company and a partnership is that a company involves limited liability, whereas a partnership does not. Limited liability means that the liability of the company is limited to the value of the company's assets. Shareholders cannot be further liable for any wrongdoing on the part of the company. Accordingly, partnerships and sole traders normally insure themselves and their operations. For example, all specialist doctors carry malpractice insurance, and engineers insure themselves against error.

Limited liability means that the liability of the company is limited to the value of the company's assets.

Corporations use capital, labour, and human expertise to produce a good, to supply a service, or to act as an intermediary. Corporations are required to produce an annual income statement that accurately describes the operation of the firm. An example is given in Table 7.1.

Table 7.1 The Regal Bank of Toronto, 2025

Total Revenue	\$	32.0b
Net income post tax	\$	4.80b
Shares outstanding		640m
Net income/share	\$	7.50
Dividends/share	\$	2.50
Share price	\$	72.0
Market capitalization	\$	46.08b

The data in Table 7.1 define the main financial characteristics of an imaginary bank: the Regal Bank of Toronto in the year 2025. "Net income post-tax" represents after-tax profits. There are 640 million shares outstanding, and thus each share could be attributed a profit of \$7.50 ($= \$4.80b \div 640m$). Of this amount, \$2.50 is distributed to shareholders in the form of dividends per share. The remainder is held by the Corporation in the form of retained earnings - to be used for future investment primarily. Each share traded at a price of \$72.00. Given that there were 640 million shares, the total market valuation of the corporation at that time stood at \$46.08 billion ($= 640m \times \72.00).

Such information is publicly available for a vast number of corporations at the 'finance' section of major search engines such as Google or Yahoo.

Retained earnings are the profits retained by a company for reinvestment and not distributed as dividends.

In Canada, the corporate sector as a whole tends to hold on to more than half of after-tax profits in the form of retained earnings. However there exists considerable variety in the behaviour of corporations, and most firms establish a pattern of how profits are allocated between dividends and retained earnings. In the Table 7.1 example, one third of profits are distributed; yet some corporations have a no-dividend policy. In these latter cases the benefit to investing in a firm must come in the form of capital gain to the owners of the shares.

This page titled [7.1: Business organization](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7.2: Profit

Ownership and corporate goals

As economists, we believe that profit maximization accurately describes a typical firm's objective. However, since large firms are not run by their owners but by their executives or agents, it is frequently hard for the shareholders to know exactly what happens within a company. Even the board of directors—the guiding managerial group—may not be fully aware of the decisions, strategies, and practices of their executives and managers. Occasionally things go wrong, sometimes as a result of managers deciding to follow their own interests rather than the interests of the company. In technical terms, the interests of the corporation and its shareholders might not be aligned with the interests of its managers. For example, managers might have a short horizon and take steps to increase their own income in the short term, knowing that they will move to another job before the long-term effects of their decisions impact the firm.

At the same time, the marketplace for the ownership of corporations exerts a certain discipline: If firms are not as productive or profitable as possible, they may become *subject to takeover* by other firms. Fear of such takeover can induce executives and boards to maximize profits.

The shareholder-manager relationship is sometimes called a principal-agent relationship, and it can give rise to a principal-agent problem. If it is costly or difficult to monitor the behaviour of an agent because the agent has additional information about his own performance, the principal may not know if the agent is working to achieve the firm's goals. This is the principal-agent problem.

Principal or owner: delegates decisions to an agent, or manager.

Agent: usually a manager who works in a corporation and is directed to follow the corporation's interests.

Principal-agent problem: arises when the principal cannot easily monitor the actions of the agent, who therefore may not act in the best interests of the principal.

In an effort to deal with such a challenge, corporate executives frequently get bonuses or stock options that are related to the overall profitability of their firm. Stock options usually take the form of an executive being allowed to purchase the company's stock in the future – but at a price that is predetermined. If the company's profits do increase, then the price of the company's stock will reflect this and increase likewise. Hence the executive has an incentive to work with the objective of increasing profits because that will enable him to buy the company stock in the future at a lower price than it will be worth.

Stock option: an option to buy the stock of the company at a future date for a fixed, predetermined price.

The threat of takeover and the structure of rewards, together, imply that the assumption of profit maximization is a reasonable one.

Application Box 7.1 The 'Sub-Prime' mortgage crisis: A principal-agent problem

With a decline in interest and mortgage rates in the early part of the twenty first century, many individuals believed they could afford to buy a house because the borrowing costs were lower than before. Employees and managers of lending companies believed likewise, and they structured loans in such a way as to provide an incentive to low-income individuals to borrow. These mortgage loans frequently enabled purchasers to buy a house with only a 5% down payment, in some cases even less, coupled with a repayment schedule that saw low repayments initially but higher repayments subsequently. The initial interest cost was so low in many of these mortgages that it was even lower than the 'prime' rate – the rate banks charge to their most prized customers.

The crisis that resulted became known as the 'sub-prime' mortgage crisis. In many cases loan officers got bonuses based on the total value of loans they oversaw, regardless of the quality or risk associated with the loan. The consequence was that they had the incentive to make loans to customers to whom they would not have lent, had these employees and managers been lending their own money, or had they been remunerated differently. The outcomes were disastrous for numerous lending institutions. When interest rates climbed, borrowers could not repay their loans. The construction industry produced a flood of houses that, combined with the sale of houses that buyers could no longer afford, sent housing prices through the floor. This in turn meant that recent house purchasers were left with negative value in their homes – the value of their property was less than what they paid for it. Many such 'owners' simply returned the keys to their bank, declared bankruptcy and walked away. Some lenders went bankrupt; some were bailed out by the government, others bought by surviving firms. This is a perfect example of the principal agent problem – the managers of the lending institutions and their loan officers did not have the incentive to act in the interest of the owners of those institutions.

The broader consequence of this lending practice was a financial collapse greater than any since the Depression of the nineteen thirties. Assets of the world's commercial and investment banks plummeted in value. Their assets included massive loans and investments both directly and indirectly to the real estate market, and when real estate values fell, so inevitably did the value of the assets based on this sector. Governments around the world had to buy up bad financial assets from financial institutions, or invest massive amounts of taxpayer money in these same institutions. Otherwise the world's financial system might have collapsed, with unknowable consequences.

Taxpayers and shareholders together bore the burden of this disastrous investment policy. Shareholders in many banks saw their shares drop in value to just a few percent of what they had been worth a year or two prior to the collapse.

Economic and accounting profit

Economists and accountants frequently differ in how they measure profits. An accountant stresses the financial flows of corporate activity; the economist is, in addition, concerned with opportunity cost. Imagine that Felicity has just inherited \$250,000 and decides to pursue her dream by opening a clothing boutique. She quits her job that pays her \$55,000 per annum, invests her inheritance in the purchase of a small retail space on the high street and launches her business. At the end of her first year she records \$110,000 in clothing sales, which she purchased from the wholesaler for \$50,000. She pays herself a salary of \$35,000 and has no other accounting costs because she owns her physical capital – the store. Her accounting profit for the year is given by the margin returned between the buying and selling price of her clothing (\$60,000) minus her incurred costs (\$35,000) in salary. Her accounting profit is thus \$25,000. Should she be content with this sum?

Felicity's economist friend, Prudence, informs Felicity that her enterprise is not returning a profit by economic standards. Prudence points out that Felicity could earn \$55,000 as an alternative to working in her own store, hence there is an additional implicit cost of \$20,000 to be considered, because Felicity only draws a salary of \$35,000. Furthermore, Felicity has invested \$250,000 in her business to avoid rent. But that sum, invested at the going interest rate of 4%, could earn her \$10,000 per annum. That too is a foregone income stream so it is an implicit cost. Altogether, the additional implicit costs, not included in the accounting flows, amount to \$30,000, and these implicit costs exceed the 'accounting profits'. Thus no economic profits are being made, because the economist includes implicit costs in her profit calculation. In economic terms Felicity would be better off by returning to her job and investing her inheritance. That strategy would generate an income of \$65,000, as opposed to the income of \$60,000 that she generates from the boutique – a salary of \$35,000 plus an accounting profit of \$25,000.

We can summarize this: Accounting profit is the difference between revenues and explicit costs. Economic profit is the difference between revenue and the sum of explicit and implicit costs. Explicit costs are the measured financial costs; Implicit costs represent the opportunity cost of the resources used in production.

Accounting profit: is the difference between revenues and explicit costs.

Economic profit: is the difference between revenue and the sum of explicit and implicit costs.

Explicit costs: are the measured financial costs.

Implicit costs: represent the opportunity cost of the resources used in production.

We will return to these concepts in the following chapters. Opportunity cost, or implicit costs, are critical in determining the long-run structure of certain sectors in the economy.

This page titled [7.2: Profit](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7.3: Risk and the investor

Firms cannot grow without investors. A successful firm's founder always arrives at a point where more investment is required if her enterprise is to expand. Frequently, she will not be able to secure a sufficiently large loan for such growth, and therefore must induce outsiders to buy shares in her firm. She may also realize that expansion carries risk, and she may want others to share in this risk. Risk plays a central role in the life of the firm and the investor. Most investors prefer to avoid risk, but are prepared to assume a limited amount of it if the anticipated rewards are sufficiently attractive.

An illustration of risk-avoidance is to be seen in the purchase of home insurance. Most home owners who run even a small risk of seeing their house burn down, being flooded, or damaged by a gas leak purchase insurance. By doing so they are avoiding risk. But how much are they willing to pay for such insurance? If the house is worth \$500,000 and the probability of its being destroyed is one in one thousand in a given year then, using an averaging perspective, individuals should be willing to pay an insurance premium of \$500 per annum. That insurance premium represents what actuaries call a 'fair' gamble: If the probability of disaster is one in one thousand, then the 'fair' premium should be one thousandth the value of the home that is being insured. If the insurance company insures millions of homes, then on average it will have to pay for the replacement of one house for every one thousand houses it insures each year. So by charging homeowners a price that exceeds \$500 the insurer will cover not only the replacement cost of homes, but in addition cover her administrative costs and perhaps make a profit. Insurers operate on the basis of what we sometimes call the 'law of large numbers'.

In fact however, most individuals are willing to pay more than this 'fair' amount, and actually do pay more. If the insurance premium is \$750 or \$1,000 the home-owner is paying more than is actuarially 'fair', but a person who dislikes risk may be willing to pay such an amount in order to avoid the risk of being uninsured.

Our challenge now is to explain why individuals who purchase home insurance on terms that are less than actuarially 'fair' in order to avoid risk are simultaneously willing to invest their retirement savings into risky companies. Companies, like homes, are risky; while they may not collapse or implode in any given year, they can have good or bad returns in any given year. Corporate returns are inherently unpredictable and therefore risky. The key to understanding the willingness of risk-averse individuals to invest in risky firms is to be found in the pooling of risks.

This page titled [7.3: Risk and the investor](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7.4: Risk pooling and diversification

Silicon Valley: from angel investors to public corporation

Risky firms frequently succeed in attracting investment through the capital market in the modern economy. A typical start-up firm in the modern economy originates in the form of an idea. The developers of *Uber* got the idea of simplifying and streamlining the ride sharing business (sometimes called the taxi business). In its simplest form the inventors developed an App that would link users to drivers. The developers of *Airbnb* got the idea that spare accommodations in individual homes could be used to satisfy the needs of travellers. The founders developed an efficient means of putting potential renters/guests in communication with suppliers of rooms, houses and condominiums. *WeWork* was founded on the belief that workers and small corporations, particularly those in the technology sector, have immediate and changing space needs. The result is that *WeWork* provides flexible work space on a 'just in time' basis, frequently on a shared overhead basis. Each of these corporations acts as an intermediary, and sells intermediation services.

Typically the initial funding for new ideas comes from a couple of founding partners who develop a model or prototype of their software or their business. Following trials, the founders may approach potential investors for 'small' amounts that will fund expansion. Such investors are frequently 'angel' investors because they are a source of funding that makes the difference between expansion and death for the venture in question. If the venture shows promise the founders seek 'round A' funding, and this funding may come from venture capitalists who specialize in new ventures. Further evidence of possible success may result in 'round B' funding, and this frequently amounts to hundreds of millions of dollars.

Venture funds are managed by partners, or capitalists, with reputations for being better able than most to predict which start-ups will ultimately see profitability. These venture capitalists invest both their own funds and the funds of individuals who entrust their accumulated savings to the investing partnership.

But extremely high risk is associated with most start-ups. Funding frequently takes place in an environment where the venture has no revenue; merely a product in the course of development. Winners in the new economy are recognized and celebrated. Bill Gates' *Microsoft*, Jeff Bezos' *Amazon*, Steve Jobs' *Apple*, and Larry Page and Sergey Brin's *Google* are corporate giants with valuations approach one trillion dollars. But Elizabeth Holmes' *Theranos*, once with an implicit valuation of several billion dollars has expired. *Theranos* hired several hundred employees with the aim of developing blood tests for scores of purposes using just a pin-prick of blood. But it was a failure, despite attracting hundreds of millions of dollars in investment. In the year 2019 Canada had dozens of cannabis-based firms listed on the Canadian Securities Exchange. None of these is earning a profit in 2020, some have negligible earnings, yet several have a value in excess of one billion dollars. It is highly improbable that all will survive.

Capital market: a set of financial institutions that funnels financing from investors into bonds and stocks.

Venture capital: investment in a business venture, where the ultimate outcome is highly unpredictable.

How can we reconcile the fact that, while these firms carry extraordinary uncertainty, investors are still willing to part with large sums of money to fund development? And the investors are not only billionaires with a good sense of the marketplace; private individuals who save for their retirement also invest in risky firms on the advice of their financial manager. Let us explore how and why.

Dealing with risk

Most (sensible) investors hold a portfolio of investments, which is a combination of different stocks and bonds. By investing in different stocks and bonds rather than concentrating in one single investment or type of investment, an individual diversifies her portfolio, which is to say she engages in risk pooling. Venture capital partnerships act in the same way. They recognize that their investments in start-ups will yield both complete failures and some roaring successes; the variation in their outcomes will exceed the variation in the outcomes of an investor who invests in 'mature' corporations.

Portfolio: a combination of assets that is designed to secure an income from investing and to reduce risk.

Risk pooling: Combining individual risks in such a way that the aggregate risk is reduced.

A rigorous theory underlies this "don't put all of your eggs in the one basket" philosophy. The essentials of diversification or pooling are illustrated in the example given in Table 7.2 below.¹ There are two risky stocks here: Natural Gas (NG) and technology (Tech). Each stock is priced at \$100, and over time it is observed that each yields a \$10 return in good times and \$0 in bad times. The investor has \$200 to invest, and each sector independently has a 50% probability ($p=0.5$) of good or bad times. This means that

each stock should yield a \$5 return *on average*: half of all outcomes will yield \$10 and half will yield zero. The challenge here is to develop an investment strategy that minimizes the risk for the investor.

At this point we need a specific working definition of risk. We define it in terms of how much variation a stock might experience in its returns from year to year. Each of NG and Tech have returns of either \$0 or \$10, with equal probability. But what if the Tech returns were either +\$20 or -\$10 with equal probability; or +\$40 or -\$30 with equal probability? In each of these alternative scenarios the average outcome remains the same: A positive average return of \$5. If the returns profile to NG remains unchanged we would say that Tech is a riskier stock (than NG) if its returns were defined by one of the alternatives here. Note that the average return is unchanged, and we are defining risk in terms of the greater spread in the possible returns around an unchanged average. The key to minimizing risk in the investor's portfolio lies in exploring how the variation in returns can be minimized by pooling risks.

Risk measurement: A higher degree of risk is associated with increased variation in the possible returns around an unchanged mean return.

Table 7.2 Investment strategies with risky assets

Strategy	Expected returns with probabilities		
\$200 in NG	220 ($p=0.5$)		200 ($p=0.5$)
\$200 in Tech	220 ($p=0.5$)		200 ($p=0.5$)
\$100 in each	220 ($p=0.25$)	210 ($p=0.5$)	200 ($p=0.25$)

The outcomes from three different investment strategies are illustrated in Table 7.2. By investing all of her \$200 in either NG or Tech, she will obtain \$220 half of the time and \$200 half of the time, as indicated in the first two outcome rows. But by diversifying through buying one of each stock, as illustrated in the final row, she reduces the variability of her portfolio. To see why note that, since the performance of each stock is independent, there is now only a one chance in four that both stocks do well, and therefore there is a 25 percent probability of earning \$220. By the same reasoning, there is a 25 percent probability of earning \$200. But there is a 50 percent chance that one stock will perform well and the other poorly. When that happens, she gets a return of \$210. In contrast to the outcomes defined in rows 1 and 2, the diversification strategy in row 3 *yields fewer extreme potential outcomes and more potential outcomes that lie closer to the mean outcome*.

Diversification reduces the total risk of a portfolio by pooling risks across several different assets whose individual returns behave independently.

Further diversification could reduce the variation in possible returns even further. To see this, imagine that, rather than having a choice between investing in one or two stocks, we could invest in four different stocks with the same returns profile as the two given in the table above. In such a case, the likelihood of getting extreme returns would be even lower than when investing in two stocks. This is because, if the returns to each stock are independent of the returns on the remaining stocks, it becomes increasingly improbable that all, or almost all, of the stocks will experience favorable (or unfavorable) returns in the same year. Now, imagine that we had 8 stocks, or 16, or 32, or 64, etc. The "magic" of diversification is that the same average return can be attained, yet variability can be reduced. If it can be reduced sufficiently by adding ever more stocks to the portfolio, then even a highly risk-averse individual can build a portfolio that is compatible with buying into risky firms.

We can conclude from this simple example that there need be no surprise over the fact that risk-averse individuals are willing, at the same time, to pay a high home-insurance premium to avoid risk, and simultaneously invest in risky ventures.

Application Box 7.2 The value of a financial advisor

The modern economy has thousands of highly-trained financial advisors. The successful ones earn huge salaries. But there is a puzzle: Why do such advisors exist? Can they predict the behaviour of the market any better than an uninformed advisor? Two insights help us answer the question.

First, Burton Malkiel wrote a best seller called *A Random Walk down Wall Street*. He provided ample evidence that a portfolio chosen on the basis of a monkey throwing darts at a list of stocks would do just as well as the average portfolio constructed by your friendly financial advisor.

Second, there are costs of transacting: An investor who builds a portfolio must devote time to the undertaking, and incur the associated financial trading cost. In recognizing this, investors may choose to invest in what they call mutual funds – a diversified

collection of stocks – or may choose to employ a financial advisor who will essentially perform the same task of building a diversified portfolio. But, on average, financial advisors cannot beat the market, even though many individual investors would like to believe otherwise.

At this point we may reasonably ask why individuals choose to invest any of their funds in a "safe" asset – perhaps cash or Canadian Government bonds. After all, if their return to bonds is lower on average than the return to stocks, and they can diversify away much of the risk associated with stocks, why not get the higher average returns associated with stocks and put little or nothing in the safer asset? The reason is that it is impossible to fully diversify. When a recession hits, for example, the *whole stock market* may take a dive, because profits fall across the whole economy. On account of this possibility, we cannot ever arrive at a portfolio where the returns to the different stocks are completely independent. As a consequence, the rational investor will decide to put some funds in bonds in order to reduce this systematic risk component that is associated with the whole market. This is not a completely risk-free strategy because such assets can depreciate in value with inflation.

To see how the whole market can change dramatically students can go to any publicly accessible financial data site – such as *Yahoo Finance* and attempt to plot the TSX for the period 2005 – present, or the NASDAQ index from the mid-nineties to the present. The year 2020 is a particularly appropriate year to examine. In that year the coronavirus pandemic struck with disastrous impacts on stock markets worldwide. Initially almost all stocks declined in value. But within a matter of days investors realized that some firms would perform better in this particular downturn: those specializing in home delivery and those specializing in home exercise equipment for example. The stock valuations of firms such as Shopify, Amazon and Peloton shot up, while the valuations of traditional auto makers languished.

Efficiency and Allocation

We have now come full circle. We started this chapter by describing the key role in economic development and growth played by firms and capital markets. Capital markets channel the funds of individual investors to risk-taking firms. Such firms—whether they are Dutch spice importers in the seventeenth century, the Hudson's Bay Company in nineteenth-century Canada, communications corporations such as *Airbnb* or *Expedia*, or some high tech start-ups in Silicon Valley—are engines of growth and play a pivotal role in an economy's development. Capital markets are what make it possible for these firms to attract the savings of risk-averse individuals. By enabling individuals to diversify their portfolios, capital markets form the link between individuals and firms.

But capital markets fulfill another function, or at least they frequently do. They are a means of funnelling financial capital into ventures that appear to have a future return. It is not possible for each individual saver to perform the research necessary on a series of existing or new corporations, or 'ventures', in order to be able to invest in a knowledgeable manner. That is one reason we have financial intermediaries. When individuals deposit their savings with a bank, or with a financial manager, these individuals are anticipating that their savings will be protected and that a return will be forthcoming. A bank may promise a return of a fixed percent if an individual deposits her money in a guaranteed investment certificate. Alternatively, if the individual saver wishes to take on some risk she can place her savings with her financial manager, an equity fund, or even a venture capitalist. These intermediaries are better at assessing risks and returns than most private individuals. This in turn means that the return to the individual from entrusting their savings to one of them should on average exceed the returns that the individual would earn herself by following some investment strategy.

If the professional investor indeed invests in more profitable ventures than an amateur investor, then that intermediary is performing an efficiency function for the whole economy: He does better at directing the economy's savings to where it is more productive on a macro level. This in turn means that the economy should have a higher growth rate than if savings are allocated towards ventures that are less likely to grow and satisfy a need or a demand in the economy.

Consider the example of *Airbnb* that we cited earlier. The original intent of this corporation was to provide the owners of unused (home) space the opportunity to earn a return on that space. *Airbnb* was thus a transformative mechanism, in that it enabled unused resources to be more fully utilized - by linking potential buyers who were willing to pay for the product, with potential sellers who were willing to supply at a price buyers were willing to pay. Unused resources became utilized, created a surplus and contributed to growth in the macro economy.

While financial intermediaries perform a valuable service to both individual savers, and the economy at large, we should not expect that intermediaries always make optimal decisions. However, these analysts have research resources available, and thus they have a comparative advantage over individuals for whom investing is a part-time activity. By being more efficient than individuals, financial intermediaries perform their broader economic allocation function, even if that is an unintended by-product of their professional activity.

At times professional investors suffer from what has been called 'irrational exuberance'. Crowd psychology creeps into the investment world from time to time, sometimes with devastating consequences. In the late 1990s tech stocks were all the rage and the stock market that specialized in trading such stocks saw the capital value of these stocks rise to stratospheric heights. The NASDAQ index stood at about 5,000 in March 2,000 but crashed to 1,300 by January of 2003. The run-up in NASDAQ valuations in the late nineties turned out to be a bubble.

This page titled [7.4: Risk pooling and diversification](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7.5: Conclusion

We next turn to examine decision making *within* the firm. Firms must make the right decisions if they are to grow and provide investors with a satisfactory return. Firms that survive the growth process and ultimately bring a product to market are the survivors of the uncertainty surrounding product development.

This page titled [7.5: Conclusion](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7.6: Key Terms

Sole proprietor is the single owner of a business and is responsible for all profits and losses.

Partnership: a business owned jointly by two or more individuals, who share in the profits and are jointly responsible for losses.

Corporation or company is an organization with a legal identity separate from its owners that produces and trades.

Shareholders invest in corporations and therefore are the owners. They have limited liability personally if the firm incurs losses.

Dividends are payments made from after-tax profits to company shareholders.

Capital gains (losses) arise from the ownership of a corporation when an individual sells a share at a price higher (lower) than when the share was purchased.

Real return on corporate stock: the sum of dividend plus capital gain, adjusted for inflation.

Real return: the nominal return minus the rate of inflation.

Limited liability means that the liability of the company is limited to the value of the company's assets.

Retained earnings are the profits retained by a company for reinvestment and not distributed as dividends.

Principal or owner: delegates decisions to an agent, or manager.

Agent: usually a manager who works in a corporation and is directed to follow the corporation's interests.

Principal-agent problem: arises when the principal cannot easily monitor the actions of the agent, who therefore may not act in the best interests of the principal.

Stock option: an option to buy the stock of the company at a future date for a fixed, predetermined price.

Accounting profit: is the difference between revenues and explicit costs.

Economic profit: is the difference between revenue and the sum of explicit and implicit costs.

Explicit costs: are the measured financial costs.

Implicit costs: represent the opportunity cost of the resources used in production.

Capital market: a set of financial institutions that funnels financing from investors into bonds and stocks.

Portfolio: a combination of assets that is designed to secure an income from investing and to reduce risk.

Risk pooling: a means of reducing risk and increasing utility by aggregating or pooling multiple independent risks.

Risk: the risk associated with an investment can be measured by the dispersion in possible outcomes. A greater dispersion in outcomes implies more risk.

Diversification reduces the total risk of a portfolio by pooling risks across several different assets whose individual returns behave independently.

This page titled [7.6: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

7.7: Exercises for Chapter 7

EXERCISE 7.1

Henry is contemplating opening a microbrewery and investing his savings of \$100,000 in it. He will quit his current job as a quality controller at Megaweiser where he is paid an annual salary of \$50,000. He plans on paying himself a salary of \$40,000 at the microbrewery. He also anticipates that his beer sales minus all costs other than his salary will yield him a surplus of \$55,000 per annum. The rate of return on savings is 7%.

- Calculate the accounting profits envisaged by Henry.
- Calculate the economic profits.
- Should Henry open the microbrewery?
- If all values except the return on savings remain the same, what rate of return would leave him indifferent between opening the brewery and not?

EXERCISE 7.2

You see an advertisement for life insurance for everyone 55 years of age and older. The advertisement says that no medical examination is required prior to purchasing insurance. If you are a very healthy 57-year old, do you think you will get a good deal from purchasing this insurance?

EXERCISE 7.3

In which of the following are risks being pooled, and in which would risks likely be spread by insurance companies?

- Insurance against Alberta's Bow River Valley flooding.
- Life insurance.
- Insurance for the voice of Avril Lavigne or Celine Dion.
- Insuring the voices of the lead vocalists in Metallica, Black Eyed Peas, Incubus, Evanescence, Green Day, and Jurassic Five.

EXERCISE 7.4

Your house has a one in five hundred probability chance of burning down in any given year. It is valued at \$350,000.

- What insurance premium would be actuarially fair for this situation?
- If the owner is willing to pay a premium of \$900, does she dislike risk or is she indifferent to risk?

EXERCISE 7.5

If individuals experience diminishing marginal utility from income it means that their utility function will resemble the total utility functions developed graphically in Section 6.2. Let us imagine specifically that if Y is income and U is utility, the individual gets utility from income according to the relation $U = \sqrt{Y}$.

- In a spreadsheet or using a calculator, calculate the amount of utility the individual gets for all income values running from \$1 to \$25.
- Graph the result with utility on the vertical axis and income on the horizontal axis, and verify from its shape that the marginal utility of income is declining.
- Using your calculations, how much utility will the individual get from \$4, \$9 and \$16?
- Suppose now that income results from a lottery and half of the time the individual gets \$4 and half of the time he gets \$16. How much utility will he get on average?
- Now suppose he gets \$10 each time with certainty. How much utility will he get from this?
- Since \$10 is exactly an average of \$4 and \$16, can you explain why \$10 with certainty gives him more utility than getting \$4 and \$16 each half of the time?

EXERCISE 7.6

In Question 7.5, suppose that the individual gets utility according to the relation $U = \frac{1}{2}Y$. Repeat the calculations for each part of the question and see if you can understand why the answers are different.

1. A different form of risk management is defined by the idea of risk spreading. Imagine that an oil supertanker has to be insured and the owner approaches one insurer. In the event of the tanker being ship-wrecked the damage caused by the resulting oil spill would be catastrophic – both to the environment and the insurance company. In this instance the insurance company may not benefit from the law of large numbers – it may not be insuring thousands of tankers and therefore would find it difficult to balance the potential claims with the annual insurance premiums. As a consequence, an insurer will spread the potential cost among other insurers – this is called risk spreading. The world's major insurers, such as *Lloyd's of London*, have hundreds of syndicates who each take on a small proportion of a big risk. These syndicates may again choose to subdivide their share among others, until the big risk becomes widely spread. In this way it is possible to insure against almost any event or possibility, no matter how large.

This page titled [7.7: Exercises for Chapter 7](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8: Production and cost

Chapter 8: Production and cost

In this chapter we will explore:

8.1	Efficient production
8.2	Time frames: The short run and the long run
8.3	Production in the short run
8.4	Costs in the short run
8.5	Fixed costs and sunk costs
8.6	Production and costs in the long run
8.7	Technological change and globalization
8.8	Clusters, externalities, learning by doing, and scope economies

8.1 Efficient production

Firms that fail to operate efficiently seldom survive. They are dominated by their competitors because the latter produce more efficiently and can sell at a lower price. The drive for profitability is everywhere present in the modern economy. Companies that promise more profit, by being more efficient, are valued more highly on the stock exchange. For example: In July of 2015 *Google* announced that, going forward, it would be more attentive to cost management in its numerous research endeavours that aim to bring new products to the marketplace. This policy, put in place by the Company's new Chief Financial Officer, was welcomed by investors who, as a result, bought up the stock. The Company's stock increased in value by 16% in one day – equivalent to about \$50 billion.

The remuneration of managers in virtually all corporations is linked to profitability. Efficient production, *a.k.a.* cost reduction, is critical to achieving this goal. In this chapter we will examine cost management and efficient production from the ground up – by exploring how a small entrepreneur brings his or her product to market in the most efficient way possible. As we shall see, efficient production and cost minimization amount to the same thing: Cost minimization is the financial reflection of efficient production.

Efficient production is critical in any budget-driven organization, not just in the private sector. Public institutions equally are, and should be, concerned with costs and efficiency.

Entrepreneurs employ factors of production (capital and labour) in order to transform raw materials and other inputs into goods or services. The relationship between output and the inputs used in the production process is called a production function. It specifies how much output can be produced with given combinations of inputs. A production function is not restricted to profit-driven organizations. Municipal road repairs are carried out with labour and capital. Students are educated with teachers, classrooms, computers, and books. Each of these is a production process.

Production function: a technological relationship that specifies how much output can be produced with specific amounts of inputs.

Economists distinguish between two concepts of efficiency: One is technological efficiency; the other is economic efficiency. To illustrate the difference, consider the case of auto assembly: the assembler could produce its vehicles either by using a large number of assembly workers and a plant that has a relatively small amount of machinery, or it could use fewer workers accompanied by more machinery in the form of robots. Each of these processes could be deemed technologically efficient, provided that there is no waste. If the workers without robots are combined with their capital to produce as much as possible, then that production process is technologically efficient. Likewise, in the scenario with robots, if the workers and capital are producing as much as possible, then that process too is efficient in the technological sense.

Technological efficiency means that the maximum output is produced with the given set of inputs.

Economic efficiency is concerned with more than just technological efficiency. Since the entrepreneur's goal is to make profit, she must consider which technologically efficient process best achieves that objective. More broadly, any budget-driven process should

focus on being economically efficient, whether in the public or private sector. An economically efficient production structure is the one that produces output at least cost.

Economic efficiency defines a production structure that produces output at least cost.

Auto-assembly plants the world over have moved to using robots during the last two decades. Why? The reason is not that robots were invented 20 years ago; they were invented long before that. The real reason is that, until recently, this technology was not economically efficient. Robots were too expensive; they were not capable of high-precision assembly. But once their cost declined and their accuracy increased they became economically efficient. The development of robots represented technological progress. When this progress reached a critical point, entrepreneurs embraced it.

To illustrate the point further, consider the case of garment assembly. There is no doubt that engineers could make robots capable of joining the pieces of fabric that form garments. This is not beyond our technological abilities. Why, then, do we not have such capital-intensive production processes for garment making, similar to the production process chosen by vehicle producers? The answer is that, while such a concept could be technologically efficient, it would not be economically efficient. It is more profitable to use large amounts of labour and relatively traditional machines to assemble garments, particularly when labour in Asia costs less and the garments can be shipped back to Canada inexpensively. Containerization and scale economies in shipping mean that a garment can be shipped to Canada from Asia for a few cents per unit.

Efficiency in production is not limited to the manufacturing sector. Farmers must choose the optimal combination of labour, capital and fertilizer to use. In the health and education sectors, efficient supply involves choices on how many high- and low-skill workers to employ, how much traditional physical capital to use, how much information technology to use, based upon the productivity and cost of each. Professors and physicians are costly inputs. When they work with new technology (capital) they become more efficient at performing their tasks: It is less costly to have a single professor teach in a 300-seat classroom that is equipped with the latest technology, than have several professors each teaching 60-seat classes with chalk and a blackboard.

8.2 The time frame

We distinguish initially between the short run and the long run. When discussing technological change, we use the term very long run. These concepts have little to do with clocks or calendars; rather, they are defined by the degree of flexibility an entrepreneur or manager has in her production process. A key decision variable is capital.

A customary assumption is that a producer can hire more labour immediately, if necessary, either by taking on new workers (since there are usually some who are unemployed and looking for work), or by getting the existing workers to work longer hours. In contrast, getting new capital in place is usually more time consuming: The entrepreneur may have to place an order for new machinery, which will involve a production and delivery time lag. Or she may have to move to a more spacious location in order to accommodate the added capital. Whether this calendar time is one week, one month, or one year is of no concern to us. We define the long run as a period of sufficient length to enable the entrepreneur to adjust her capital stock, whereas in the short run at least one factor of production is fixed. Note that it matters little whether it is labour or capital that is fixed in the short run. A software development company may be able to install new capital (computing power) instantaneously but have to train new developers. In such a case capital is variable and labour is fixed in the short run. The definition of the short run is that one of the factors is fixed, and in our examples we will assume that it is capital.

Short run: a period during which at least one factor of production is fixed. If capital is fixed, then more output is produced by using additional labour.

Long run: a period of time that is sufficient to enable all factors of production to be adjusted.

Very long run: a period sufficiently long for new technology to develop.

8.3 Production in the short run

Black Diamond Snowboards (BDS) is a start-up snowboard producing enterprise. Its founder has invented a new lamination process that gives extra strength to his boards. He has set up a production line in his garage that has four workstations: Laminating, attaching the steel edge, waxing, and packing.

With this process in place, he must examine how productive his firm can be. After extensive testing, he has determined exactly how his productivity depends upon the number of workers. If he employs only one worker, then that worker must perform several tasks, and will encounter 'down time' between workstations. Extra workers would therefore not only increase the total output; they could, in addition, increase output *per worker*. He also realizes that once he has employed a critical number of workers, additional

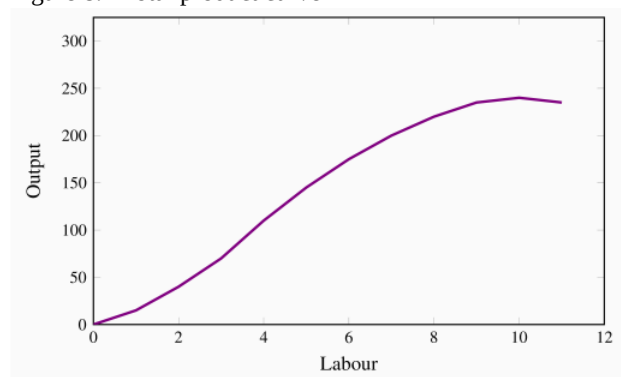
workers may not be so productive: Because they will have to share the fixed amount of machinery in his garage, they may have to wait for another worker to finish using a machine. At such a point, the productivity of his plant will begin to fall off, and he may want to consider capital expansion. But for the moment he is constrained to using this particular assembly plant. Testing leads him to formulate the relationship between workers and output that is described in Table 8.1.

Table 8.1 Snowboard production and productivity

1	2	3	4	5
Workers	Output	Marginal	Average	Stages of
	(<i>TP</i>)	product	product	production
		(<i>MP_L</i>)	(<i>AP_L</i>)	
0	0			<i>MP_L increasing</i>
1	15	15	15	
2	40	25	20	
3	70	30	23.3	
4	110	40	27.5	<i>MP_L positive and declining</i>
5	145	35	29	
6	175	30	29.2	
7	200	25	28.6	
8	220	20	27.5	
9	235	15	26.1	<i>MP_L negative</i>
10	240	5	24.0	
11	235	-5	21.4	

By increasing the number of workers in the plant, BDS produces more boards. The relationship between these two variables in columns 1 and 2 in the table is plotted in Figure 8.1. This is called the total product function (*TP*), and it defines the output produced with different amounts of labour in a plant of fixed size.

Figure 8.1 Total product curve



Output increases with the amount of labour used. Initially the increase in output due to using more labour is high, subsequently it is lower. The initial phase characterizes increasing productivity, the later phase defines declining productivity.

Total product is the relationship between total output produced and the number of workers employed, for a given amount of capital.

This relationship is positive, indicating that more workers produce more boards. But the curve has an interesting pattern. In the initial expansion of employment it becomes progressively steeper – its curvature is slightly convex; following this phase the

function's increase becomes progressively less steep – its curvature is concave. These different stages in the TP curve tell us a great deal about productivity in BDS. To see this, consider the additional number of boards produced by each worker. The first worker produces 15. When a second worker is hired, the total product rises to 40, so the additional product attributable to the second worker is 25. A third worker increases output by 30 units, and so on. We refer to this additional output as the marginal product (MP) of an additional worker, because it defines the incremental, or marginal, contribution of the worker. These values are entered in column 3.

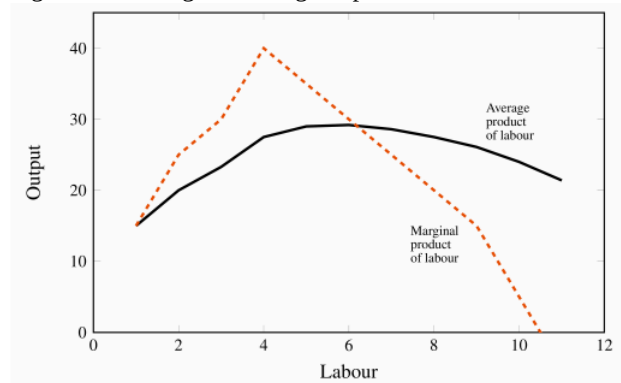
More generally the MP of labour is defined as the change in output divided by the change in the number of units of labour employed. Using, as before, the Greek capital delta (Δ) to denote a change, we can define

$$MP_L = \frac{\text{Change in output produced}}{\text{Change in labour employed}} = \frac{\Delta Q}{\Delta L}$$

In this example the change in labour is one unit at each stage and hence the marginal product of labour is simply the corresponding change in output. It is also the case that the MP_L is the slope of the TP curve – the change in the value on the vertical axis due to a change in the value of the variable on the horizontal axis.

Marginal product of labour is the addition to output produced by each additional worker. It is also the slope of the total product curve.

Figure 8.2 Average and marginal product curves



The productivity curves initially rise and then decline, reflecting increasing and decreasing productivity. The MP_L curves must intersect the AP_L curve at the maximum of the AP_L : The average must increase if the marginal exceeds the average and must decline if the marginal is less than the average.

During the initial stage of production expansion, the marginal product of each worker is increasing. It increases from 15 to 40 as BDS moves from having one employee to four employees. This increasing MP is made possible by the fact that each worker is able to spend more time at his workstation, and less time moving between tasks. But, at a certain point in the employment expansion, the MP reaches a maximum and then begins to tail off. At this stage – in the concave region of the TP curve – additional workers continue to produce additional output, but at a diminishing rate. For example, while the fourth worker adds 40 units to output, the fifth worker adds 35, the sixth worker 30, and so on. This declining MP is due to the constraint of a fixed number of machines: All workers must share the same capital. The MP function is plotted in Figure 8.2.

The phenomenon we have just described has the status of a law in economics: The law of diminishing returns states that, in the face of a fixed amount of capital, the contribution of additional units of a variable factor must eventually decline.

Law of diminishing returns: when increments of a variable factor (labour) are added to a fixed amount of another factor (capital), the marginal product of the variable factor must eventually decline.

The relationship between Figures 8.1 and 8.2 should be noted. First, the MP_L reaches a maximum at an output of 4 units – where the slope of the TP curve is greatest. The MP_L curve remains positive beyond this output, but declines: The TP curve reaches a maximum when the tenth unit of labour is employed. An eleventh unit actually reduces total output; therefore, the MP of this eleventh worker is negative! In Figure 8.2, the MP curve becomes negative at this point. The garage is now so crowded with workers that they are beginning to obstruct the operation of the production process. Thus the producer would never employ an eleventh unit of labour.

Next, consider the information in the fourth column of the table. It defines the average product of labour (AP_L)—the amount of output produced, on average, by workers at different employment levels:

$$AP_L = \frac{\text{Total output produced}}{\text{Total amount of labour employed}} = \frac{Q}{L}.$$

This function is also plotted in Figure 8.2. Referring to the table: The *AP* column indicates, for example, that when two units of labour are employed and forty units of output are produced, the average production level of each worker is 20 units ($=40/2$). When three workers produce 70 units, their average production is 23.3 ($=70/3$), and so forth. Like the *MP* function, this one also increases and subsequently decreases, reflecting exactly the same productivity forces that are at work on the *MP* curve.

Average product of labour is the number of units of output produced per unit of labour at different levels of employment.

The *AP* and *MP* functions intersect at the point where the *AP* is at its peak. This is no accident, and has a simple explanation. Imagine a softball player who is batting .280 coming into today's game—she has been hitting her way onto base 28 percent of the time when batting, so far this season. This is her average product, *AP*.

In today's game, if she bats .500 (hits her way to base on half of her at-bats), then she will improve her average. Today's batting (*MP*) at .500 therefore pulls up the season's *AP*. Accordingly, whenever the *MP* exceeds the *AP*, the *AP* is pulled up. By the same reasoning, if her *MP* is less than the season average, her average will be pulled down. It follows that the two functions must intersect at the peak of the *AP* curve. To summarize:

If the MP exceeds the AP, then the AP increases;
If the MP is less than the AP, then the AP declines.

While the owner of BDS may understand his productivity relations, his ultimate goal is to make profit, and for this he must figure out how productivity translates into cost.

8.4 Costs in the short run

The cost structure for the production of snowboards at Black Diamond is illustrated in Table 8.2. Employees are skilled and are paid a weekly wage of \$1,000. The cost of capital is \$3,000 and it is fixed, which means that it does not vary with output. As in Table 8.1, the number of employees and the output are given in the first two columns. The following three columns define the capital costs, the labour costs, and the sum of these in producing different levels of output. We use the terms fixed, variable, and total costs to define the cost structure of a firm. Fixed costs do not vary with output, whereas variable costs do, and total costs are the sum of fixed and variable costs. To keep this example as simple as possible, we will ignore the cost of raw materials. We could add an additional column of costs, but doing so will not change the conclusions.

Table 8.2 Snowboard production costs

Workers	Output	Capital	Labour	Total	Average	Average	Average	Marginal
		cost	cost	costs	fixed	variable	total	cost
		fixed	variable		cost	cost	cost	
0	0	3,000	0	3,000				
1	15	3,000	1,000	4,000	200.0	66.7	266.7	66.7
2	40	3,000	2,000	5,000	75.0	50.0	125.0	40.0
3	70	3,000	3,000	6,000	42.9	42.9	85.7	33.3
4	110	3,000	4,000	7,000	27.3	36.4	63.6	25.0
5	145	3,000	5,000	8,000	20.7	34.5	55.2	28.6
6	175	3,000	6,000	9,000	17.1	34.3	51.4	33.3
7	200	3,000	7,000	10,000	15.0	35.0	50.0	40.0
8	220	3,000	8,000	11,000	13.6	36.4	50.0	50.0
9	235	3,000	9,000	12,000	12.8	38.3	51.1	66.7
10	240	3,000	10,000	13,000	12.5	41.7	54.2	200.0

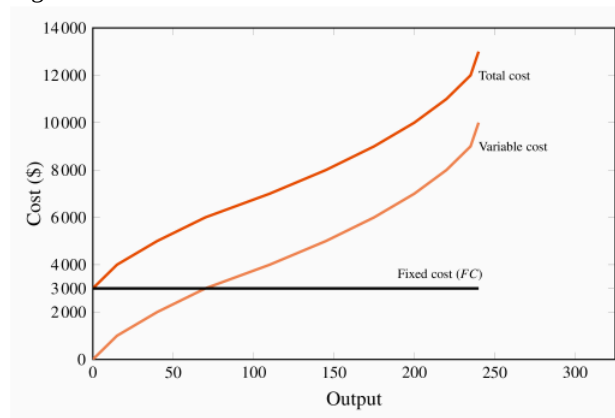
Fixed costs are costs that are independent of the level of output.

Variable costs are related to the output produced.

Total cost is the sum of fixed cost and variable cost.

Total costs are illustrated in Figure 8.3 as the vertical sum of variable and fixed costs. For example, Table 8.2 indicates that the total cost of producing 220 units of output is the sum of \$3,000 in fixed costs plus \$8,000 in variable costs. Therefore, at the output level 220 on the horizontal axis in Figure 8.3, the sum of the cost components yields a value of \$11,000 that forms one point on the total cost curve. Performing a similar calculation for every possible output yields a series of points that together form the complete total cost curve.

Figure 8.3 Total cost curves



Total cost is the vertical sum of the variable and fixed costs.

Average costs are given in the next three columns of Table 8.2. Average cost is the cost per unit of output, and we can define an average cost corresponding to each of the fixed, variable, and total costs defined above. Average fixed cost (*AFC*) is the total fixed cost divided by output; average variable cost (*AVC*) is the total variable cost divided by output; and average total cost (*ATC*) is the total cost divided by output.

<i>AFC</i>	$= (\text{Fixed cost})/Q = FC/Q$
<i>AVC</i>	$= (\text{Total variable costs})/Q = TVC/Q$
<i>ATC</i>	$= AFC + AVC$

Average fixed cost is the total fixed cost per unit of output.

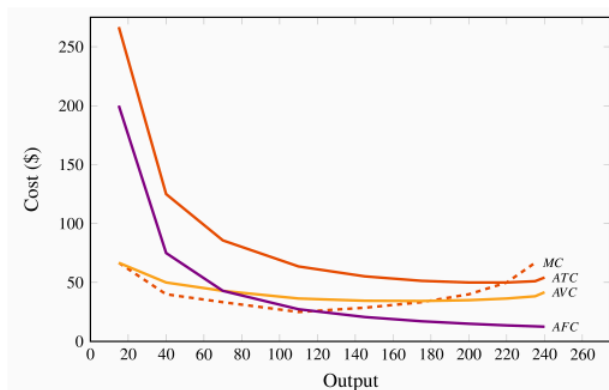
Average variable cost is the total variable cost per unit of output.

Average total cost is the sum of all costs per unit of output.

The productivity-cost relationship

Consider the **average variable cost - average product relationship**, as developed in column 7 of Table 8.2; its corresponding variable cost curve is plotted in Figure 8.4. In this example, *AVC* first decreases and then increases. The intuition behind its shape is straightforward (and realistic) if you have understood why productivity varies in the short run: The variable cost, which represents the cost of labour, is constant per unit of labour, because the wage paid to each worker does not change. However, each worker's productivity varies. Initially, when we hire more workers, they become more productive, perhaps because they have less 'down time' in switching between tasks. This means that the labour costs per snowboard must decline. At some point, however, the law of diminishing returns sets in: As before, each additional worker is paid a constant amount, but as productivity declines the labour cost per snowboard increases.

Figure 8.4 Average and marginal cost curves



The MC intersects the ATC and AVC at their minimum values. The AFC declines indefinitely as fixed costs are spread over a greater output.

In this numerical example the AP is at a maximum when six units of labour are employed and output is 175. This is also the point where the AVC is at a minimum. This maximum/minimum relationship is also illustrated in Figures 8.2 and 8.4.

Now consider the **marginal cost - marginal product relationship**. The marginal cost (MC) defines the cost of producing one more unit of output. In Table 8.2, the marginal cost of output is given in the final column. It is the additional cost of production divided by the additional number of units produced. For example, in going from 15 units of output to 40, total costs increase from \$4,000 to \$5,000. The MC is the cost of those additional units divided by the number of additional units. In this range of output, MC is $\$1,000/25 = \40 . We could also calculate the MC as the addition to variable costs rather than the addition to total costs, because the *addition* to each is the same—fixed costs are fixed. Hence:

MC	$= \frac{\text{Change in total costs}}{\text{Change in output produced}} = \frac{\Delta TC}{\Delta Q}$
	$= \frac{\text{Change in variable costs}}{\text{Change in output produced}} = \frac{\Delta TVC}{\Delta Q}$

Marginal cost of production is the cost of producing each additional unit of output.

Just as the behaviour of the AVC curve is determined by the AP curve, so too the behaviour of the MC is determined by the MP curve. When the MP of an additional worker exceeds the MP of the previous worker, this implies that the cost of the additional output produced by the last worker hired must be declining. To summarize:

- If the marginal product of labour increases, then the marginal cost of output declines;*
- If the marginal product of labour declines, then the marginal cost of output increases.*

In our example, the MP_L reaches a maximum when the fourth unit of labour is employed (or 110 units of output are produced), and this also is where the MC is at a minimum. This illustrates that the *marginal cost reaches a minimum at the output level where the marginal product reaches a maximum*.

The average total cost is the sum of the fixed cost per unit of output and the variable cost per unit of output. Typically, fixed costs are the dominant component of total costs at low output levels, but become less dominant at higher output levels. Unlike average variable costs, note that the average fixed cost must always decline with output, because a fixed cost is being spread over more units of output. Hence, when the ATC curve eventually increases, it is because the increasing variable cost component eventually dominates the declining AFC component. In our example, this occurs when output increases from 220 units (8 workers) to 235 (9 workers).

Finally, observe the interrelationship between the MC curve on the one hand and the ATC and AVC on the other. Note from Figure 8.4 that the MC cuts the AVC and the ATC at the minimum point of each of the latter. The logic behind this pattern is analogous to the logic of the relationship between marginal and average product curves: When the cost of an additional unit of output is less than the average, this reduces the average cost; whereas, if the cost of an additional unit of output is above the average, this raises the average cost. This must hold true regardless of whether we relate the MC to the ATC or the AVC .

*When the marginal cost is less than the average cost, the average cost must decline;
When the marginal cost exceeds the average cost, the average cost must increase.*

Notation: We use both the abbreviations ATC and AC to denote average total cost. The term 'average cost' is understood in economics to include both fixed and variable costs.

Teams and services

The choice faced by the producer in the example above is slightly 'stylized', yet it still provides an appropriate rule for analyzing hiring decisions. In practice, it is quite difficult to isolate or identify the marginal product of an individual worker. One reason is that individuals work in teams within organizations. The accounting department, the marketing department, the sales department, the assembly unit, the chief executive's unit are all composed of teams. Adding one more person to human resources may have no impact on the number of units of output produced by the company in a measurable way, but it may influence worker morale and hence longer-term productivity. Nonetheless, if we consider expanding, or contracting, any one department within an organization, management can attempt to estimate the net impact of additional hires (or layoffs) on the contribution of each team to the firm's profitability. Adding a person in marketing may increase sales, laying off a person in research and development may reduce costs by more than it reduces future value to the firm. In practice this is what firms do: they attempt to assess the contribution of each team in their organization to costs and revenues, and on that basis determine the appropriate number of employees.

The manufacturing sector of the macro economy is dominated, sizewise, by the services sector. But the logic that drives hiring decisions, as developed above, applies equally to services. For example, how does a law firm determine the optimal number of paralegals to employ per lawyer? How many nurses are required to support a surgeon? How many university professors are required to teach a given number of students?

All of these employment decisions involve optimization at the margin. The goal of the decision maker is not always profit, but she should attempt to estimate the cost and value of adding personnel at the margin.

8.5 Fixed costs and sunk costs

The distinction between fixed and variable costs is important for producers who are not making a profit. If a producer has committed himself to setting up a plant, then he has made a decision to incur a fixed cost. Having done this, he must now decide on a production strategy that will maximize profit. However, the price that consumers are willing to pay may not be sufficient to yield a profit. So, if Black Diamond Snowboards cannot make a profit, should it shut down? The answer is that if it can cover its variable costs, *having already incurred its fixed costs*, it should stay in production, at least temporarily. By covering the variable cost of its operation, Black Diamond is at least earning some return. A sunk cost is a fixed cost that has already been incurred and cannot be recovered. But if the pressures of the marketplace are so great that the total costs cannot be covered in the longer run, then this is not a profitable business and the firm should close its doors.

Is a fixed cost always a sunk cost? No: Any production that involves capital will incur a fixed cost component. Such capital can be financed in several ways however: It might be financed on a very short-term lease basis, or it might have been purchased by the entrepreneur. If it is leased on a month-to-month basis, an unprofitable entrepreneur who can only cover variable costs (and who does not foresee better market conditions ahead) can exit the industry quickly – by not renewing the lease on the capital. But an individual who has actually purchased equipment that cannot readily be resold has essentially sunk money into the fixed cost component of his production. This entrepreneur should continue to produce as long as he can cover variable costs.

Sunk cost is a fixed cost that has already been incurred and cannot be recovered, even by producing a zero output.

R & D as a sunk cost

Sunk costs in the modern era are frequently in the form of research and development costs, not the cost of building a plant or purchasing machinery. The prototypical example is the pharmaceutical industry, where it is becoming progressively more challenging to make new drug breakthroughs – both because the 'easier' breakthroughs have already been made, and because it is necessary to meet tighter safety conditions attaching to new drugs. Research frequently leads to drugs that are not sufficiently effective in meeting their target. As a consequence, the pharmaceutical sector regularly writes off hundreds of millions of dollars of lost sunk costs – unfruitful research and development.

Finally, we need to keep in mind the opportunity costs of running the business. The owner pays himself a salary, and ultimately he must recognize that the survival of the business should not depend upon his drawing a salary that is less than his opportunity cost.

As developed in Section 7.2, if he underpays himself in order to avoid shutting down, he might be better off in the long run to close the business and earn his opportunity cost elsewhere in the marketplace.

A dynamic setting

We need to ask why it might be possible to cover all costs in a longer run horizon, while in the near-term costs are not covered. The principal reason is that demand may grow, particularly for a new product. For example, in 2019 numerous cannabis producing firms were listed on the Canadian Securities Exchange, and collectively were valued at about fifty billion dollars. None had revenues that covered costs, yet investors poured money into this sector. Investors evidently envisaged that the market for legal cannabis would grow. As of 2020 it appears that these investors were excessively optimistic. Sales growth has been slow and stock valuations have plummeted.

8.6 Long-run production and costs

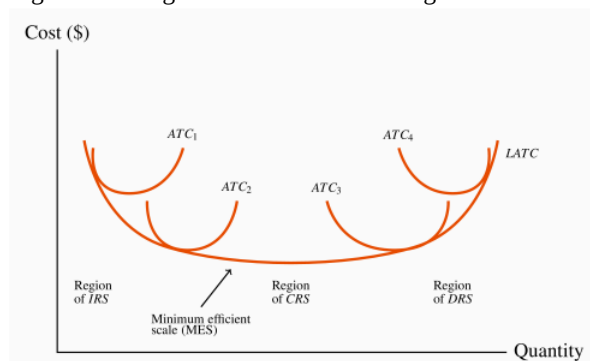
The snowboard manufacturer we portray produces a relatively low level of output; in reality, millions of snowboards are produced each year in the global market. Black Diamond Snowboards may have hoped to get a start by going after a local market—the "free-ride" teenagers at Mont Sainte Anne in Quebec or at Fernie in British Columbia. If this business takes off, the owner must increase production, take the business out of his garage and set up a larger-scale operation. But how will this affect his cost structure? Will he be able to produce boards at a lower cost than when he was producing a very limited number of boards each season? Real-world experience would indicate yes.

Production costs almost always decline when the *scale* of the operation initially increases. We refer to this phenomenon simply as economies of scale. There are several reasons why scale economies are encountered. One is that production flows can be organized in a more efficient manner when more is being produced. Another is that the opportunity to make greater use of task specialization presents itself; for example, Black Diamond Snowboards may be able to subdivide tasks within the laminating and packaging stations. With a larger operating scale the replacement of labor with capital may be economically efficient. If scale economies do define the real world, then a bigger plant—one that is geared to produce a higher level of output—should have an average total cost curve that is "lower" than the cost curve corresponding to the smaller scale of operation we considered in the example above.

Average costs in the long run

Figure 8.5 illustrates a possible relationship between the ATC curves for four different scales of operation. ATC_1 is the average total cost curve associated with a small-sized plant; think of it as the plant built in the entrepreneur's garage. ATC_2 is associated with a somewhat larger plant, perhaps one she has put together in a rented industrial or commercial space. The further a cost curve is located to the right of the diagram the larger the production facility it defines, given that output is measured on the horizontal axis. If there are economies associated with a larger scale of operation, then the average costs associated with producing larger outputs in a larger plant should be lower than the average costs associated with lower outputs in a smaller plant, assuming that the plants are producing the output levels they were designed to produce. For this reason, the cost curve ATC_2 and the cost curve ATC_3 each have a segment that is lower than the lowest segment on ATC_1 . However, in Figure 8.5 the cost curve ATC_4 has moved upwards. What behaviours are implied here?

Figure 8.5 Long-run and short-run average costs



The long-run ATC curve, $LATC$, is the lower envelope of all short-run ATC curves. It defines the least cost per unit of output when all inputs are variable. Minimum efficient scale is that output level at which the $LATC$ is a minimum, indicating that further increases in the scale of production will not reduce unit costs.

In many production environments, beyond some large scale of operation, it becomes increasingly difficult to reap further cost reductions from specialization, organizational economies, or marketing economies. At such a point, the scale economies are effectively exhausted, and larger plant sizes no longer give rise to lower (short-run) ATC curves. This is reflected in the similarity of the ATC_2 and the ATC_3 curves. The pattern suggests that we have almost exhausted the possibilities of further scale advantages once we build a plant size corresponding to ATC_2 . Consider next what is implied by the position of the ATC_4 curve relative to the ATC_2 and ATC_3 curves. The relatively higher position of the ATC_4 curve implies that unit costs will be higher in a yet larger plant. Stated differently: If we increase the scale of this firm to extremely high output levels, we are actually encountering diseconomies of scale. Diseconomies of scale imply that unit costs increase as a result of the firm's becoming too large: Perhaps co-ordination difficulties have set in at the very high output levels, or quality-control monitoring costs have risen. These coordination and management difficulties are reflected in increasing unit costs in the long run.

The terms increasing, constant, and decreasing returns to scale underlie the concepts of scale economies and diseconomies: Increasing returns to scale (IRS) implies that, when all inputs are increased by a given proportion, output increases more than proportionately. Constant returns to scale (CRS) implies that output increases in direct proportion to an equal proportionate increase in all inputs. Decreasing returns to scale (DRS) implies that an equal proportionate increase in all inputs leads to a less than proportionate increase in output.

Increasing returns to scale implies that, when all inputs are increased by a given proportion, output increases more than proportionately.

Constant returns to scale implies that output increases in direct proportion to an equal proportionate increase in all inputs.

Decreasing returns to scale implies that an equal proportionate increase in all inputs leads to a less than proportionate increase in output.

These are pure production function relationships, but, if the prices of inputs are fixed for producers, they translate directly into the various cost structures illustrated in Figure 8.5. For example, if a 40% increase in capital and labour use allows for better production flows than when in the smaller plant, and therefore yields more than a 40% increase in output, this implies that the cost per snowboard produced must fall in the new plant. In contrast, if a 40% increase in capital and labour leads to say just a 30% increase in output, then the cost per snowboard in the new larger plant must be higher. Between these extremes, there may be a range of relatively constant unit costs, corresponding to where the production relation is subject to constant returns to scale. In Figure 8.5, the falling unit costs output region has increasing returns to scale, the region that has relatively constant unit costs has constant returns to scale, and the increasing cost region has decreasing returns to scale.

Increasing returns to scale characterize businesses with large initial costs and relatively low costs of producing each unit of output. Computer chip manufacturers, pharmaceutical manufacturers, vehicle rental agencies, booking agencies such as *booking.com* or *hotels.com*, intermediaries such as *airbnb.com*, even brewers, all benefit from scale economies. In the beer market, brewing, bottling and shipping are all low-cost operations relative to the capital cost of setting up a brewery. Consequently, we observe surprisingly few breweries in any brewing company, even in large land-mass economies such as Canada or the US.

In addition to the four short-run average total cost curves, Figure 8.5 contains a curve that forms an envelope around the bottom of these short-run average cost curves. This envelope is the long-run average total cost ($LATC$) curve, because it defines average cost as we move from one plant size to another. Remember that in the long run both labour and capital are variable, and as we move from one short-run average cost curve to another, that is exactly what happens—all factors of production are variable. Hence, the collection of short-run cost curves in Figure 8.5 provides the ingredients for a long-run average total cost curve¹.

$$LATC = (\text{Long-run total costs})/Q = LTC/Q$$

Long-run average total cost is the lower envelope of all the short-run ATC curves.

The particular range of output on the $LATC$ where it begins to flatten out is called the range of minimum efficient scale. This is an important concept in industrial policy, as we shall see in later chapters. At such an output level, the producer has expanded sufficiently to take advantage of virtually all the scale economies available.

Minimum efficient scale defines a threshold size of operation such that scale economies are almost exhausted.

In view of this discussion and the shape of the $LATC$ in Figure 8.5, it is obvious that economies of scale can also be defined in terms of the curvature of the $LATC$. Where the $LATC$ declines there are IRS, where the $LATC$ is flat there are CRS, where the $LATC$ slopes upward there are DRS.

Table 8.3 LATC elements for two plants (thousands \$)

Q	AFC_1	$MC_1 = AVC_1$	ATC_1	AFC_2	$MC_2 = AVC_2$	ATC_2
20	50	30	80	100	25	125
40	25	30	55	50	25	75
60	16.67	30	46.67	33.33	25	58.33
80	12.5	30	42.5	25	25	50
100	10	30	40	20	25	45
120	8.33	30	38.33	16.67	25	41.67
140	7.14	30	37.14	14.29	25	39.29
160	6.25	30	36.25	12.5	25	37.5
180	5.56	30	35.56	11.11	25	36.11
200	5	30	35	10	25	35
220	4.55	30	34.55	9.09	25	34.09
240	4.17	30	34.17	8.33	25	33.33
260	3.85	30	33.85	7.69	25	32.69
280	3.57	30	33.57	7.14	25	32.14

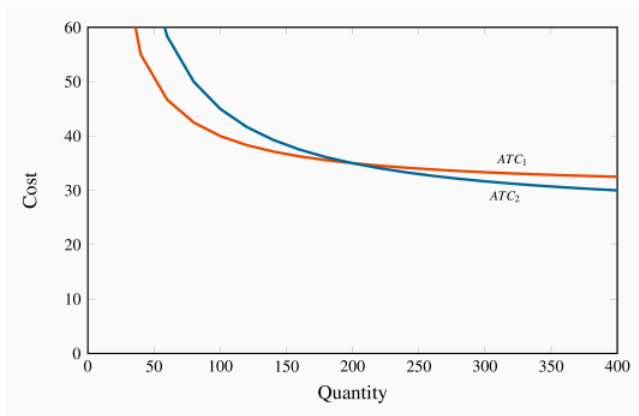
Plant 1 $FC = \$1m$. Plant 2 $FC = \$2m$. For $Q < 200$, $ATC_1 < ATC_2$; for $Q > 200$, $ATC_1 > ATC_2$; and for $Q = 200$, $ATC_1 = ATC_2$. LATC defined by data in bold font.

Long-run costs – a simple numerical example

Kitt is an automobile designer specializing in the production of off-road vehicles sold to a small clientele. He has a choice of two (and only two) plant sizes; one involving mainly labour and the other employing robots extensively. The set-up (i.e. fixed) costs of these two assembly plants are \$1 million and \$2 million respectively. The advantage to having the more costly plant is that the pure production costs (variable costs) are less. The cost components are defined in Table 8.3. The variable cost (equal to the marginal cost here) is \$30,000 in the plant that relies primarily on labour, and \$25,000 in the plant that has robots. The ATC for each plant size is the sum of AFC and AVC. The AFC declines as the fixed cost is spread over more units produced. The variable cost per unit is constant in each case. By comparing the fourth and final columns, it is clear that the robot-intensive plant has lower costs if it produces a large number of vehicles. At an output of 200 vehicles the average costs in each plant are identical: The higher fixed costs associated with the robots are exactly offset by the lower variable costs at this output level.

The ATC curve corresponding to each plant size is given in Figure 8.6. There are two short-run ATC curves. The positions of these curves indicate that if the manufacturer believes he can produce at least 200 vehicles his unit costs will be less with the plant involving robots; but at output levels less than this his unit costs would be less in the labour-intensive plant.

Figure 8.6 LATC for two plants in \$000



The long-run average cost curve for this producer is the lower envelope of these two cost curves: ATC_1 up to output 200 and ATC_2 thereafter. Two features of this example are to be noted. First we do not encounter decreasing returns – the $LATC$ curve never increases. ATC_1 tends asymptotically to a lower bound of \$30, while ATC_2 tends towards \$25. Second, in the interests of simplicity we have assumed just two plant sizes are possible. With more possibilities on the introduction of robots we could imagine more short-run ATC curves which would form the lower-envelope $LATC$.

8.7 Technological change: globalization and localization

Technological change represents innovation that can reduce the cost of production or bring new products on line. As stated earlier, the very long run is a period that is sufficiently long for new technology to evolve and be implemented.

Technological change represents innovation that can reduce the cost of production or bring new products on line.

Technological change has had an enormous impact on economic life for several centuries. It is not something that is defined in terms of the recent telecommunications revolution. The industrial revolution began in eighteenth century Britain. It was accompanied by a less well-recognized, but equally important, agricultural revolution. The improvement in cultivation technology, and ensuing higher yields, freed up enough labour to populate the factories that were the core of the industrial revolution². The development and spread of mechanical power dominated the nineteenth century, and the mass production line of Henry Ford in autos or Andrew Carnegie in steel heralded in the twentieth century.

Globalization

The modern communications revolution has reduced costs, just like its predecessors. But it has also greatly sped up globalization, the increasing integration of national markets.

Globalization is the tendency for international markets to be ever more integrated.

Globalization has several drivers: lower transportation and communication costs; reduced barriers to trade and capital mobility; the spread of new technologies that facilitate cost and quality control; different wage rates between developed and less developed economies. New technology and better communications have been critical in both increasing the minimum efficient scale of operation and reducing diseconomies of scale; they facilitate the efficient management of large companies.

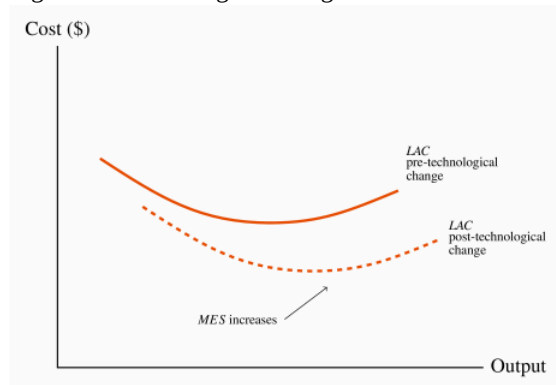
The continued reduction in trade barriers in the post-World War II era has also meant that the effective marketplace has become the globe rather than the national economy for many products. Companies like *Apple*, *Microsoft*, and *Facebook* are visible worldwide. Globalization has been accompanied by the collapse of the Soviet Union, the adoption of an outward looking philosophy on the part of China, and an increasing role for the market place in India. These developments together have facilitated the outsourcing of much of the West's manufacturing to lower-wage economies.

But new technology not only helps existing companies grow large; it also enables new ones to start up. It is now cheaper for small producers to manage their inventories and maintain contact with their own suppliers.

The impact of technology is to reduce the cost of production, hence it will lower the average cost curve in both the short and long run. First, decreasing returns to scale become less probable due to improved communications, so the upward sloping section of the $LRAC$ curve may disappear altogether. Second, capital costs are now lower than in earlier times, because much modern technology has transformed some fixed costs into variable cost: Both software and hardware functions can be subcontracted to specialty firms, who in turn may use cloud computing services, and it is thus no longer necessary to have a substantial in-house computing

department. The use of almost-free software such as *Skype*, *Hangouts* and *WhatsApp* reduces communication costs. Advertising on social media is more effective and less costly than in traditional hard-print form. Hiring may be cheaper through *LinkedIn* than through a traditional human resources department. These developments may actually reduce the minimum efficient scale of operation because they reduce the need for large outlays on fixed capital. On the other hand, changes in technology may induce producers to use more capital and less labor, with the passage of time. That would increase the minimum efficient scale. An example of this phenomenon is in mining, or tunnel drilling, where capital investment per worker is greater than when technology was less developed. A further example is the introduction of robotic assistants in *Amazon* warehouses. In these scenarios the minimum efficient scale should increase, and that is illustrated in Figure 8.7.

Figure 8.7 Technological change and LAC



Technological change reduces the unit production cost for any output produced and may also increase the minimum efficient scale (MES) threshold.

The local diffusion of technology

The impacts of technological change are not just evident in a global context. Technological change impacts every sector of the domestic economy. For example, the modern era in dentistry sees specialists in root canals (endodontists) performing root canals in the space of a single hour with the help of new technology; dental implants into bone, as an alternative to dentures, are commonplace; crowns can be machined with little human intervention; and X-rays are now performed with about one hundredth of the power formerly required. These technologies spread and are adopted through several channels. Dental practices do not usually compete on the basis of price, but if they do not adopt best practices and new technologies, then community word of mouth will see patients shifting to more efficient operators.

Some technological developments are protected by patents. But patent protection rarely inhibits new and more efficient practices that in some way mimic patent breakthroughs.

8.8 Clusters, learning by doing, scope economies

Clusters

The phenomenon of a grouping of firms that specialize in producing related products is called a cluster. For example, Ottawa has more than its share of software development firms; Montreal has a disproportionate share of Canada's pharmaceutical producers and electronic game developers; Calgary has its 'oil patch'; Hollywood has movies; Toronto is Canada's financial capital, San Francisco and Seattle are leaders in new electronic products. Provincial and state capitals have most of their province's bureaucracy. Clusters give rise to externalities, frequently in the form of ideas that flow between firms, which in turn result in cost reductions and new products.

Cluster: a group of firms producing similar products, or engaged in similar research.

The most famous example of clustering is Silicon Valley, surrounding San Francisco, in California, the original high-tech cluster. The presence of a large group of firms with a common focus serves as a signal to workers with the right skill set that they are in demand in such a region. Furthermore, if these clusters are research oriented, as they frequently are, then knowledge spillovers benefit virtually all of the contiguous firms; when workers change employers, they bring their previously-learned skills with them; on social occasions, friends may chat about their work and interests and share ideas. This is a positive externality.

Learning by doing

Learning from production-related experiences frequently reduces costs: The accumulation of knowledge that is associated with having produced a large volume of output over a considerable time period enables managers to implement more efficient production methods and avoid errors. We give the term learning by doing to this accumulation of knowledge.

Examples abound, but the best known may be the continual improvement in the capacity of computer chips, whose efficiency has doubled about every eighteen months for several decades – a phenomenon known as Moore's Law. As *Intel Corporation* continues to produce chips it learns how to produce each succeeding generation of chips at lower cost. Past experience is key. Economies of scale and learning by doing therefore may not be independent: Large firms usually require time to grow or to attain a dominant role in their market, and this time and experience enables them to produce at lower cost. This lower cost in turn can solidify their market position further.

Learning by doing can reduce costs. A longer history of production enables firms to accumulate knowledge and thereby implement more efficient production processes.

Economies of scope

Economies of scope define a production process if the production of multiple products results in lower unit costs per product than if those products were produced alone. Scope economies, therefore, define the returns or cost reductions associated with broadening a firm's product range.

Corporations like Proctor and Gamble do not produce a single product in their health line; rather, they produce first aid, dental care, and baby care products. Cable companies offer their customers TV, high-speed Internet, and telephone services either individually or packaged. A central component of some new-economy multi-product firms is a technology platform that can be used for multiple purposes. We shall analyze the operation of these firms in more detail in Chapter 11.

Economies of scope occur if the unit cost of producing particular products is less when combined with the production of other products than when produced alone.

A platform is a hardware-cum-software capital installation that has multiple production capabilities

Conclusion

Efficient production is critical to the survival of firms. Firms that do not adopt the most efficient production methods are likely to be left behind by their competitors. Efficiency translates into cost considerations, and the structure of costs in turn has a major impact on market type. Some sectors of the economy have very many firms (the restaurant business or the dry-cleaning business), whereas other sectors have few (internet providers or airlines). We will see in the following chapters how market structures depend critically upon the concept of scale economies that we have developed here.

Key Terms

Production function: a technological relationship that specifies how much output can be produced with specific amounts of inputs.

Technological efficiency means that the maximum output is produced with the given set of inputs.

Economic efficiency defines a production structure that produces output at least cost.

Short run: a period during which at least one factor of production is fixed. If capital is fixed, then more output is produced by using additional labour.

Long run: a period of time that is sufficient to enable all factors of production to be adjusted.

Very long run: a period sufficiently long for new technology to develop.

Total product is the relationship between total output produced and the number of workers employed, for a given amount of capital.

Marginal product of labour is the addition to output produced by each additional worker. It is also the slope of the total product curve.

Law of diminishing returns: when increments of a variable factor (labour) are added to a fixed amount of another factor (capital), the marginal product of the variable factor must eventually decline.

Average product of labour is the number of units of output produced per unit of labour at different levels of employment.

Fixed costs are costs that are independent of the level of output.

Variable costs are related to the output produced.

Total cost is the sum of fixed cost and variable cost.

Average fixed cost is the total fixed cost per unit of output.

Average variable cost is the total variable cost per unit of output.

Average total cost is the sum of all costs per unit of output.

Marginal cost of production is the cost of producing each additional unit of output.

Sunk cost is a fixed cost that has already been incurred and cannot be recovered, even by producing a zero output.

Increasing returns to scale implies that, when all inputs are increased by a given proportion, output increases more than proportionately.

Constant returns to scale implies that output increases in direct proportion to an equal proportionate increase in all inputs.

Decreasing returns to scale implies that an equal proportionate increase in all inputs leads to a less than proportionate increase in output.

Long-run average total cost is the lower envelope of all the short-run *ATC* curves.

Minimum efficient scale defines a threshold size of operation such that scale economies are almost exhausted.

Long-run marginal cost is the increment in cost associated with producing one more unit of output when all inputs are adjusted in a cost minimizing manner.

Technological change represents innovation that can reduce the cost of production or bring new products on line.

Globalization is the tendency for international markets to be ever more integrated.

Cluster: a group of firms producing similar products, or engaged in similar research.

Learning by doing can reduce costs. A longer history of production enables firms to accumulate knowledge and thereby implement more efficient production processes.

Economies of scope occur if the unit cost of producing particular products is less when combined with the production of other products than when produced alone.

A platform is a hardware-cum-software capital installation that has multiple production capabilities

Exercises for Chapter 8

EXERCISE 8.1

The relationship between output Q and the single variable input L is given by the form $Q = 5\sqrt{L}$. Capital is fixed. This relationship is given in the table below for a range of L values.

L	1	2	3	4	5	6	7	8	9	10	11	12
Q	5	7.07	8.66	10	11.18	12.25	13.23	14.14	15	15.81	16.58	17.32

1. Add a row to this table and compute the *MP*.
2. Draw the total product (*TP*) curve to scale, either on graph paper or in a spreadsheet.
3. Inspect your graph to see if it displays diminishing *MP*.

EXERCISE 8.2

The *TP* for different output levels for Primitive Products is given in the table below.

Q	1	6	12	20	30	42	53	60	66	70
L	1	2	3	4	5	6	7	8	9	10

1. Graph the TP curve to scale.
2. Add a row to the table and enter the values of the *MP* of labour. Graph this in a separate diagram.
3. Add a further row and compute the *AP* of labour. Add it to the graph containing the *MP* of labour.
4. By inspecting the *AP* and *MP* graph, can you tell if you have drawn the curves correctly? How?

EXERCISE 8.3

A short-run relationship between output and total cost is given in the table below.

Output	0	1	2	3	4	5	6	7	8	9
Total Cost	12	27	40	51	61	70	80	91	104	120

1. What is the total fixed cost of production in this example?
2. Add four rows to the table and compute the *TVC*, *AFC*, *AVC* and *ATC* values for each level of output.
3. Add one more row and compute the *MC* of producing additional output levels.
4. Graph the *MC* and *AC* curves using the information you have developed.

EXERCISE 8.4

Consider the long-run total cost structure for the two firms A and B below.

Output	1	2	3	4	5	6	7
Total cost A	40	52	65	80	97	119	144
Total cost B	30	40	50	60	70	80	90

1. Compute the long-run *ATC* curve for each firm.
2. Plot these curves and examine the type of scale economies each firm experiences at different output levels.

EXERCISE 8.5

Use the data in Exercise 8.4,

1. Calculate the long-run *MC* at each level of output for the two firms.
2. Verify in a graph that these *LMC* values are consistent with the *LAC* values.

EXERCISE 8.6

Optional: Suppose you are told that a firm of interest has a long-run average total cost that is defined by the relationship $LATC=4+48/q$.

1. In a table, compute the *LATC* for output values ranging from 1...24. Plot the resulting *LATC* curve.
2. What kind of returns to scale does this firm never experience?
3. By examining your graph, what will be the numerical value of the *LATC* as output becomes very large?
4. Can you guess what the form of the long-run *MC* curve is?

1. Note that the long-run average total cost is not the collection of minimum points from each short-run average cost curve. The envelope of the short-run curves will pick up mainly points that are not at the minimum, as you will see if you try to draw the outcome. The intuition behind the definition is this: With increasing returns to scale, it may be better to build a plant size that operates with some spare capacity than to build one that is geared to producing a smaller output level. In building the larger plant, we can take greater advantage of the scale economies, and it may prove less costly to produce in such a plant than to produce with a smaller plant that has less unused capacity and does not exploit the underlying scale economies. Conversely, in the presence of decreasing returns to scale, it may be less costly to produce output in a plant that is used "overtime" than to use a larger plant that suffers from scale diseconomies.

2. Two English farmers changed the lives of millions in the 18th century through their technological genius. The first was Charles Townsend, who introduced the concept of crop rotation. He realized that continual use of soil for a single purpose drained the land of key nutrients. This led him ultimately to propose a five-year rotation that included tillage, vegetables, and sheep. This rotation reduced the time during which land would lie fallow, and therefore increased the productivity of land. A second game-changing technological innovation saw the introduction of the seed plough, a tool that facilitated seed sowing in rows, rather than scattering it randomly. This line sowing meant that weeds could be controlled more easily—weeds that would otherwise smother much of the crop. The genius here was a gentleman by the name of Jethro Tull.

This page titled [8: Production and cost](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.1: Efficient production

Firms that fail to operate efficiently seldom survive. They are dominated by their competitors because the latter produce more efficiently and can sell at a lower price. The drive for profitability is everywhere present in the modern economy. Companies that promise more profit, by being more efficient, are valued more highly on the stock exchange. For example: In July of 2015 *Google* announced that, going forward, it would be more attentive to cost management in its numerous research endeavours that aim to bring new products to the marketplace. This policy, put in place by the Company's new Chief Financial Officer, was welcomed by investors who, as a result, bought up the stock. The Company's stock increased in value by 16% in one day – equivalent to about \$50 billion.

The remuneration of managers in virtually all corporations is linked to profitability. Efficient production, *a.k.a.* cost reduction, is critical to achieving this goal. In this chapter we will examine cost management and efficient production from the ground up – by exploring how a small entrepreneur brings his or her product to market in the most efficient way possible. As we shall see, efficient production and cost minimization amount to the same thing: Cost minimization is the financial reflection of efficient production.

Efficient production is critical in any budget-driven organization, not just in the private sector. Public institutions equally are, and should be, concerned with costs and efficiency.

Entrepreneurs employ factors of production (capital and labour) in order to transform raw materials and other inputs into goods or services. The relationship between output and the inputs used in the production process is called a production function. It specifies how much output can be produced with given combinations of inputs. A production function is not restricted to profit-driven organizations. Municipal road repairs are carried out with labour and capital. Students are educated with teachers, classrooms, computers, and books. Each of these is a production process.

Production function: a technological relationship that specifies how much output can be produced with specific amounts of inputs.

Economists distinguish between two concepts of efficiency: One is technological efficiency; the other is economic efficiency. To illustrate the difference, consider the case of auto assembly: the assembler could produce its vehicles either by using a large number of assembly workers and a plant that has a relatively small amount of machinery, or it could use fewer workers accompanied by more machinery in the form of robots. Each of these processes could be deemed technologically efficient, provided that there is no waste. If the workers without robots are combined with their capital to produce as much as possible, then that production process is technologically efficient. Likewise, in the scenario with robots, if the workers and capital are producing as much as possible, then that process too is efficient in the technological sense.

Technological efficiency means that the maximum output is produced with the given set of inputs.

Economic efficiency is concerned with more than just technological efficiency. Since the entrepreneur's goal is to make profit, she must consider which technologically efficient process best achieves that objective. More broadly, any budget-driven process should focus on being economically efficient, whether in the public or private sector. An economically efficient production structure is the one that produces output at least cost.

Economic efficiency defines a production structure that produces output at least cost.

Auto-assembly plants the world over have moved to using robots during the last two decades. Why? The reason is not that robots were invented 20 years ago; they were invented long before that. The real reason is that, until recently, this technology was not economically efficient. Robots were too expensive; they were not capable of high-precision assembly. But once their cost declined and their accuracy increased they became economically efficient. The development of robots represented technological progress. When this progress reached a critical point, entrepreneurs embraced it.

To illustrate the point further, consider the case of garment assembly. There is no doubt that engineers could make robots capable of joining the pieces of fabric that form garments. This is not beyond our technological abilities. Why, then, do we not have such capital-intensive production processes for garment making, similar to the production process chosen by vehicle producers? The answer is that, while such a concept could be technologically efficient, it would not be economically efficient. It is more profitable to use large amounts of labour and relatively traditional machines to assemble garments, particularly when labour in Asia costs less and the garments can be shipped back to Canada inexpensively. Containerization and scale economies in shipping mean that a garment can be shipped to Canada from Asia for a few cents per unit.

Efficiency in production is not limited to the manufacturing sector. Farmers must choose the optimal combination of labour, capital and fertilizer to use. In the health and education sectors, efficient supply involves choices on how many high- and low-skill workers

to employ, how much traditional physical capital to use, how much information technology to use, based upon the productivity and cost of each. Professors and physicians are costly inputs. When they work with new technology (capital) they become more efficient at performing their tasks: It is less costly to have a single professor teach in a 300-seat classroom that is equipped with the latest technology, than have several professors each teaching 60-seat classes with chalk and a blackboard.

This page titled [8.1: Efficient production](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.2: The time frame

We distinguish initially between the short run and the long run. When discussing technological change, we use the term very long run. These concepts have little to do with clocks or calendars; rather, they are defined by the degree of flexibility an entrepreneur or manager has in her production process. A key decision variable is capital.

A customary assumption is that a producer can hire more labour immediately, if necessary, either by taking on new workers (since there are usually some who are unemployed and looking for work), or by getting the existing workers to work longer hours. In contrast, getting new capital in place is usually more time consuming: The entrepreneur may have to place an order for new machinery, which will involve a production and delivery time lag. Or she may have to move to a more spacious location in order to accommodate the added capital. Whether this calendar time is one week, one month, or one year is of no concern to us. We define the long run as a period of sufficient length to enable the entrepreneur to adjust her capital stock, whereas in the short run at least one factor of production is fixed. Note that it matters little whether it is labour or capital that is fixed in the short run. A software development company may be able to install new capital (computing power) instantaneously but have to train new developers. In such a case capital is variable and labour is fixed in the short run. The definition of the short run is that one of the factors is fixed, and in our examples we will assume that it is capital.

Short run: a period during which at least one factor of production is fixed. If capital is fixed, then more output is produced by using additional labour.

Long run: a period of time that is sufficient to enable all factors of production to be adjusted.

Very long run: a period sufficiently long for new technology to develop.

This page titled [8.2: The time frame](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis](#) and [Ian Irvine](#) ([Lyryx](#)) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.3: Production in the short run

Black Diamond Snowboards (BDS) is a start-up snowboard producing enterprise. Its founder has invented a new lamination process that gives extra strength to his boards. He has set up a production line in his garage that has four workstations: Laminating, attaching the steel edge, waxing, and packing.

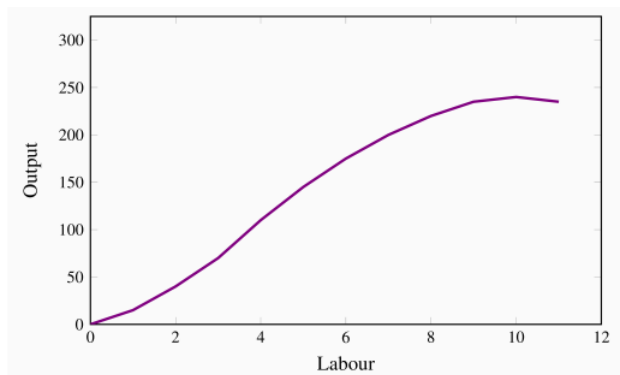
With this process in place, he must examine how productive his firm can be. After extensive testing, he has determined exactly how his productivity depends upon the number of workers. If he employs only one worker, then that worker must perform several tasks, and will encounter 'down time' between workstations. Extra workers would therefore not only increase the total output; they could, in addition, increase output *per worker*. He also realizes that once he has employed a critical number of workers, additional workers may not be so productive: Because they will have to share the fixed amount of machinery in his garage, they may have to wait for another worker to finish using a machine. At such a point, the productivity of his plant will begin to fall off, and he may want to consider capital expansion. But for the moment he is constrained to using this particular assembly plant. Testing leads him to formulate the relationship between workers and output that is described in Table 8.1.

Table 8.1 Snowboard production and productivity

1	2	3	4	5
Workers	Output	Marginal	Average	Stages of
	(<i>TP</i>)	product	product	production
		(MP_L)	(AP_L)	
0	0			<i>MP_L increasing</i>
1	15	15	15	
2	40	25	20	
3	70	30	23.3	
4	110	40	27.5	<i>MP_L positive and declining</i>
5	145	35	29	
6	175	30	29.2	
7	200	25	28.6	
8	220	20	27.5	
9	235	15	26.1	
10	240	5	24.0	<i>MP_L negative</i>
11	235	-5	21.4	

By increasing the number of workers in the plant, BDS produces more boards. The relationship between these two variables in columns 1 and 2 in the table is plotted in Figure 8.1. This is called the total product function (*TP*), and it defines the output produced with different amounts of labour in a plant of fixed size.

Figure 8.1 Total product curve



Output increases with the amount of labour used. Initially the increase in output due to using more labour is high, subsequently it is lower. The initial phase characterizes increasing productivity, the later phase defines declining productivity.

Total product is the relationship between total output produced and the number of workers employed, for a given amount of capital.

This relationship is positive, indicating that more workers produce more boards. But the curve has an interesting pattern. In the initial expansion of employment it becomes progressively steeper – its curvature is slightly convex; following this phase the function's increase becomes progressively less steep – its curvature is concave. These different stages in the *TP* curve tell us a great deal about productivity in BDS. To see this, consider the additional number of boards produced by each worker. The first worker produces 15. When a second worker is hired, the total product rises to 40, so the additional product attributable to the second worker is 25. A third worker increases output by 30 units, and so on. We refer to this additional output as the marginal product (*MP*) of an additional worker, because it defines the incremental, or marginal, contribution of the worker. These values are entered in column 3.

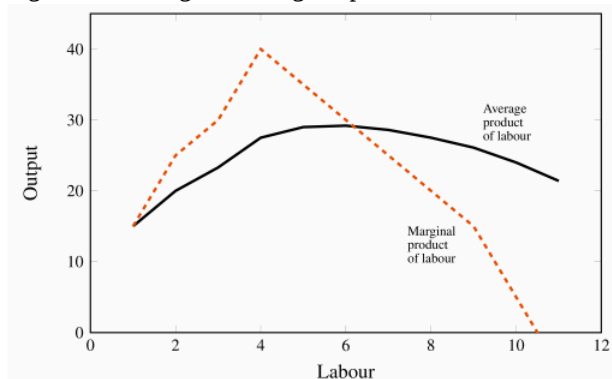
More generally the *MP* of labour is defined as the change in output divided by the change in the number of units of labour employed. Using, as before, the Greek capital delta (Δ) to denote a change, we can define

$$MP_L = \frac{\text{Change in output produced}}{\text{Change in labour employed}} = \frac{\Delta Q}{\Delta L}$$

In this example the change in labour is one unit at each stage and hence the marginal product of labour is simply the corresponding change in output. It is also the case that the MP_L is the slope of the *TP* curve – the change in the value on the vertical axis due to a change in the value of the variable on the horizontal axis.

Marginal product of labour is the addition to output produced by each additional worker. It is also the slope of the total product curve.

Figure 8.2 Average and marginal product curves



The productivity curves initially rise and then decline, reflecting increasing and decreasing productivity. The MP_L curves must intersect the AP_L curve at the maximum of the AP_L : The average must increase if the marginal exceeds the average and must decline if the marginal is less than the average.

During the initial stage of production expansion, the marginal product of each worker is increasing. It increases from 15 to 40 as BDS moves from having one employee to four employees. This increasing *MP* is made possible by the fact that each worker is able

to spend more time at his workstation, and less time moving between tasks. But, at a certain point in the employment expansion, the MP reaches a maximum and then begins to tail off. At this stage – in the concave region of the TP curve – additional workers continue to produce additional output, but at a diminishing rate. For example, while the fourth worker adds 40 units to output, the fifth worker adds 35, the sixth worker 30, and so on. This declining MP is due to the constraint of a fixed number of machines: All workers must share the same capital. The MP function is plotted in Figure 8.2.

The phenomenon we have just described has the status of a law in economics: The law of diminishing returns states that, in the face of a fixed amount of capital, the contribution of additional units of a variable factor must eventually decline.

Law of diminishing returns: when increments of a variable factor (labour) are added to a fixed amount of another factor (capital), the marginal product of the variable factor must eventually decline.

The relationship between Figures 8.1 and 8.2 should be noted. First, the MP_L reaches a maximum at an output of 4 units – where the slope of the TP curve is greatest. The MP_L curve remains positive beyond this output, but declines: The TP curve reaches a maximum when the tenth unit of labour is employed. An eleventh unit actually reduces total output; therefore, the MP of this eleventh worker is negative! In Figure 8.2, the MP curve becomes negative at this point. The garage is now so crowded with workers that they are beginning to obstruct the operation of the production process. Thus the producer would never employ an eleventh unit of labour.

Next, consider the information in the fourth column of the table. It defines the average product of labour (AP_L)—the amount of output produced, on average, by workers at different employment levels:

$$AP_L = \frac{\text{Total output produced}}{\text{Total amount of labour employed}} = \frac{Q}{L}.$$

This function is also plotted in Figure 8.2. Referring to the table: The AP column indicates, for example, that when two units of labour are employed and forty units of output are produced, the average production level of each worker is 20 units ($=40/2$). When three workers produce 70 units, their average production is 23.3 ($=70/3$), and so forth. Like the MP function, this one also increases and subsequently decreases, reflecting exactly the same productivity forces that are at work on the MP curve.

Average product of labour is the number of units of output produced per unit of labour at different levels of employment.

The AP and MP functions intersect at the point where the AP is at its peak. This is no accident, and has a simple explanation. Imagine a softball player who is batting .280 coming into today's game—she has been hitting her way onto base 28 percent of the time when batting, so far this season. This is her average product, AP .

In today's game, if she bats .500 (hits her way to base on half of her at-bats), then she will improve her average. Today's batting (MP) at .500 therefore pulls up the season's AP . Accordingly, whenever the MP exceeds the AP , the AP is pulled up. By the same reasoning, if her MP is less than the season average, her average will be pulled down. It follows that the two functions must intersect at the peak of the AP curve. To summarize:

*If the MP exceeds the AP , then the AP increases;
If the MP is less than the AP , then the AP declines.*

While the owner of BDS may understand his productivity relations, his ultimate goal is to make profit, and for this he must figure out how productivity translates into cost.

This page titled [8.3: Production in the short run](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.4: Costs in the short run

The cost structure for the production of snowboards at Black Diamond is illustrated in Table 8.2. Employees are skilled and are paid a weekly wage of \$1,000. The cost of capital is \$3,000 and it is fixed, which means that it does not vary with output. As in Table 8.1, the number of employees and the output are given in the first two columns. The following three columns define the capital costs, the labour costs, and the sum of these in producing different levels of output. We use the terms fixed, variable, and total costs to define the cost structure of a firm. Fixed costs do not vary with output, whereas variable costs do, and total costs are the sum of fixed and variable costs. To keep this example as simple as possible, we will ignore the cost of raw materials. We could add an additional column of costs, but doing so will not change the conclusions.

Table 8.2 Snowboard production costs

Workers	Output	Capital	Labour	Total	Average	Average	Average	Marginal
		cost	cost	costs	fixed	variable	total	cost
		fixed	variable		cost	cost	cost	
0	0	3,000	0	3,000				
1	15	3,000	1,000	4,000	200.0	66.7	266.7	66.7
2	40	3,000	2,000	5,000	75.0	50.0	125.0	40.0
3	70	3,000	3,000	6,000	42.9	42.9	85.7	33.3
4	110	3,000	4,000	7,000	27.3	36.4	63.6	25.0
5	145	3,000	5,000	8,000	20.7	34.5	55.2	28.6
6	175	3,000	6,000	9,000	17.1	34.3	51.4	33.3
7	200	3,000	7,000	10,000	15.0	35.0	50.0	40.0
8	220	3,000	8,000	11,000	13.6	36.4	50.0	50.0
9	235	3,000	9,000	12,000	12.8	38.3	51.1	66.7
10	240	3,000	10,000	13,000	12.5	41.7	54.2	200.0

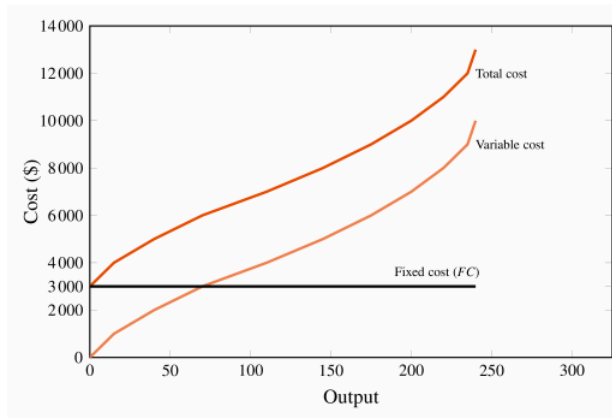
Fixed costs are costs that are independent of the level of output.

Variable costs are related to the output produced.

Total cost is the sum of fixed cost and variable cost.

Total costs are illustrated in Figure 8.3 as the vertical sum of variable and fixed costs. For example, Table 8.2 indicates that the total cost of producing 220 units of output is the sum of \$3,000 in fixed costs plus \$8,000 in variable costs. Therefore, at the output level 220 on the horizontal axis in Figure 8.3, the sum of the cost components yields a value of \$11,000 that forms one point on the total cost curve. Performing a similar calculation for every possible output yields a series of points that together form the complete total cost curve.

Figure 8.3 Total cost curves



Total cost is the vertical sum of the variable and fixed costs.

Average costs are given in the next three columns of Table 8.2. Average cost is the cost per unit of output, and we can define an average cost corresponding to each of the fixed, variable, and total costs defined above. Average fixed cost (*AFC*) is the total fixed cost divided by output; average variable cost (*AVC*) is the total variable cost divided by output; and average total cost (*ATC*) is the total cost divided by output.

<i>AFC</i>	$= (\text{Fixed cost})/Q = FC/Q$
<i>AVC</i>	$= (\text{Total variable costs})/Q = TVC/Q$
<i>ATC</i>	$= AFC + AVC$

Average fixed cost is the total fixed cost per unit of output.

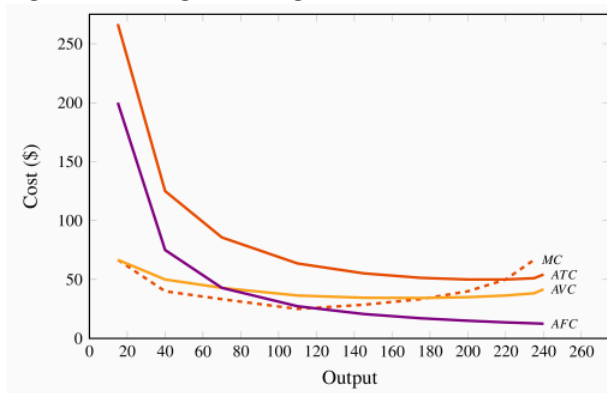
Average variable cost is the total variable cost per unit of output.

Average total cost is the sum of all costs per unit of output.

The productivity-cost relationship

Consider the **average variable cost - average product relationship**, as developed in column 7 of Table 8.2; its corresponding variable cost curve is plotted in Figure 8.4. In this example, *AVC* first decreases and then increases. The intuition behind its shape is straightforward (and realistic) if you have understood why productivity varies in the short run: The variable cost, which represents the cost of labour, is constant per unit of labour, because the wage paid to each worker does not change. However, each worker's productivity varies. Initially, when we hire more workers, they become more productive, perhaps because they have less 'down time' in switching between tasks. This means that the labour costs per snowboard must decline. At some point, however, the law of diminishing returns sets in: As before, each additional worker is paid a constant amount, but as productivity declines the labour cost per snowboard increases.

Figure 8.4 Average and marginal cost curves



The *MC* intersects the *ATC* and *AVC* at their minimum values. The *AFC* declines indefinitely as fixed costs are spread over a greater output.

In this numerical example the AP is at a maximum when six units of labour are employed and output is 175. This is also the point where the AVC is at a minimum. This maximum/minimum relationship is also illustrated in Figures 8.2 and 8.4.

Now consider the **marginal cost - marginal product relationship**. The marginal cost (MC) defines the cost of producing one more unit of output. In Table 8.2, the marginal cost of output is given in the final column. It is the additional cost of production divided by the additional number of units produced. For example, in going from 15 units of output to 40, total costs increase from \$4,000 to \$5,000. The MC is the cost of those additional units divided by the number of additional units. In this range of output, MC is $\$1,000/25 = \40 . We could also calculate the MC as the addition to variable costs rather than the addition to total costs, because the *addition* to each is the same—fixed costs are fixed. Hence:

MC	$= \frac{\text{Change in total costs}}{\text{Change in output produced}} = \frac{\Delta TC}{\Delta Q}$
	$= \frac{\text{Change in variable costs}}{\text{Change in output produced}} = \frac{\Delta TVC}{\Delta Q}$

Marginal cost of production is the cost of producing each additional unit of output.

Just as the behaviour of the AVC curve is determined by the AP curve, so too the behaviour of the MC is determined by the MP curve. When the MP of an additional worker exceeds the MP of the previous worker, this implies that the cost of the additional output produced by the last worker hired must be declining. To summarize:

- If the marginal product of labour increases, then the marginal cost of output declines;*
- If the marginal product of labour declines, then the marginal cost of output increases.*

In our example, the MP_L reaches a maximum when the fourth unit of labour is employed (or 110 units of output are produced), and this also is where the MC is at a minimum. This illustrates that the *marginal cost reaches a minimum at the output level where the marginal product reaches a maximum*.

The average total cost is the sum of the fixed cost per unit of output and the variable cost per unit of output. Typically, fixed costs are the dominant component of total costs at low output levels, but become less dominant at higher output levels. Unlike average variable costs, note that the average fixed cost must always decline with output, because a fixed cost is being spread over more units of output. Hence, when the ATC curve eventually increases, it is because the increasing variable cost component eventually dominates the declining AFC component. In our example, this occurs when output increases from 220 units (8 workers) to 235 (9 workers).

Finally, observe the interrelationship between the MC curve on the one hand and the ATC and AVC on the other. Note from Figure 8.4 that the MC cuts the AVC and the ATC at the minimum point of each of the latter. The logic behind this pattern is analogous to the logic of the relationship between marginal and average product curves: When the cost of an additional unit of output is less than the average, this reduces the average cost; whereas, if the cost of an additional unit of output is above the average, this raises the average cost. This must hold true regardless of whether we relate the MC to the ATC or the AVC .

- When the marginal cost is less than the average cost, the average cost must decline;*
- When the marginal cost exceeds the average cost, the average cost must increase.*

Notation: We use both the abbreviations ATC and AC to denote average total cost. The term 'average cost' is understood in economics to include both fixed and variable costs.

Teams and services

The choice faced by the producer in the example above is slightly 'stylized', yet it still provides an appropriate rule for analyzing hiring decisions. In practice, it is quite difficult to isolate or identify the marginal product of an individual worker. One reason is that individuals work in teams within organizations. The accounting department, the marketing department, the sales department, the assembly unit, the chief executive's unit are all composed of teams. Adding one more person to human resources may have no impact on the number of units of output produced by the company in a measurable way, but it may influence worker morale and hence longer-term productivity. Nonetheless, if we consider expanding, or contracting, any one department within an organization, management can attempt to estimate the net impact of additional hires (or layoffs) on the contribution of each team to the firm's profitability. Adding a person in marketing may increase sales, laying off a person in research and development may reduce costs

by more than it reduces future value to the firm. In practice this is what firms do: they attempt to assess the contribution of each team in their organization to costs and revenues, and on that basis determine the appropriate number of employees.

The manufacturing sector of the macro economy is dominated, sizewise, by the services sector. But the logic that drives hiring decisions, as developed above, applies equally to services. For example, how does a law firm determine the optimal number of paralegals to employ per lawyer? How many nurses are required to support a surgeon? How many university professors are required to teach a given number of students?

All of these employment decisions involve optimization at the margin. The goal of the decision maker is not always profit, but she should attempt to estimate the cost and value of adding personnel at the margin.

This page titled [8.4: Costs in the short run](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.5: Fixed costs and sunk costs

The distinction between fixed and variable costs is important for producers who are not making a profit. If a producer has committed himself to setting up a plant, then he has made a decision to incur a fixed cost. Having done this, he must now decide on a production strategy that will maximize profit. However, the price that consumers are willing to pay may not be sufficient to yield a profit. So, if Black Diamond Snowboards cannot make a profit, should it shut down? The answer is that if it can cover its variable costs, *having already incurred its fixed costs*, it should stay in production, at least temporarily. By covering the variable cost of its operation, Black Diamond is at least earning some return. A sunk cost is a fixed cost that has already been incurred and cannot be recovered. But if the pressures of the marketplace are so great that the total costs cannot be covered in the longer run, then this is not a profitable business and the firm should close its doors.

Is a fixed cost always a sunk cost? No: Any production that involves capital will incur a fixed cost component. Such capital can be financed in several ways however: It might be financed on a very short-term lease basis, or it might have been purchased by the entrepreneur. If it is leased on a month-to-month basis, an unprofitable entrepreneur who can only cover variable costs (and who does not foresee better market conditions ahead) can exit the industry quickly – by not renewing the lease on the capital. But an individual who has actually purchased equipment that cannot readily be resold has essentially sunk money into the fixed cost component of his production. This entrepreneur should continue to produce as long as he can cover variable costs.

Sunk cost is a fixed cost that has already been incurred and cannot be recovered, even by producing a zero output.

R & D as a sunk cost

Sunk costs in the modern era are frequently in the form of research and development costs, not the cost of building a plant or purchasing machinery. The prototypical example is the pharmaceutical industry, where it is becoming progressively more challenging to make new drug breakthroughs – both because the 'easier' breakthroughs have already been made, and because it is necessary to meet tighter safety conditions attaching to new drugs. Research frequently leads to drugs that are not sufficiently effective in meeting their target. As a consequence, the pharmaceutical sector regularly writes off hundreds of millions of dollars of lost sunk costs – unfruitful research and development.

Finally, we need to keep in mind the opportunity costs of running the business. The owner pays himself a salary, and ultimately he must recognize that the survival of the business should not depend upon his drawing a salary that is less than his opportunity cost. As developed in Section 7.2, if he underpays himself in order to avoid shutting down, he might be better off in the long run to close the business and earn his opportunity cost elsewhere in the marketplace.

A dynamic setting

We need to ask why it might be possible to cover all costs in a longer run horizon, while in the near-term costs are not covered. The principal reason is that demand may grow, particularly for a new product. For example, in 2019 numerous cannabis producing firms were listed on the Canadian Securities Exchange, and collectively were valued at about fifty billion dollars. None had revenues that covered costs, yet investors poured money into this sector. Investors evidently envisaged that the market for legal cannabis would grow. As of 2020 it appears that these investors were excessively optimistic. Sales growth has been slow and stock valuations have plummeted.

This page titled [8.5: Fixed costs and sunk costs](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine](#) ([Lyryx](#)) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.6: Long-run production and costs

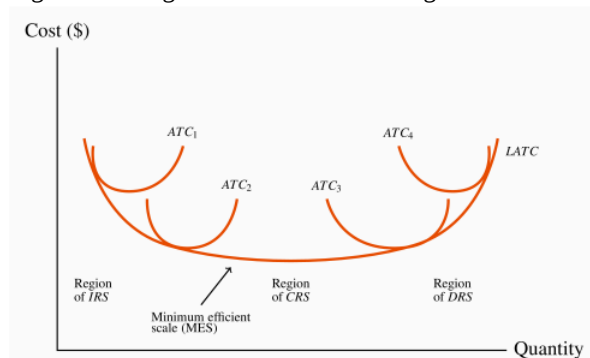
The snowboard manufacturer we portray produces a relatively low level of output; in reality, millions of snowboards are produced each year in the global market. Black Diamond Snowboards may have hoped to get a start by going after a local market—the "free-ride" teenagers at Mont Sainte Anne in Quebec or at Fernie in British Columbia. If this business takes off, the owner must increase production, take the business out of his garage and set up a larger-scale operation. But how will this affect his cost structure? Will he be able to produce boards at a lower cost than when he was producing a very limited number of boards each season? Real-world experience would indicate yes.

Production costs almost always decline when the *scale* of the operation initially increases. We refer to this phenomenon simply as economies of scale. There are several reasons why scale economies are encountered. One is that production flows can be organized in a more efficient manner when more is being produced. Another is that the opportunity to make greater use of task specialization presents itself; for example, Black Diamond Snowboards may be able to subdivide tasks within the laminating and packaging stations. With a larger operating scale the replacement of labor with capital may be economically efficient. If scale economies do define the real world, then a bigger plant—one that is geared to produce a higher level of output—should have an average total cost curve that is "lower" than the cost curve corresponding to the smaller scale of operation we considered in the example above.

Average costs in the long run

Figure 8.5 illustrates a possible relationship between the ATC curves for four different scales of operation. ATC_1 is the average total cost curve associated with a small-sized plant; think of it as the plant built in the entrepreneur's garage. ATC_2 is associated with a somewhat larger plant, perhaps one she has put together in a rented industrial or commercial space. The further a cost curve is located to the right of the diagram the larger the production facility it defines, given that output is measured on the horizontal axis. If there are economies associated with a larger scale of operation, then the average costs associated with producing larger outputs in a larger plant should be lower than the average costs associated with lower outputs in a smaller plant, assuming that the plants are producing the output levels they were designed to produce. For this reason, the cost curve ATC_2 and the cost curve ATC_3 each have a segment that is lower than the lowest segment on ATC_1 . However, in Figure 8.5 the cost curve ATC_4 has moved upwards. What behaviours are implied here?

Figure 8.5 Long-run and short-run average costs



The long-run ATC curve, $LATC$, is the lower envelope of all short-run ATC curves. It defines the least cost per unit of output when all inputs are variable. Minimum efficient scale is that output level at which the $LATC$ is a minimum, indicating that further increases in the scale of production will not reduce unit costs.

In many production environments, beyond some large scale of operation, it becomes increasingly difficult to reap further cost reductions from specialization, organizational economies, or marketing economies. At such a point, the scale economies are effectively exhausted, and larger plant sizes no longer give rise to lower (short-run) ATC curves. This is reflected in the similarity of the ATC_2 and the ATC_3 curves. The pattern suggests that we have almost exhausted the possibilities of further scale advantages once we build a plant size corresponding to ATC_2 . Consider next what is implied by the position of the ATC_4 curve relative to the ATC_2 and ATC_3 curves. The relatively higher position of the ATC_4 curve implies that unit costs will be higher in a yet larger plant. Stated differently: If we increase the scale of this firm to extremely high output levels, we are actually encountering diseconomies of scale. Diseconomies of scale imply that unit costs increase as a result of the firm's becoming too large: Perhaps co-ordination difficulties have set in at the very high output levels, or quality-control monitoring costs have risen. These coordination and management difficulties are reflected in increasing unit costs in the long run.

The terms increasing, constant, and decreasing returns to scale underlie the concepts of scale economies and diseconomies: Increasing returns to scale (IRS) implies that, when all inputs are increased by a given proportion, output increases more than proportionately. Constant returns to scale (CRS) implies that output increases in direct proportion to an equal proportionate increase in all inputs. Decreasing returns to scale (DRS) implies that an equal proportionate increase in all inputs leads to a less than proportionate increase in output.

Increasing returns to scale implies that, when all inputs are increased by a given proportion, output increases more than proportionately.

Constant returns to scale implies that output increases in direct proportion to an equal proportionate increase in all inputs.

Decreasing returns to scale implies that an equal proportionate increase in all inputs leads to a less than proportionate increase in output.

These are pure production function relationships, but, if the prices of inputs are fixed for producers, they translate directly into the various cost structures illustrated in Figure 8.5. For example, if a 40% increase in capital and labour use allows for better production flows than when in the smaller plant, and therefore yields more than a 40% increase in output, this implies that the cost per snowboard produced must fall in the new plant. In contrast, if a 40% increase in capital and labour leads to say just a 30% increase in output, then the cost per snowboard in the new larger plant must be higher. Between these extremes, there may be a range of relatively constant unit costs, corresponding to where the production relation is subject to constant returns to scale. In Figure 8.5, the falling unit costs output region has increasing returns to scale, the region that has relatively constant unit costs has constant returns to scale, and the increasing cost region has decreasing returns to scale.

Increasing returns to scale characterize businesses with large initial costs and relatively low costs of producing each unit of output. Computer chip manufacturers, pharmaceutical manufacturers, vehicle rental agencies, booking agencies such as *booking.com* or *hotels.com*, intermediaries such as *airbnb.com*, even brewers, all benefit from scale economies. In the beer market, brewing, bottling and shipping are all low-cost operations relative to the capital cost of setting up a brewery. Consequently, we observe surprisingly few breweries in any brewing company, even in large land-mass economies such as Canada or the US.

In addition to the four short-run average total cost curves, Figure 8.5 contains a curve that forms an envelope around the bottom of these short-run average cost curves. This envelope is the long-run average total cost (*LATC*) curve, because it defines average cost as we move from one plant size to another. Remember that in the long run both labour and capital are variable, and as we move from one short-run average cost curve to another, that is exactly what happens—all factors of production are variable. Hence, the collection of short-run cost curves in Figure 8.5 provides the ingredients for a long-run average total cost curve¹.

$$LATC = (\text{Long-run total costs})/Q = LTC/Q$$

Long-run average total cost is the lower envelope of all the short-run *ATC* curves.

The particular range of output on the *LATC* where it begins to flatten out is called the range of minimum efficient scale. This is an important concept in industrial policy, as we shall see in later chapters. At such an output level, the producer has expanded sufficiently to take advantage of virtually all the scale economies available.

Minimum efficient scale defines a threshold size of operation such that scale economies are almost exhausted.

In view of this discussion and the shape of the *LATC* in Figure 8.5, it is obvious that economies of scale can also be defined in terms of the curvature of the *LATC*. Where the *LATC* declines there are IRS, where the *LATC* is flat there are CRS, where the *LATC* slopes upward there are DRS.

Table 8.3 *LATC* elements for two plants (thousands \$)

<i>Q</i>	<i>AFC</i> ₁	<i>MC</i> ₁ = <i>AVC</i> ₁	<i>ATC</i> ₁	<i>AFC</i> ₂	<i>MC</i> ₂ = <i>AVC</i> ₂	<i>ATC</i> ₂
20	50	30	80	100	25	125
40	25	30	55	50	25	75
60	16.67	30	46.67	33.33	25	58.33
80	12.5	30	42.5	25	25	50
100	10	30	40	20	25	45

120	8.33	30	38.33	16.67	25	41.67
140	7.14	30	37.14	14.29	25	39.29
160	6.25	30	36.25	12.5	25	37.5
180	5.56	30	35.56	11.11	25	36.11
200	5	30	35	10	25	35
220	4.55	30	34.55	9.09	25	34.09
240	4.17	30	34.17	8.33	25	33.33
260	3.85	30	33.85	7.69	25	32.69
280	3.57	30	33.57	7.14	25	32.14

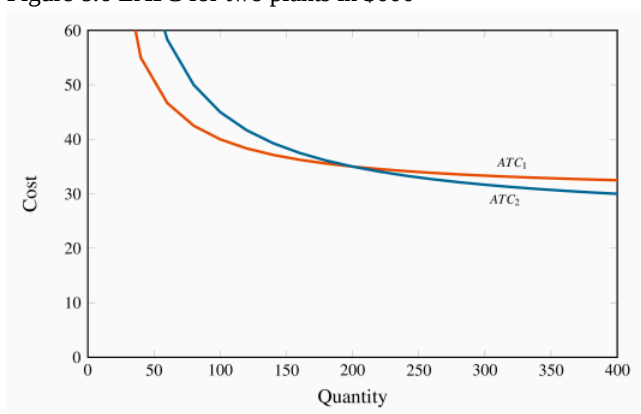
Plant 1 $FC = \$1\text{m}$. Plant 2 $FC = \$2\text{m}$. For $Q < 200$, $ATC_1 < ATC_2$; for $Q > 200$, $ATC_1 > ATC_2$; and for $Q = 200$, $ATC_1 = ATC_2$. $LATC$ defined by data in bold font.

Long-run costs – a simple numerical example

Kitt is an automobile designer specializing in the production of off-road vehicles sold to a small clientele. He has a choice of two (and only two) plant sizes; one involving mainly labour and the other employing robots extensively. The set-up (i.e. fixed) costs of these two assembly plants are \$1 million and \$2 million respectively. The advantage to having the more costly plant is that the pure production costs (variable costs) are less. The cost components are defined in Table 8.3. The variable cost (equal to the marginal cost here) is \$30,000 in the plant that relies primarily on labour, and \$25,000 in the plant that has robots. The ATC for each plant size is the sum of AFC and AVC . The AFC declines as the fixed cost is spread over more units produced. The variable cost per unit is constant in each case. By comparing the fourth and final columns, it is clear that the robot-intensive plant has lower costs if it produces a large number of vehicles. At an output of 200 vehicles the average costs in each plant are identical: The higher fixed costs associated with the robots are exactly offset by the lower variable costs at this output level.

The ATC curve corresponding to each plant size is given in Figure 8.6. There are two short-run ATC curves. The positions of these curves indicate that if the manufacturer believes he can produce at least 200 vehicles his unit costs will be less with the plant involving robots; but at output levels less than this his unit costs would be less in the labour-intensive plant.

Figure 8.6 $LATC$ for two plants in \$000



The long-run average cost curve for this producer is the lower envelope of these two cost curves: ATC_1 up to output 200 and ATC_2 thereafter. Two features of this example are to be noted. First we do not encounter decreasing returns – the $LATC$ curve never increases. ATC_1 tends asymptotically to a lower bound of \$30, while ATC_2 tends towards \$25. Second, in the interests of simplicity we have assumed just two plant sizes are possible. With more possibilities on the introduction of robots we could imagine more short-run ATC curves which would form the lower-envelope $LATC$.

This page titled [8.6: Long-run production and costs](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.7: Technological change- globalization and localization

Technological change represents innovation that can reduce the cost of production or bring new products on line. As stated earlier, the very long run is a period that is sufficiently long for new technology to evolve and be implemented.

Technological change represents innovation that can reduce the cost of production or bring new products on line.

Technological change has had an enormous impact on economic life for several centuries. It is not something that is defined in terms of the recent telecommunications revolution. The industrial revolution began in eighteenth century Britain. It was accompanied by a less well-recognized, but equally important, agricultural revolution. The improvement in cultivation technology, and ensuing higher yields, freed up enough labour to populate the factories that were the core of the industrial revolution². The development and spread of mechanical power dominated the nineteenth century, and the mass production line of Henry Ford in autos or Andrew Carnegie in steel heralded in the twentieth century.

Globalization

The modern communications revolution has reduced costs, just like its predecessors. But it has also greatly sped up globalization, the increasing integration of national markets.

Globalization is the tendency for international markets to be ever more integrated.

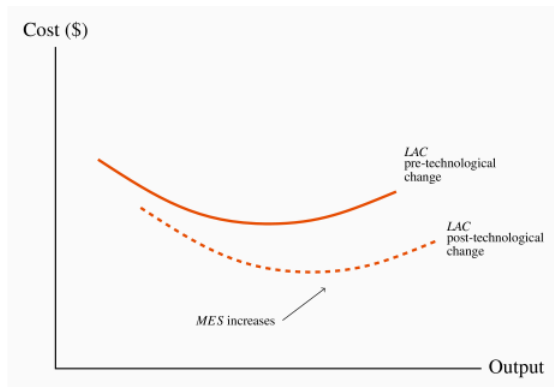
Globalization has several drivers: lower transportation and communication costs; reduced barriers to trade and capital mobility; the spread of new technologies that facilitate cost and quality control; different wage rates between developed and less developed economies. New technology and better communications have been critical in both increasing the minimum efficient scale of operation and reducing diseconomies of scale; they facilitate the efficient management of large companies.

The continued reduction in trade barriers in the post-World War II era has also meant that the effective marketplace has become the globe rather than the national economy for many products. Companies like *Apple*, *Microsoft*, and *Facebook* are visible worldwide. Globalization has been accompanied by the collapse of the Soviet Union, the adoption of an outward looking philosophy on the part of China, and an increasing role for the market place in India. These developments together have facilitated the outsourcing of much of the West's manufacturing to lower-wage economies.

But new technology not only helps existing companies grow large; it also enables new ones to start up. It is now cheaper for small producers to manage their inventories and maintain contact with their own suppliers.

The impact of technology is to reduce the cost of production, hence it will lower the average cost curve in both the short and long run. First, decreasing returns to scale become less probable due to improved communications, so the upward sloping section of the LRAC curve may disappear altogether. Second, capital costs are now lower than in earlier times, because much modern technology has transformed some fixed costs into variable cost: Both software and hardware functions can be subcontracted to specialty firms, who in turn may use cloud computing services, and it is thus no longer necessary to have a substantial in-house computing department. The use of almost-free software such as *Skype*, *Hangouts* and *WhatsApp* reduces communication costs. Advertising on social media is more effective and less costly than in traditional hard-print form. Hiring may be cheaper through *LinkedIn* than through a traditional human resources department. These developments may actually reduce the minimum efficient scale of operation because they reduce the need for large outlays on fixed capital. On the other hand, changes in technology may induce producers to use more capital and less labor, with the passage of time. That would increase the minimum efficient scale. An example of this phenomenon is in mining, or tunnel drilling, where capital investment per worker is greater than when technology was less developed. A further example is the introduction of robotic assistants in *Amazon* warehouses. In these scenarios the minimum efficient scale should increase, and that is illustrated in Figure 8.7.

Figure 8.7 Technological change and LAC



Technological change reduces the unit production cost for any output produced and may also increase the minimum efficient scale (MES) threshold.

The local diffusion of technology

The impacts of technological change are not just evident in a global context. Technological change impacts every sector of the domestic economy. For example, the modern era in dentistry sees specialists in root canals (endodontists) performing root canals in the space of a single hour with the help of new technology; dental implants into bone, as an alternative to dentures, are commonplace; crowns can be machined with little human intervention; and X-rays are now performed with about one hundredth of the power formerly required. These technologies spread and are adopted through several channels. Dental practices do not usually compete on the basis of price, but if they do not adopt best practices and new technologies, then community word of mouth will see patients shifting to more efficient operators.

Some technological developments are protected by patents. But patent protection rarely inhibits new and more efficient practices that in some way mimic patent breakthroughs.

This page titled [8.7: Technological change- globalization and localization](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.8: Clusters, learning by doing, scope economics

Clusters

The phenomenon of a grouping of firms that specialize in producing related products is called a cluster. For example, Ottawa has more than its share of software development firms; Montreal has a disproportionate share of Canada's pharmaceutical producers and electronic game developers; Calgary has its 'oil patch'; Hollywood has movies; Toronto is Canada's financial capital, San Francisco and Seattle are leaders in new electronic products. Provincial and state capitals have most of their province's bureaucracy. Clusters give rise to externalities, frequently in the form of ideas that flow between firms, which in turn result in cost reductions and new products.

Cluster: a group of firms producing similar products, or engaged in similar research.

The most famous example of clustering is Silicon Valley, surrounding San Francisco, in California, the original high-tech cluster. The presence of a large group of firms with a common focus serves as a signal to workers with the right skill set that they are in demand in such a region. Furthermore, if these clusters are research oriented, as they frequently are, then knowledge spillovers benefit virtually all of the contiguous firms; when workers change employers, they bring their previously-learned skills with them; on social occasions, friends may chat about their work and interests and share ideas. This is a positive externality.

Learning by doing

Learning from production-related experiences frequently reduces costs: The accumulation of knowledge that is associated with having produced a large volume of output over a considerable time period enables managers to implement more efficient production methods and avoid errors. We give the term learning by doing to this accumulation of knowledge.

Examples abound, but the best known may be the continual improvement in the capacity of computer chips, whose efficiency has doubled about every eighteen months for several decades – a phenomenon known as Moore's Law. As *Intel Corporation* continues to produce chips it learns how to produce each succeeding generation of chips at lower cost. Past experience is key. Economies of scale and learning by doing therefore may not be independent: Large firms usually require time to grow or to attain a dominant role in their market, and this time and experience enables them to produce at lower cost. This lower cost in turn can solidify their market position further.

Learning by doing can reduce costs. A longer history of production enables firms to accumulate knowledge and thereby implement more efficient production processes.

Economies of scope

Economies of scope define a production process if the production of multiple products results in lower unit costs per product than if those products were produced alone. Scope economies, therefore, define the returns or cost reductions associated with broadening a firm's product range.

Corporations like Proctor and Gamble do not produce a single product in their health line; rather, they produce first aid, dental care, and baby care products. Cable companies offer their customers TV, high-speed Internet, and telephone services either individually or packaged. A central component of some new-economy multi-product firms is a technology platform that can be used for multiple purposes. We shall analyze the operation of these firms in more detail in Chapter 11.

Economies of scope occur if the unit cost of producing particular products is less when combined with the production of other products than when produced alone.

A platform is a hardware-cum-software capital installation that has multiple production capabilities

This page titled [8.8: Clusters, learning by doing, scope economics](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.9: Conclusion

Efficient production is critical to the survival of firms. Firms that do not adopt the most efficient production methods are likely to be left behind by their competitors. Efficiency translates into cost considerations, and the structure of costs in turn has a major impact on market type. Some sectors of the economy have very many firms (the restaurant business or the dry-cleaning business), whereas other sectors have few (internet providers or airlines). We will see in the following chapters how market structures depend critically upon the concept of scale economies that we have developed here.

This page titled [8.9: Conclusion](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.10: Key Terms

Production function: a technological relationship that specifies how much output can be produced with specific amounts of inputs.

Technological efficiency means that the maximum output is produced with the given set of inputs.

Economic efficiency defines a production structure that produces output at least cost.

Short run: a period during which at least one factor of production is fixed. If capital is fixed, then more output is produced by using additional labour.

Long run: a period of time that is sufficient to enable all factors of production to be adjusted.

Very long run: a period sufficiently long for new technology to develop.

Total product is the relationship between total output produced and the number of workers employed, for a given amount of capital.

Marginal product of labour is the addition to output produced by each additional worker. It is also the slope of the total product curve.

Law of diminishing returns: when increments of a variable factor (labour) are added to a fixed amount of another factor (capital), the marginal product of the variable factor must eventually decline.

Average product of labour is the number of units of output produced per unit of labour at different levels of employment.

Fixed costs are costs that are independent of the level of output.

Variable costs are related to the output produced.

Total cost is the sum of fixed cost and variable cost.

Average fixed cost is the total fixed cost per unit of output.

Average variable cost is the total variable cost per unit of output.

Average total cost is the sum of all costs per unit of output.

Marginal cost of production is the cost of producing each additional unit of output.

Sunk cost is a fixed cost that has already been incurred and cannot be recovered, even by producing a zero output.

Increasing returns to scale implies that, when all inputs are increased by a given proportion, output increases more than proportionately.

Constant returns to scale implies that output increases in direct proportion to an equal proportionate increase in all inputs.

Decreasing returns to scale implies that an equal proportionate increase in all inputs leads to a less than proportionate increase in output.

Long-run average total cost is the lower envelope of all the short-run *ATC* curves.

Minimum efficient scale defines a threshold size of operation such that scale economies are almost exhausted.

Long-run marginal cost is the increment in cost associated with producing one more unit of output when all inputs are adjusted in a cost minimizing manner.

Technological change represents innovation that can reduce the cost of production or bring new products on line.

Globalization is the tendency for international markets to be ever more integrated.

Cluster: a group of firms producing similar products, or engaged in similar research.

Learning by doing can reduce costs. A longer history of production enables firms to accumulate knowledge and thereby implement more efficient production processes.

Economies of scope occur if the unit cost of producing particular products is less when combined with the production of other products than when produced alone.

A platform is a hardware-cum-software capital installation that has multiple production capabilities

This page titled [8.10: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

8.11: Exercises for Chapter 8

EXERCISE 8.1

The relationship between output Q and the single variable input L is given by the form $Q = 5\sqrt{L}$. Capital is fixed. This relationship is given in the table below for a range of L values.

L	1	2	3	4	5	6	7	8	9	10	11	12
Q	5	7.07	8.66	10	11.18	12.25	13.23	14.14	15	15.81	16.58	17.32

- Add a row to this table and compute the MP .
- Draw the total product (TP) curve to scale, either on graph paper or in a spreadsheet.
- Inspect your graph to see if it displays diminishing MP .

EXERCISE 8.2

The TP for different output levels for Primitive Products is given in the table below.

Q	1	6	12	20	30	42	53	60	66	70
L	1	2	3	4	5	6	7	8	9	10

- Graph the TP curve to scale.
- Add a row to the table and enter the values of the MP of labour. Graph this in a separate diagram.
- Add a further row and compute the AP of labour. Add it to the graph containing the MP of labour.
- By inspecting the AP and MP graph, can you tell if you have drawn the curves correctly? How?

EXERCISE 8.3

A short-run relationship between output and total cost is given in the table below.

Output	0	1	2	3	4	5	6	7	8	9
Total Cost	12	27	40	51	61	70	80	91	104	120

- What is the total fixed cost of production in this example?
- Add four rows to the table and compute the TVC , AFC , AVC and ATC values for each level of output.
- Add one more row and compute the MC of producing additional output levels.
- Graph the MC and AC curves using the information you have developed.

EXERCISE 8.4

Consider the long-run total cost structure for the two firms A and B below.

Output	1	2	3	4	5	6	7
Total cost A	40	52	65	80	97	119	144
Total cost B	30	40	50	60	70	80	90

- Compute the long-run ATC curve for each firm.
- Plot these curves and examine the type of scale economies each firm experiences at different output levels.

EXERCISE 8.5

Use the data in Exercise 8.4,

- a. Calculate the long-run MC at each level of output for the two firms.
- b. Verify in a graph that these LMC values are consistent with the LAC values.

EXERCISE 8.6

Optional: Suppose you are told that a firm of interest has a long-run average total cost that is defined by the relationship $LATC=4+48/q$.

- a. In a table, compute the $LATC$ for output values ranging from 1...24. Plot the resulting $LATC$ curve.
 - b. What kind of returns to scale does this firm never experience?
 - c. By examining your graph, what will be the numerical value of the $LATC$ as output becomes very large?
 - d. Can you guess what the form of the long-run MC curve is?
1. Note that the long-run average total cost is not the collection of minimum points from each short-run average cost curve. The envelope of the short-run curves will pick up mainly points that are not at the minimum, as you will see if you try to draw the outcome. The intuition behind the definition is this: With increasing returns to scale, it may be better to build a plant size that operates with some spare capacity than to build one that is geared to producing a smaller output level. In building the larger plant, we can take greater advantage of the scale economies, and it may prove less costly to produce in such a plant than to produce with a smaller plant that has less unused capacity and does not exploit the underlying scale economies. Conversely, in the presence of decreasing returns to scale, it may be less costly to produce output in a plant that is used "overtime" than to use a larger plant that suffers from scale diseconomies.
 2. Two English farmers changed the lives of millions in the 18th century through their technological genius. The first was Charles Townsend, who introduced the concept of crop rotation. He realized that continual use of soil for a single purpose drained the land of key nutrients. This led him ultimately to propose a five-year rotation that included tillage, vegetables, and sheep. This rotation reduced the time during which land would lie fallow, and therefore increased the productivity of land. A second game-changing technological innovation saw the introduction of the seed plough, a tool that facilitated seed sowing in rows, rather than scattering it randomly. This line sowing meant that weeds could be controlled more easily—weeds that would otherwise smother much of the crop. The genius here was a gentleman by the name of Jethro Tull.

This page titled [8.11: Exercises for Chapter 8](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

SECTION OVERVIEW

Unit 4: Market Structures

Markets are all around us and they come in many forms. Some are on-line, others are physical. Some involve goods such as food and vehicles; others involve health provision or financial advice. Markets differ also by the degree of competition associated with each. For example the wholesale egg market is very homogeneous in that the product has minimal variation. The restaurant market offers food, but product variation is high. Each market has many suppliers.

In contrast to the egg and restaurant market, many other markets are characterized by just a few suppliers or in some cases just one. For example, passenger train services may have just a single supplier, and this supplier is therefore a monopolist. Pharmaceuticals tend to be supplied by a limited number of large international corporations plus a group of generic drug manufacturers. The internet and communications services market usually has only a handful of providers.

In this part we examine the reasons why markets take on different forms and display a variety of patterns of behaviour. We delve into the working of each market structure to understand why these markets retain their structure.

9: Perfect competition

- 9.1: The perfect competition paradigm
- 9.2: Market characteristics
- 9.3: The firm's supply decision
- 9.4: Dynamics- Entry and exit
- 9.5: Long-run industry supply
- 9.6: Globalization and technological change
- 9.7: Efficient resource allocation
- 9.8: Key Terms
- 9.9: Exercises for Chapter 9

10: Monopoly

- 10.1: Monopolies
- 10.2: Profit maximizing behaviour
- 10.3: Long-run choices
- 10.4: Output inefficiency
- 10.5: Price discrimination
- 10.6: Cartels- Acting like a monopolist
- 10.7: Invention, innovation and rent seeking
- 10.8: Conclusion
- 10.9: Key Terms
- 10.10: Exercises for Chapter 10

11: Imperfect competition

- 11.1: The principle ideas
- 11.2: Imperfect competitors
- 11.3: Imperfect competitors- measures of structure and market power
- 11.4: Imperfect competition- monopolistic competition
- 11.5: Imperfect competition- economies of scope and platforms
- 11.6: Strategic behaviour- Oligopoly and games
- 11.7: Strategic behaviour- Duopoly and Cournot games

[11.8: Strategic behaviour- Entry, exit and potential competition](#)

[11.9: Matching markets- design](#)

[11.10: Conclusion](#)

[11.11: Key Terms](#)

[11.12: Exercises for Chapter 11](#)

This page titled [Unit 4: Market Structures](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9: Perfect competition

Chapter 9: Perfect competition

In this chapter we will explore:

9.1	The competitive marketplace
9.2	Market characteristics
9.3	Supply in the short run
9.4	Dynamics: Entry and exit
9.5	Industry supply in the long run
9.6	Globalization and technological change
9.7	Perfect competition and market efficiency

9.1 The perfect competition paradigm

A competitive market is one that encompasses a very large number of suppliers, each producing a similar or identical product. Each supplier produces an output that forms a small part of the total market, and the sum of all of these individual outputs represents the production of that sector of the economy. Florists, barber shops, corner stores and dry cleaners all fit this description.

At the other extreme, a market that has just a single supplier is a monopolist. For example, the National Hockey League is the sole supplier of top-quality professional hockey games in North America; Hydro Quebec is a monopoly electricity distributor in Quebec; Via Rail is the only supplier of passenger rail services between Windsor, Ontario and the city of Quebec.

We use the word 'paradigm' in the title to this section: It implies that we will develop a *model* of supply behaviour for a market in which there are many small suppliers, producing essentially the same product, competing with one-another to meet the demands of consumers.

The structures that we call perfect competition and monopoly are extremes in the market place. Most sectors of the economy lie somewhere between these limiting cases. For example, the market for internet services usually contains several providers in any area – some provide using a fibre cable, others by satellite. The market for smart-phones in North America is dominated by two major players – *Apple* and *Samsung* (although there are several others). Hence, while these markets that have a limited number of suppliers are competitive in that they freely and perhaps fiercely compete for the buyer's expenditure, these are not perfectly competitive markets, because they do not have a very large number of suppliers.

In all of the models we develop in this chapter we will assume that the objective of firms is to maximize profit – the difference between revenues and costs.

A **perfectly competitive** industry is one in which many suppliers, producing an identical product, face many buyers, and no one participant can influence the market.

Profit maximization is the goal of competitive suppliers – they seek to maximize the difference between revenues and costs.

The presence of so many sellers in perfect competition means that each firm recognizes its own small size in relation to the total market, and that its actions have no perceptible impact on the market price for the good or service being traded. Each firm is therefore a *price taker*—in contrast to a monopolist, who is a *price setter*.

The same 'smallness' characteristic was assumed when we examined the demands of individuals earlier. Each buyer takes the price as given. He or she is not big enough to be able to influence the price. In contrast, when international airlines purchase or lease aircraft from *Boeing* or *Airbus*, they negotiate over the price and other conditions of supply. The market models underlying these types of transactions are examined in Chapter 11.

Hence, when we describe a market as being perfectly competitive we do not mean that other market types are not competitive; all market structure are competitive in the sense that the suppliers wish to make profit, and they produce as efficiently as possible in order to meet that goal.

9.2 Market characteristics

The key attributes of a perfectly competitive market are the following:

1. There must be *many firms*, each one so small that it cannot influence price or quantity in the industry, and powerless relative to the entire industry.
2. The *product must be standardized*. Barber shops offer a standard product, but a Lexus differs from a Ford. Barbers tend to be price takers, but Lexus does not charge the same price as Ford, and is a price setter.
3. Buyers are assumed to have *full information* about the product and its pricing. For example, buyers know that the products of different suppliers really are the same in quality.
4. There are *many buyers*.
5. There is *free entry and exit* of firms.

In terms of the demand curve that suppliers face, these market characteristics imply that the demand curve facing the perfectly competitive firm is horizontal, or infinitely elastic, as we defined in Chapter 4. In contrast, the demand curve facing the whole industry is downward sloping. The demand curve facing a firm is represented in Figure 9.1. It implies that the supplier can sell any output he chooses at the going price P_0 . He is a small player in the market, and variations in his output have no perceptible impact in the marketplace. But what quantity should he choose, or what quantity will maximize his profit? The profit-maximizing choice is his target, and the MC curve plays a key role in this decision.

9.3 The firm's supply decision

The concept of marginal revenue is key to analyzing the supply decision of an individual firm. We have used marginal analysis at several points to date. In consumer theory, we saw how consumers balance the utility per dollar *at the margin* in allocating their budget. Marginal revenue is the additional revenue accruing to the firm from the sale of one more unit of output.

Marginal revenue is the additional revenue accruing to the firm resulting from the sale of one more unit of output.

In perfect competition, a firm's marginal revenue (MR) is the price of the good. Since the price is constant for the individual supplier, each additional unit sold at the price P brings in the same additional revenue. Therefore, $P=MR$. For example, whether a dry cleaning business launders 10 shirts or 100 shirts per day, the price charged to customers is the same. This equality holds in no other market structure, as we shall see in the following chapters.

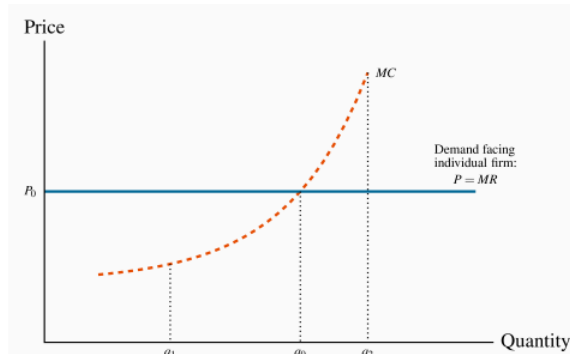
Supply in the short run

Recall how we defined the short run in the previous chapter: Each firm's plant size is fixed in the short run, so too is the number of firms in an industry. In the long run, each individual firm can change its scale of operation, and at the same time new firms can enter or existing firms can leave the industry.

Perfectly competitive suppliers face the choice of how much to produce at the going market price: That is, the amount that will maximize their profit. We abstract for the moment on how the price in the marketplace is determined. We shall see later in this chapter that it emerges as the value corresponding to the intersection of the supply and demand curves for the whole market – as described in Chapter 3.

The firm's MC curve is critical in defining the optimal amount to supply at any price. In Figure 9.1, MC is the firm's marginal cost curve in the short run. At the price P_0 the optimal amount to supply is q_0 , the amount determined by the intersection of the MC and the demand. To see why, imagine that the producer chose to supply the quantity q_1 . Such an output would leave the opportunity for further profit untapped. By producing one additional unit beyond q_1 , the supplier would get P_0 in additional revenue and incur an additional cost that is less than P_0 in producing this unit. In fact, on every unit between q_1 and q_0 he can make a profit, because the MR exceeds the associated cost, MC . By the same argument, it makes no sense to increase output beyond q_0 , to q_2 for example, because the cost of such additional units of output, MC , exceeds the revenue from them. *The MC therefore defines an optimal supply response.*

Figure 9.1 The competitive firm's optimal output



Here, q_0 represents the optimal supply decision when the price is P_0 . At output q_1 the cost of additional units is less than the revenue from such units and therefore it is profitable to increase output beyond q_1 . Conversely, at q_2 the MC of production exceeds the revenue obtained, and so output should be reduced.

Application Box 9.1 The law of one price

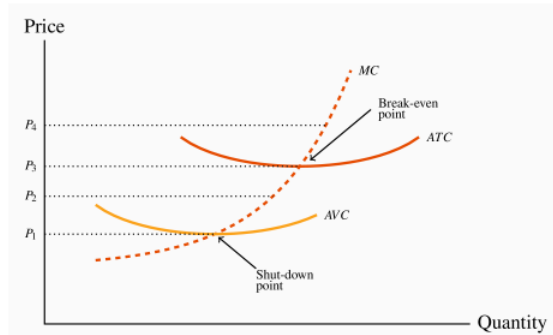
If information does not flow then prices in different parts of a market may differ and potential entrants may not know to enter a profitable market.

Consider the fishermen off the coast of Kerala, India in the late 1990s. Their market was studied by Robert Jensen, a development economist. Prior to 1997, fishermen tended to bring their fish to their home market or port. This was cheaper than venturing to other ports, particularly if there was no certainty regarding price. This practice resulted in prices that were high in some local markets and low in others – depending upon the daily catch. Frequently fish was thrown away in low-price markets even though it might have found a favourable price in another village's fish market.

This all changed with the advent of cell phones. Rather than head automatically to their home port, fishermen began to phone several different markets in the hope of finding a good price for their efforts. They began to form agreements with buyers before even bringing their catch to port. Economist Jensen observed a major decline in price variation between the markets that he surveyed. In effect the 'law of one price' came into being for sardines as a result of the introduction of cheap technology and the relatively free flow of information.

While the choice of the output q_0 is the best choice for the producer, Figure 9.1 does not tell us anything about profit. For that we need more information on costs. Accordingly, in Figure 9.2 the firm's AVC and ATC curves have been added to Figure 9.1. As explained in the previous chapter, the ATC curve includes both fixed and variable cost components, and the MC curve cuts the AVC and the ATC at their minima.

Figure 9.2 Short-run supply for the competitive firm



A price below P_1 does not cover variable costs, so the firm should shut down. Between prices P_1 and P_3 , the producer can cover variable, but not total, costs and therefore should produce in the short run if fixed costs are 'sunk'. In the long run the firm must close if the price does not reach P_3 . Profits are made if the price exceeds P_3 . The short-run supply curve is the quantity supplied at each price. It is therefore the MC curve above P_1 .

First, note that any price below P_3 , which corresponds to the minimum of the ATC curve, yields no profit, since it does not enable the producer to cover all of his costs. This price is therefore called the break-even price. Second, any price below P_1 , which corresponds to the minimum of the AVC , does not even enable the producer to cover variable costs. What about a price such as P_2 , that lies between these? The answer is that, if the supplier has already incurred some fixed costs, he should continue to produce,

provided he can cover his variable cost. But in the long run he must cover all of his costs, fixed and variable. Therefore, if the price falls below P_1 , he should shut down, even in the short run. This price is therefore called the shut-down price. If a price at least equal to P_1 cannot be sustained in the long run, he should leave the industry. But at a price such as P_2 he can cover variable costs and therefore should continue to produce in the short run. His optimal output at P_2 is defined by the intersection of the P_2 line with the MC curve. The firm's short-run supply curve is, therefore, that portion of the MC curve above the minimum of the AVC .

To illustrate this more concretely, consider again the example of our snowboard producer, and imagine that he is producing in a perfectly competitive marketplace. How should he behave in response to different prices? Table 9.1 reproduces the data from Table 8.2.

Table 9.1 Profit maximization in the short run

Labour	Output	Total	Average	Average	Marginal	Total	Profit
		Revenue \$	Variable	Total Cost	Cost \$	Cost \$	
			Cost	\$			
L	Q	TR	AVC	ATC	MC	TC	TR-TC
0	0					3,000	
1	15	1,050	66.67	266.67	66.67	4,000	-2,950
2	40	2,800	50.0	125.0	40.0	5,000	-2,200
3	70	4,900	42.86	85.71	33.33	6,000	-1,100
4	110	7,700	36.36	63.64	25.0	7,000	700
5	145	10,150	34.48	55.17	28.57	8,000	2,150
6	175	12,250	34.29	51.43	33.33	9,000	3,250
7	200	14,000	35.0	50.0	40.0	10,000	4,000
8	220	15,400	36.36	50.0	50.0	11,000	4,400
9	235	16,450	38.30	51.06	66.67	12,000	4,450
10	240	16,800	41.67	54.17	200.0	13,000	3,800

Output Price=\$70; Wage=\$1,000; Fixed Cost=\$3,000. The shut-down point occurs at a price of \$34.3, where the AVC attains a minimum. Hence no production, even in the short run, takes place unless the price exceeds this value. The break-even level of output occurs at a price of \$50, where the ATC attains a minimum.

The **shut-down price** corresponds to the minimum value of the AVC curve.

The **break-even price** corresponds to the minimum of the ATC curve.

The firm's **short-run supply curve** is that portion of the MC curve above the minimum of the AVC .

Suppose that the price is \$70. How many boards should he produce? The answer is defined by the behaviour of the MC curve. For any output less than or equal to 235, the MC is less than the price. For example, at $L=9$ and $Q=235$, the MC is \$66.67. At this output level, he makes a profit on the marginal unit produced, because the MC is less than the revenue he gets (\$70) from selling it.

But, at outputs above this, he registers a loss on the marginal units because the MC exceeds the revenue. For example, at $L=10$ and $Q=240$, the MC is \$200. Clearly, 235 snowboards is the optimum. To produce more would generate a loss on each additional unit, because the additional cost would exceed the additional revenue. Furthermore, to produce fewer snowboards would mean not availing of the potential for profit on additional boards.

His profit is based on the difference between revenue per unit and cost per unit at this output: $(P-ATC)$. Since the ATC for the 235 units produced by the nine workers is \$51.06, his profit margin is $\$70 - \$51.06 = \$18.94$ per board, and total profit is therefore $235 \times \$18.94 = \$4,450$.

Let us establish two other key outputs and prices for the producer. First, the shut-down point is the minimum of his AVC curve. Table 9.1 indicates that the price must be at least \$34.29 for him to be willing to supply *any* output, since that is the value of the AVC at its minimum. Second, the minimum of his ATC is at \$50. Accordingly, provided the price exceeds \$50, he will cover both

variable and fixed costs and make a maximum profit when he chooses an output where $P=MC$, above $P = \$50$. It follows that the short-run supply curve for Black Diamond Snowboards is the segment of the MC curve in Figure 8.4 above the AVC curve.

Given that we have developed the individual firm's supply curve, the next task is to develop the industry supply curve.

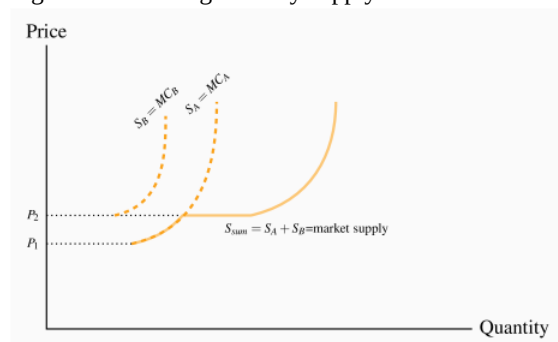
Industry supply in the short run

In Chapter 3 it was demonstrated that individual demands can be aggregated into an industry demand by summing them horizontally. The industry supply is obtained in exactly the same manner—by summing the firms' supply quantities across all firms in the industry.

To illustrate, imagine we have many firms, possibly operating at different scales of output and therefore having different short-run MC curves. The MC curves of two of these firms are illustrated in Figure 9.3. The MC of A is below the MC of B; therefore, B likely has a smaller scale of plant than A. Consider first the supply decisions in the price range P_1 to P_2 . At any price between these limits, only firm A will supply output – firm B does not cover its AVC in this price range. Therefore, the joint contribution to industry supply of firms A and B is given by the MC curve of firm A. But once a price of P_2 is attained, firm B is now willing to supply. The S_{sum} schedule is the horizontal addition of their supply quantities. Adding the supplies of every firm in the industry in this way yields the industry supply.

Industry supply (short run) in perfect competition is the horizontal sum of all firms' supply curves.

Figure 9.3 Deriving industry supply



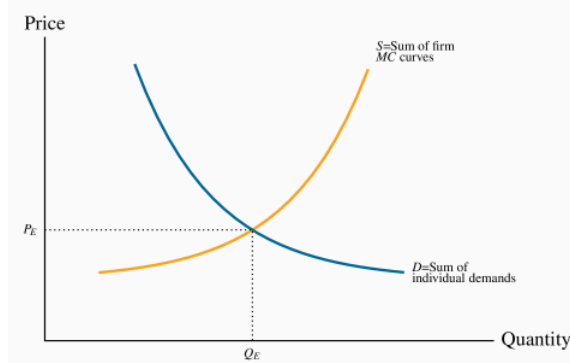
The marginal cost curves for firms A and B indicate that at any price below P_1 production is unprofitable and supply is therefore zero for both firms. At prices between P_1 and P_2 firm A is willing to supply, but not firm B. Consequently the market supply comes only from A. At prices above P_2 both firms are willing to supply. Therefore the market supply is the horizontal sum of each firm's supply.

Industry equilibrium

Consider next the industry equilibrium. Since the industry supply is the sum of the individual supplies, and the industry demand curve is the sum of individual demands, an equilibrium price and quantity (P_E, Q_E) are defined by the intersection of these industry-level curves, as in Figure 9.4. Here, each firm takes P_E as given (it is so small that it cannot influence the going price), and supplies an amount determined by the intersection of this price with its MC curve. The sum of such quantities is therefore Q_E .

Short-run equilibrium in perfect competition occurs when each firm maximizes profit by producing a quantity where $P=MC$, provided the price exceeds the minimum of the average variable cost.

Figure 9.4 Market equilibrium



The market supply curve S is the sum of each firm's supply or MC curve above the shut-down price. D is the sum of individual demands. The market equilibrium price and quantity are defined by P_E and Q_E .

9.4 Dynamics: Entry and exit

We have now described the market and firm-level equilibrium in the short run. However, this equilibrium may be only temporary; whether it can be sustained or not depends upon whether profits (or losses) are being incurred, or whether all participant firms are making what are termed normal profits. Such profits are considered an essential part of a firm's operation. They reflect the opportunity cost of the resources used in production. Firms do not operate if they cannot make a minimal, or normal, profit level. Above such profits are economic profits (also called supernormal profits), and these are what entice entry into the industry.

Recall from Chapter 7 that accounting and economic profits are different. The economist includes opportunity costs in determining profit, whereas the accountant considers actual revenues and costs. In the example developed in Section 7.2 the entrepreneur recorded accounting profit, but not economic profit. Suppose now that the numbers were slightly different, and are as defined in Table 9.2: Felicity invests \$250,000 in her business in the form of capital, as before. But she now has gross revenues of \$165,000 and incurs a cost of \$90,000 to buy the clothing wholesale that she then sells retail. She pays herself a salary of \$35,000. If these numbers represent her balance sheet, then she records an accounting profit of \$40,000.

Table 9.2 Economic profits

Sales		\$165,000
Materials costs		\$90,000
Wage costs		\$35,000
Accounting profit		\$40,000
Capital invested	\$250,000	
Implicit return on capital at 4%		\$10,000
Additional implicit wage costs		\$20,000
Total implicit costs		\$30,000
Economic profit		\$10,000

Her economic profit calculation must include opportunity costs. The opportunity cost of tying up \$250,000 of capital, if the interest rate is 4%, amounts to \$10,000. In addition, if Felicity could earn \$55,000 in her best alternative job then an additional implicit cost of \$20,000 must be considered. When these two opportunity (or implicit) costs are added to the balance sheet, her profit is reduced to \$10,000. This is her economic profit. If Felicity's economic profit is representative of the retail clothing sector of the economy, then that profitability should attract new entrepreneurs. Our conclusion is that this sector of the economy *should experience new entrants and hence an outward shift of the supply curve*. In contrast, in the numerical example considered in Section 7.2, Felicity was experiencing losses (negative economic profits), and in the longer term she would have to consider leaving the business. If she and other suppliers exited, then the market supply curve would shift back to the left – representing a reduction in supply.

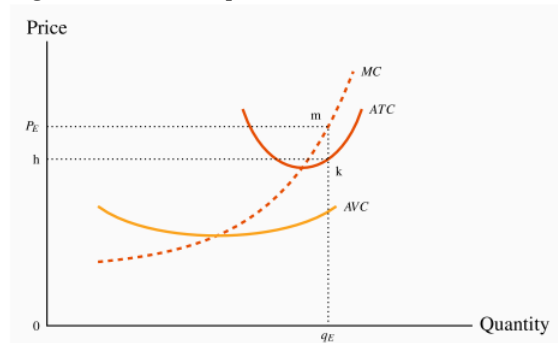
The critical point in this distinction between accounting and economic cost is that the decision to enter or leave a market in the longer term is based on what the entrepreneur can earn in the wider market place. That is, economic profits rather than accounting

profits will determine the equilibrium number of firms in the long term. In terms of our cost curves, we will assume that the full economic costs are included in the various curves that we use. Consequently any profits (or losses) that arise are based upon the full economic costs of the firm's operation.

Economic (supernormal) profits are those profits above normal profits that induce firms to enter an industry.

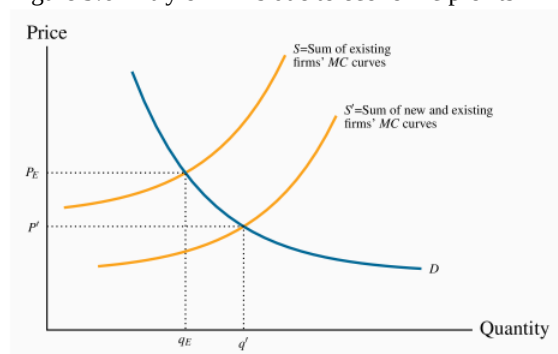
Let us return to our graphical analysis, and begin by supposing that the market equilibrium described in Figure 9.4 results in profits being made by some firms. Such an outcome is described in Figure 9.5, where the price exceeds the ATC . At the price P_E , a profit-making firm supplies the quantity q_E , as determined by its MC curve. On average, the cost of producing each unit of output, q_E , is defined by the point on the ATC at that output level, point k . Profit per unit is thus given by the value $(m-k)$ – the difference between revenue per unit and cost per unit. Total (economic) profit is therefore the area $P_E m k h$, which is quantity times profit per unit.

Figure 9.5 Short-run profits for the firm



At the price P_E , determined by the intersection of market demand and market supply, an individual firm produces the amount Q_E . The ATC of this output is k and therefore profit per unit is mk . Total profit is therefore $P_E m k h = Q_E \times mk = TR - TC$.

Figure 9.6 Entry of firms due to economic profits



If economic profits result from the price P_E new firms enter the industry. This entry increases the market supply to S' and the equilibrium price falls to P' . Entry continues as long as economic profits are present. Eventually the price is driven to a level where only normal profits are made, and entry ceases.

While q_E represents an equilibrium for the firm, it is only a short-run, or temporary, equilibrium for the industry. The assumption of free entry and exit implies that the presence of economic profits will induce new entrepreneurs to enter and start producing. The impact of this dynamic is illustrated in Figure 9.6. An increased number of firms shifts supply rightwards to become S' , thereby increasing the amount supplied at any price. The impact on price of this supply shift is evident: With an unchanged demand, the equilibrium price must fall.

How far will the price fall, and how many new firms will enter this profitable industry? As long as economic profits exist new firms will enter and the resulting increase in supply will continue to drive the price downwards. But, once the price has been driven down to the minimum of the ATC of a representative firm, there is no longer an incentive for new entrepreneurs to enter. Therefore, the long-run industry equilibrium is where the market price equals the minimum point of a firm's ATC . This generates normal profits, and there is no incentive for firms to enter or exit.

A **long-run equilibrium** in a competitive industry requires a price equal to the minimum point of a firm's ATC . At this point, only normal profits exist, and there is no incentive for firms to enter or exit.

In developing this dynamic, we began with a situation in which economic profits were present. However, we could have equally started from a position of losses. With a market price between the minimum of the AVC and the minimum of the ATC in Figure 9.5, revenues per unit would exceed variable costs but not total costs per unit. When firms cannot cover their ATC in the long run, they will cease production. Such closures must reduce aggregate supply; consequently the market supply curve contracts, rather than expands as it did in Figure 9.6. The reduced supply drives up the price of the good. This process continues as long as firms are making losses. A final industry equilibrium is attained only when the price reaches a level where firms can make a normal profit. Again, this will be at the minimum of the typical firm's ATC .

Accordingly, the long-run equilibrium is the same, regardless of whether we begin from a position in which firms are incurring losses, or where they are making profits.

Application Box 9.2 Entry and exit: Oil rigs

Oil drilling is a competitive market. There are a large number of suppliers, information is ubiquitous, and entry and exit are relatively free.

In the years 2012 and 2013 the price of crude oil was around \$100 US per barrel. Towards the end of 2014 the price of oil began to drop on world markets, and by early 2015 it fluctuated around \$50. The response of drillers in the US was substantial and immediate. The number of active rigs declined dramatically. In terms of our economic model, certain suppliers exited; they moth-balled their rigs and waited for the price of oil to recover.

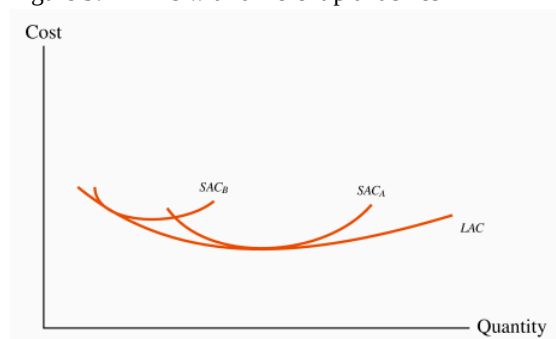
Another such cycle, even more pronounced, occurred in 2020. With the coronavirus pandemic, the demand for oil dropped and its price plummeted. Again, many firms shut down their rigs and had no choice but to sit out the price decline. A highly informative graphic is presented at tradingeconomics.com/united-states/crude-oil-rigs.

In addition to the decline in traditional oil recovery rigs, the number of operating shale crews declined by even greater amounts. Details at <https://www.forbes.com/sites/davidblackmon/2020/05/12/a-grim-earnings-season-for-the-us-shale-business/#6f55a95a1cf2>

9.5 Long-run industry supply

When aggregating the firm-level supply curves, as illustrated in Figure 9.3, we did not assume that all firms were identical. In that example, firm A has a cost structure with a lower AVC curve, since its supply curve starts at a lower dollar value. This indicates that firm A may have a larger plant size than firm B – one that puts A closer to the minimum efficient scale region of its long-run ATC curve.

Figure 9.7 Firms with different plant sizes



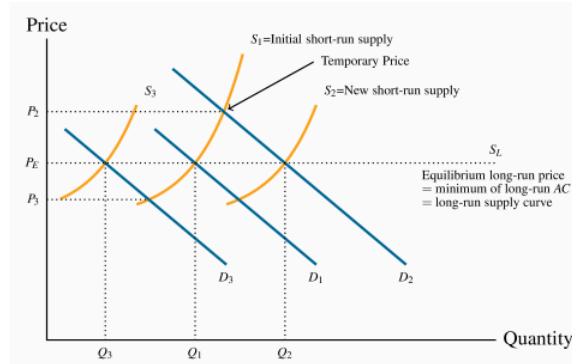
Firm B cannot compete with Firm A in the long run given that B has a less efficient plant size than firm A. The equilibrium long-run price equals the minimum of the LAC . At this price firm B must move to a more efficient plant size or make losses.

Can firm B survive with his current scale of operation in the long run? Our industry dynamics indicate that it cannot. The reason is that, provided *some* firms are making economic profits, new entrepreneurs will enter the industry and drive the price down to the minimum of the ATC curve of those firms who are operating with the lowest cost plant size. B-type firms will therefore be forced either to leave the industry or to adjust to the least-cost plant size—corresponding to the lowest point on its long-run ATC curve. Remember that the same technology is available to all firms; they each have the same long-run ATC curve, and may choose different scales of operation in the short run, as illustrated in Figure 9.7. But in the long run they must all produce using the minimum-cost plant size, or else they will be driven from the market.

This behaviour enables us to define a long-run industry supply. The long run involves the entry and exit of firms, and leads to a price corresponding to the minimum of the long-run ATC curve. Therefore, if the long-run equilibrium price corresponds to this minimum, *the long-run supply curve of the industry is defined by a particular price value—it is horizontal at the price corresponding to the minimum of the LATC*. More or less output is produced as a result of firms entering or leaving the industry, with those present always producing at the same unit cost in a long-run equilibrium.

Industry supply in the long run in perfect competition is horizontal at a price corresponding to the minimum of the representative firm's long-run ATC curve.

Figure 9.8 Long-run dynamics



The LR equilibrium price P_E is disturbed by a shift in demand from D_1 to D_2 . With a fixed number of firms, P_2 results. Profits accrue at this price and entry occurs. Therefore the SR supply shifts outwards until these profits are eroded and the new equilibrium output is Q_2 . If, instead, D falls to D_3 then firms exit because they make losses, S shifts back until the price is driven up sufficiently to restore normal profits. Different outputs are supplied in the long run at the same price P_E , therefore the long-run supply is horizontal at P_E .

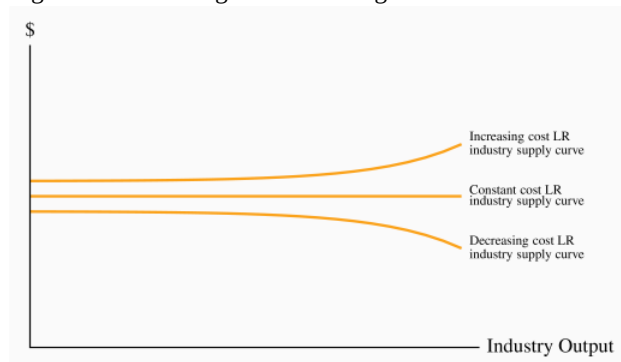
This industry's long-run supply curve, S_L , and a particular short-run supply are illustrated in Figure 9.8. Different points on S_L are attained when demand shifts. Suppose that, from an initial equilibrium Q_1 , defined by the intersection of D_1 and S_1 , demand increases from D_1 to D_2 because of a growth in income. With a fixed number of firms, the additional demand can be met only at a higher price (P_2), where each existing firm produces more using their existing plant size. The economic profits that result induce new operators to produce. This addition to the industry's production capacity shifts the short-run supply outwards and price declines until normal profits are once again being made. The new long-run equilibrium is at Q_2 , with more firms each producing at the minimum of their long-run ATC curve, P_E .

The same dynamic would describe the industry reaction to a decline in demand—price would fall, some firms would exit, and the resulting contraction in supply would force the price back up to the long-run equilibrium level. This is illustrated by a decline in demand from D_1 to D_3 .

Increasing and decreasing cost industries

While a horizontal long-run supply is the norm for perfect competition, in some industries costs increase with the scale of industry output; in others they decrease. This may be because all of the producers use a particular input that itself becomes more or less costly, depending upon the amount supplied.

Figure 9.9 Increasing and decreasing cost industries



When individual-supplier costs rise as the output of the industry increases we have an increasing cost supply curve for the industry in the long run. Conversely, when the costs of individual suppliers fall with the scale of the industry, we have a decreasing cost industry.

Decreasing cost sectors are those that benefit from a decline in the prices of their inputs as the size of their market expands. This is frequently because the suppliers of the inputs themselves can benefit from scale economies as a result of expansion in the market for the final good. A case in point has been the computer market, or the tablet market: As output in these markets has grown, the producers of videocards and random-access memory have benefited from scale economies and thus been able to sell these components at a lower price to the manufacturers of the final goods. An example of an increasing cost market is the market for landings and take-offs at airports. Airports are frequently limited in their ability to expand their size and build additional runways. In such markets, as use grows, planes about to land may have to adopt a circling holding pattern, while those departing encounter clearance delays. Such delays increase the time costs to passengers and the fuel and labour costs to the suppliers. Decreasing and increasing industry costs are reflected in the long-run industry supply curve by a downward-sloping segment or an upward sloping segment, as illustrated in Figure 9.9.

Increasing (decreasing) cost industry is one where costs rise (fall) for each firm because of the scale of industry operation.

9.6 Globalization and technological change

Globalization and technological change have had a profound impact on the way goods and services are produced and brought to market in the modern world. The cost structure of many firms has been reduced by *outsourcing to lower-wage economies*. Furthermore, the advent of the communications revolution has effectively *increased the minimum efficient scale for many industries*, as illustrated in Chapter 8 (Figure 8.7). Larger firms are less difficult to manage nowadays, and the *LAC* curve may not slope upwards until very high output levels are attained. The consequence is that some industries may not have sufficient "production space" to sustain a large number of firms. In order to reap the advantages of scale economies, firms become so large that they can supply a significant part of the market. They are no longer so small as to have no impact on the price.

Outsourcing and easier communications have in many cases simply eliminated many industries in the developed world. Garment making is an example. Some decades ago Quebec was Canada's main garment maker: Brokers dealt with 'cottage-type' garment assemblers outside Montreal and Quebec City. But ultimately the availability of cheaper labour in the developing world combined with efficient communications undercut the local manufacture. Most of Canada's garments are now imported. Other North American and European industries have been impacted in similar ways. Displaced labour has had to reskill, retool, reeducate itself, and either seek alternative employment in the manufacturing sector, or move to the service sector of the economy, or retire.

Globalization has had a third impact on the domestic economy, in so far as it *reduces the cost of components*. Even industries that continue to operate within national boundaries see a reduction in their cost structure on account of globalization's impact on input costs. This is particularly in evidence in the computing industry, where components are produced in numerous low-wage economies, imported to North America and assembled into computers domestically. Such components are termed *intermediate goods*.

9.7 Efficient resource allocation

Economists have a particular liking for competitive markets. The reason is not, as is frequently thought, that we love competitive battles; it really concerns resource allocation in the economy at large. In Chapter 5 we explained why markets are frequently an excellent vehicle for transporting the economy's resources to where they are most valued: A perfectly competitive marketplace in which there are no externalities results in resources being used up to the point where the demand and supply prices are equal. If demand is a measure of marginal benefit and supply is a measure of marginal cost, then a perfectly competitive market ensures that this condition will hold in equilibrium. *Perfect competition, therefore, results in resources being used efficiently.*

Our initial reaction to this perspective may be: If market equilibrium is such that the quantity supplied always equals the quantity demanded, is not every market efficient? The answer is no. As we shall see in the next chapter on monopoly, the monopolist's supply decision does not reflect the marginal cost of resources used in production, and therefore does not result in an efficient allocation in the economy.

Key Terms

Perfect competition: an industry in which many suppliers, producing an identical product, face many buyers, and no one participant can influence the market.

Profit maximization is the goal of competitive suppliers – they seek to maximize the difference between revenues and costs.

Marginal revenue is the additional revenue accruing to the firm resulting from the sale of one more unit of output.

Shut-down price corresponds to the minimum value of the *AVC* curve.

Break-even price corresponds to the minimum of the *ATC* curve.

Short-run supply curve for perfect competitor: the portion of the *MC* curve above the minimum of the *AVC*.

Industry supply (short run) in perfect competition is the horizontal sum of all firms' supply curves.

Short-run equilibrium in perfect competition occurs when each firm maximizes profit by producing a quantity where $P=MC$.

Economic (supernormal) profits are those profits above normal profits that induce firms to enter an industry. Economic profits are based on the opportunity cost of the resources used in production.

Long-run equilibrium in a competitive industry requires a price equal to the minimum point of a firm's *ATC*. At this point, only normal profits exist, and there is no incentive for firms to enter or exit.

Industry supply in the long run in perfect competition is horizontal at a price corresponding to the minimum of the representative firm's long-run *ATC* curve.

Increasing (decreasing) cost industry is one where costs rise (fall) for each firm because of the scale of industry operation.

Exercises for Chapter 9

EXERCISE 9.1

Wendy's Window Cleaning is a small local operation. Wendy presently cleans the outside windows in her neighbours' houses for \$36 per house. She does ten houses per day. She is incurring total costs of \$420, and of this amount \$100 is fixed. The cost per house is constant.

1. What is the marginal cost associated with cleaning the windows of one house – we know it is constant?
2. At a price of \$36, what is her break-even level of output (number of houses)?
3. If the fixed cost is 'sunk' and she cannot increase her output in the short run, should she shut down?

EXERCISE 9.2

A manufacturer of vacuum cleaners incurs a constant variable cost of production equal to \$80. She can sell the appliances to a wholesaler for \$130. Her annual fixed costs are \$200,000. How many vacuums must she sell in order to cover her total costs?

EXERCISE 9.3

For the vacuum cleaner producer in Exercise 9.2:

1. Draw the *MC* curve.
2. Next, draw her *AFC* and her *AVC* curves.
3. Finally, draw her *ATC* curve.
4. In order for this cost structure to be compatible with a perfectly competitive industry, what must happen to her *MC* curve at some output level?

EXERCISE 9.4

Consider the supply curves of two firms in a competitive industry: $P=q_A$ and $P=2q_B$.

1. On a diagram, draw these two supply curves, marking their intercepts and slopes numerically (remember that they are really *MC* curves).
2. Now draw a supply curve that represents the combined supply of these two firms.

EXERCISE 9.5

Amanda's Apple Orchard Productions Limited produces 10,000 kilograms of apples per month. Her total production costs at this output level are \$8,000. Two of her many competitors have larger-scale operations and produce 12,000 and 15,000 kilos at total costs of \$9,500 and \$11,000 respectively. If this industry is competitive, on what segment of the *LAC* curve are these producers producing?

EXERCISE 9.6

Consider the data in the table below. TC is total cost, TR is total revenue, and Q is output.

Q	0	1	2	3	4	5	6	7	8	9	10
TC	10	18	24	31	39	48	58	69	82	100	120
TR	0	11	22	33	44	55	66	77	88	99	110

1. Add some extra rows to the table and for each level of output calculate the MR , the MC and total profit.
2. Next, compute AFC , AVC , and ATC for each output level, and draw these three cost curves on a diagram.
3. What is the profit-maximizing output?
4. How can you tell that this firm is in a competitive industry?

EXERCISE 9.7

Optional: The market demand and supply curves in a perfectly competitive industry are given by: $Q_d = 30,000 - 600P$ and $Q_s = 200P - 2000$.

1. Draw these functions on a diagram, and calculate the equilibrium price of output in this industry.
2. Now assume that an additional firm is considering entering. This firm has a short-run MC curve defined by $MC = 10 + 0.5q$, where q is the firm's output. If this firm enters the industry and it knows the equilibrium price in the industry, what output should it produce?

EXERCISE 9.8

Optional: Consider two firms in a perfectly competitive industry. They have the same MC curves and differ only in having higher and lower fixed costs. Suppose the ATC curves are of the form: $400/q + 10 + (1/4)q$ and $225/q + 10 + (1/4)q$. The MC for each is a straight line: $MC = 10 + (1/2)q$.

1. In the first column of a spreadsheet enter quantity values of 1, 5, 10, 15, 20,..., 50. In the following columns compute the ATC curves for each quantity value.
2. Compute the MC at each output in the next column, and plot all three curves.
3. Compute the break-even price for each firm.
4. Explain why both of these firms cannot continue to produce in the long run in a perfectly competitive market.

This page titled 9: Perfect competition is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis](#) and [Ian Irvine](#) (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.1: The perfect competition paradigm

A competitive market is one that encompasses a very large number of suppliers, each producing a similar or identical product. Each supplier produces an output that forms a small part of the total market, and the sum of all of these individual outputs represents the production of that sector of the economy. Florists, barber shops, corner stores and dry cleaners all fit this description.

At the other extreme, a market that has just a single supplier is a monopolist. For example, the National Hockey League is the sole supplier of top-quality professional hockey games in North America; Hydro Quebec is a monopoly electricity distributor in Quebec; Via Rail is the only supplier of passenger rail services between Windsor, Ontario and the city of Quebec.

We use the word 'paradigm' in the title to this section: It implies that we will develop a *model* of supply behaviour for a market in which there are many small suppliers, producing essentially the same product, competing with one-another to meet the demands of consumers.

The structures that we call perfect competition and monopoly are extremes in the market place. Most sectors of the economy lie somewhere between these limiting cases. For example, the market for internet services usually contains several providers in any area – some provide using a fibre cable, others by satellite. The market for smart-phones in North America is dominated by two major players – *Apple* and *Samsung* (although there are several others). Hence, while these markets that have a limited number of suppliers are competitive in that they freely and perhaps fiercely compete for the buyer's expenditure, these are not perfectly competitive markets, because they do not have a very large number of suppliers.

In all of the models we develop in this chapter we will assume that the objective of firms is to maximize profit – the difference between revenues and costs.

A **perfectly competitive** industry is one in which many suppliers, producing an identical product, face many buyers, and no one participant can influence the market.

Profit maximization is the goal of competitive suppliers – they seek to maximize the difference between revenues and costs.

The presence of so many sellers in perfect competition means that each firm recognizes its own small size in relation to the total market, and that its actions have no perceptible impact on the market price for the good or service being traded. Each firm is therefore a *price taker*—in contrast to a monopolist, who is a *price setter*.

The same 'smallness' characteristic was assumed when we examined the demands of individuals earlier. Each buyer takes the price as given. He or she is not big enough to be able to influence the price. In contrast, when international airlines purchase or lease aircraft from *Boeing* or *Airbus*, they negotiate over the price and other conditions of supply. The market models underlying these types of transactions are examined in Chapter 11.

Hence, when we describe a market as being perfectly competitive we do not mean that other market types are not competitive; all market structure are competitive in the sense that the suppliers wish to make profit, and they produce as efficiently as possible in order to meet that goal.

This page titled [9.1: The perfect competition paradigm](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.2: Market characteristics

The key attributes of a perfectly competitive market are the following:

1. There must be *many firms*, each one so small that it cannot influence price or quantity in the industry, and powerless relative to the entire industry.
2. The *product must be standardized*. Barber shops offer a standard product, but a Lexus differs from a Ford. Barbers tend to be price takers, but Lexus does not charge the same price as Ford, and is a price setter.
3. Buyers are assumed to have *full information* about the product and its pricing. For example, buyers know that the products of different suppliers really are the same in quality.
4. There are *many buyers*.
5. There is *free entry and exit* of firms.

In terms of the demand curve that suppliers face, these market characteristics imply that the demand curve facing the perfectly competitive firm is horizontal, or infinitely elastic, as we defined in Chapter 4. In contrast, the demand curve facing the whole industry is downward sloping. The demand curve facing a firm is represented in Figure 9.1. It implies that the supplier can sell any output he chooses at the going price P_0 . He is a small player in the market, and variations in his output have no perceptible impact in the marketplace. But what quantity should he choose, or what quantity will maximize his profit? The profit-maximizing choice is his target, and the *MC* curve plays a key role in this decision.

This page titled [9.2: Market characteristics](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.3: The firm's supply decision

The concept of marginal revenue is key to analyzing the supply decision of an individual firm. We have used marginal analysis at several points to date. In consumer theory, we saw how consumers balance the utility per dollar *at the margin* in allocating their budget. Marginal revenue is the additional revenue accruing to the firm from the sale of one more unit of output.

Marginal revenue is the additional revenue accruing to the firm resulting from the sale of one more unit of output.

In perfect competition, a firm's marginal revenue (MR) is the price of the good. Since the price is constant for the individual supplier, each additional unit sold at the price P brings in the same additional revenue. Therefore, $P=MR$. For example, whether a dry cleaning business launders 10 shirts or 100 shirts per day, the price charged to customers is the same. This equality holds in no other market structure, as we shall see in the following chapters.

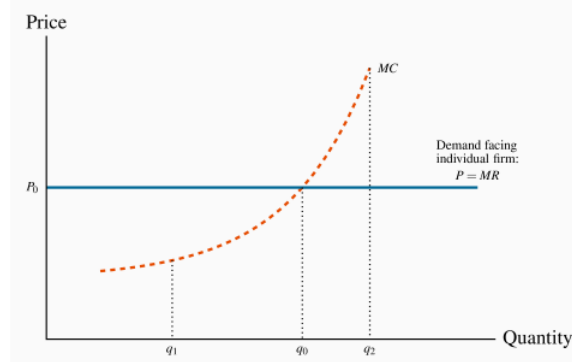
Supply in the short run

Recall how we defined the short run in the previous chapter: Each firm's plant size is fixed in the short run, so too is the number of firms in an industry. In the long run, each individual firm can change its scale of operation, and at the same time new firms can enter or existing firms can leave the industry.

Perfectly competitive suppliers face the choice of how much to produce at the going market price: That is, the amount that will maximize their profit. We abstract for the moment on how the price in the marketplace is determined. We shall see later in this chapter that it emerges as the value corresponding to the intersection of the supply and demand curves for the whole market – as described in Chapter 3.

The firm's MC curve is critical in defining the optimal amount to supply at any price. In Figure 9.1, MC is the firm's marginal cost curve in the short run. At the price P_0 the optimal amount to supply is q_0 , the amount determined by the intersection of the MC and the demand. To see why, imagine that the producer chose to supply the quantity q_1 . Such an output would leave the opportunity for further profit untapped. By producing one additional unit beyond q_1 , the supplier would get P_0 in additional revenue and incur an additional cost that is less than P_0 in producing this unit. In fact, on every unit between q_1 and q_0 he can make a profit, because the MR exceeds the associated cost, MC . By the same argument, it makes no sense to increase output beyond q_0 , to q_2 for example, because the cost of such additional units of output, MC , exceeds the revenue from them. *The MC therefore defines an optimal supply response.*

Figure 9.1 The competitive firm's optimal output



Here, q_0 represents the optimal supply decision when the price is P_0 . At output q_1 the cost of additional units is less than the revenue from such units and therefore it is profitable to increase output beyond q_1 . Conversely, at q_2 the MC of production exceeds the revenue obtained, and so output should be reduced.

Application Box 9.1 The law of one price

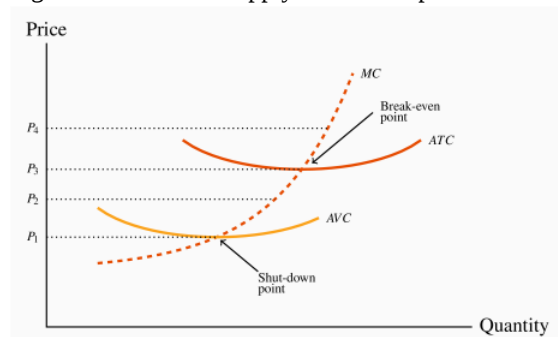
If information does not flow then prices in different parts of a market may differ and potential entrants may not know to enter a profitable market.

Consider the fishermen off the coast of Kerala, India in the late 1990s. Their market was studied by Robert Jensen, a development economist. Prior to 1997, fishermen tended to bring their fish to their home market or port. This was cheaper than venturing to other ports, particularly if there was no certainty regarding price. This practice resulted in prices that were high in some local markets and low in others – depending upon the daily catch. Frequently fish was thrown away in low-price markets even though it might have found a favourable price in another village's fish market.

This all changed with the advent of cell phones. Rather than head automatically to their home port, fishermen began to phone several different markets in the hope of finding a good price for their efforts. They began to form agreements with buyers before even bringing their catch to port. Economist Jensen observed a major decline in price variation between the markets that he surveyed. In effect the 'law of one price' came into being for sardines as a result of the introduction of cheap technology and the relatively free flow of information.

While the choice of the output q_0 is the best choice for the producer, Figure 9.1 does not tell us anything about profit. For that we need more information on costs. Accordingly, in Figure 9.2 the firm's AVC and ATC curves have been added to Figure 9.1. As explained in the previous chapter, the ATC curve includes both fixed and variable cost components, and the MC curve cuts the AVC and the ATC at their minima.

Figure 9.2 Short-run supply for the competitive firm



A price below P_1 does not cover variable costs, so the firm should shut down. Between prices P_1 and P_3 , the producer can cover variable, but not total, costs and therefore should produce in the short run if fixed costs are 'sunk'. In the long run the firm must close if the price does not reach P_3 . Profits are made if the price exceeds P_3 . The short-run supply curve is the quantity supplied at each price. It is therefore the MC curve above P_1 .

First, note that any price below P_3 , which corresponds to the minimum of the ATC curve, yields no profit, since it does not enable the producer to cover all of his costs. This price is therefore called the break-even price. Second, any price below P_1 , which corresponds to the minimum of the AVC , does not even enable the producer to cover variable costs. What about a price such as P_2 , that lies between these? The answer is that, if the supplier has already incurred some fixed costs, he should continue to produce, provided he can cover his variable cost. But in the long run he must cover all of his costs, fixed and variable. Therefore, if the price falls below P_1 , he should shut down, even in the short run. This price is therefore called the shut-down price. If a price at least equal to P_3 cannot be sustained in the long run, he should leave the industry. But at a price such as P_2 he can cover variable costs and therefore should continue to produce in the short run. His optimal output at P_2 is defined by the intersection of the P_2 line with the MC curve. The firm's short-run supply curve is, therefore, that portion of the MC curve above the minimum of the AVC .

To illustrate this more concretely, consider again the example of our snowboard producer, and imagine that he is producing in a perfectly competitive marketplace. How should he behave in response to different prices? Table 9.1 reproduces the data from Table 8.2.

Table 9.1 Profit maximization in the short run

Labour	Output	Total	Average	Average	Marginal	Total	Profit
		Revenue \$	Variable	Total Cost	Cost \$	Cost \$	
			Cost	\$			
L	Q	TR	AVC	ATC	MC	TC	TR-TC
0	0					3,000	
1	15	1,050	66.67	266.67	66.67	4,000	-2,950
2	40	2,800	50.0	125.0	40.0	5,000	-2,200
3	70	4,900	42.86	85.71	33.33	6,000	-1,100
4	110	7,700	36.36	63.64	25.0	7,000	700
5	145	10,150	34.48	55.17	28.57	8,000	2,150

6	175	12,250	34.29	51.43	33.33	9,000	3,250
7	200	14,000	35.0	50.0	40.0	10,000	4,000
8	220	15,400	36.36	50.0	50.0	11,000	4,400
9	235	16,450	38.30	51.06	66.67	12,000	4,450
10	240	16,800	41.67	54.17	200.0	13,000	3,800

Output Price=\$70; Wage=\$1,000; Fixed Cost=\$3,000. The shut-down point occurs at a price of \$34.3, where the *AVC* attains a minimum. Hence no production, even in the short run, takes place unless the price exceeds this value. The break-even level of output occurs at a price of \$50, where the *ATC* attains a minimum.

The **shut-down price** corresponds to the minimum value of the *AVC* curve.

The **break-even price** corresponds to the minimum of the *ATC* curve.

The firm's **short-run supply curve** is that portion of the *MC* curve above the minimum of the *AVC*.

Suppose that the price is \$70. How many boards should he produce? The answer is defined by the behaviour of the *MC* curve. For any output less than or equal to 235, the *MC* is less than the price. For example, at $L=9$ and $Q=235$, the *MC* is \$66.67. At this output level, he makes a profit on the marginal unit produced, because the *MC* is less than the revenue he gets (\$70) from selling it.

But, at outputs above this, he registers a loss on the marginal units because the *MC* exceeds the revenue. For example, at $L=10$ and $Q=240$, the *MC* is \$200. Clearly, 235 snowboards is the optimum. To produce more would generate a loss on each additional unit, because the additional cost would exceed the additional revenue. Furthermore, to produce fewer snowboards would mean not availing of the potential for profit on additional boards.

His profit is based on the difference between revenue per unit and cost per unit at this output: $(P-ATC)$. Since the *ATC* for the 235 units produced by the nine workers is \$51.06, his profit margin is $\$70 - \$51.06 = \$18.94$ per board, and total profit is therefore $235 \times \$18.94 = \$4,450$.

Let us establish two other key outputs and prices for the producer. First, the shut-down point is the minimum of his *AVC* curve. Table 9.1 indicates that the price must be at least \$34.29 for him to be willing to supply *any* output, since that is the value of the *AVC* at its minimum. Second, the minimum of his *ATC* is at \$50. Accordingly, provided the price exceeds \$50, he will cover both variable and fixed costs and make a maximum profit when he chooses an output where $P=MC$, above $P = \$50$. It follows that the short-run supply curve for Black Diamond Snowboards is the segment of the *MC* curve in Figure 8.4 above the *AVC* curve.

Given that we have developed the individual firm's supply curve, the next task is to develop the industry supply curve.

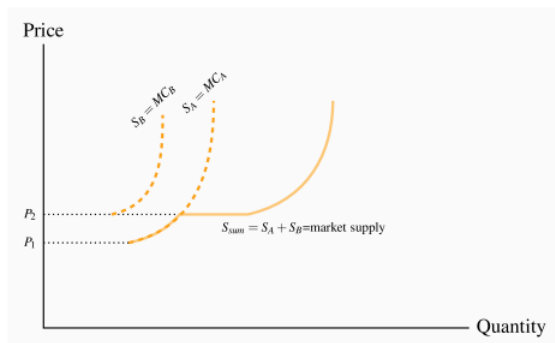
Industry supply in the short run

In Chapter 3 it was demonstrated that individual demands can be aggregated into an industry demand by summing them horizontally. The industry supply is obtained in exactly the same manner—by summing the firms' supply quantities across all firms in the industry.

To illustrate, imagine we have many firms, possibly operating at different scales of output and therefore having different short-run *MC* curves. The *MC* curves of two of these firms are illustrated in Figure 9.3. The *MC* of A is below the *MC* of B; therefore, B likely has a smaller scale of plant than A. Consider first the supply decisions in the price range P_1 to P_2 . At any price between these limits, only firm A will supply output – firm B does not cover its *AVC* in this price range. Therefore, the joint contribution to industry supply of firms A and B is given by the *MC* curve of firm A. But once a price of P_2 is attained, firm B is now willing to supply. The S_{sum} schedule is the horizontal addition of their supply quantities. Adding the supplies of every firm in the industry in this way yields the industry supply.

Industry supply (short run) in perfect competition is the horizontal sum of all firms' supply curves.

Figure 9.3 Deriving industry supply



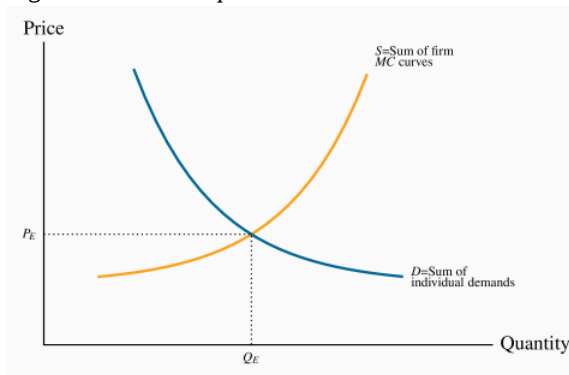
The marginal cost curves for firms A and B indicate that at any price below P_1 production is unprofitable and supply is therefore zero for both firms. At prices between P_1 and P_2 firm A is willing to supply, but not firm B. Consequently the market supply comes only from A. At prices above P_2 both firms are willing to supply. Therefore the market supply is the horizontal sum of each firm's supply.

Industry equilibrium

Consider next the industry equilibrium. Since the industry supply is the sum of the individual supplies, and the industry demand curve is the sum of individual demands, an equilibrium price and quantity (P_E, Q_E) are defined by the intersection of these industry-level curves, as in Figure 9.4. Here, each firm takes P_E as given (it is so small that it cannot influence the going price), and supplies an amount determined by the intersection of this price with its MC curve. The sum of such quantities is therefore Q_E .

Short-run equilibrium in perfect competition occurs when each firm maximizes profit by producing a quantity where $P=MC$, provided the price exceeds the minimum of the average variable cost.

Figure 9.4 Market equilibrium



The market supply curve S is the sum of each firm's supply or MC curve above the shut-down price. D is the sum of individual demands. The market equilibrium price and quantity are defined by P_E and Q_E .

This page titled 9.3: The firm's supply decision is shared under a [CC BY-NC-SA 4.0](https://creativecommons.org/licenses/by-nc-sa/4.0/) license and was authored, remixed, and/or curated by Douglas Curtis and Ian Irvine (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.4: Dynamics- Entry and exit

We have now described the market and firm-level equilibrium in the short run. However, this equilibrium may be only temporary; whether it can be sustained or not depends upon whether profits (or losses) are being incurred, or whether all participant firms are making what are termed normal profits. Such profits are considered an essential part of a firm's operation. They reflect the opportunity cost of the resources used in production. Firms do not operate if they cannot make a minimal, or normal, profit level. Above such profits are economic profits (also called supernormal profits), and these are what entice entry into the industry.

Recall from Chapter 7 that accounting and economic profits are different. The economist includes opportunity costs in determining profit, whereas the accountant considers actual revenues and costs. In the example developed in Section 7.2 the entrepreneur recorded accounting profit, but not economic profit. Suppose now that the numbers were slightly different, and are as defined in Table 9.2: Felicity invests \$250,000 in her business in the form of capital, as before. But she now has gross revenues of \$165,000 and incurs a cost of \$90,000 to buy the clothing wholesale that she then sells retail. She pays herself a salary of \$35,000. If these numbers represent her balance sheet, then she records an accounting profit of \$40,000.

Table 9.2 Economic profits

Sales		\$165,000
Materials costs		\$90,000
Wage costs		\$35,000
Accounting profit		\$40,000
Capital invested	\$250,000	
Implicit return on capital at 4%		\$10,000
Additional implicit wage costs		\$20,000
Total implicit costs		\$30,000
Economic profit		\$10,000

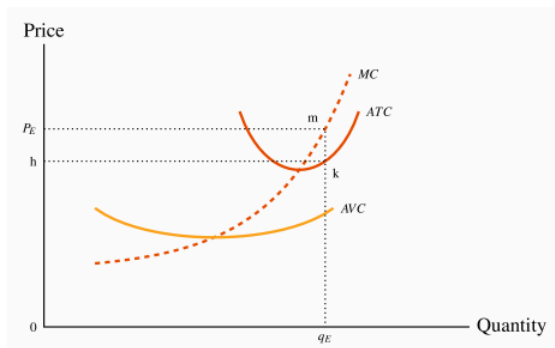
Her economic profit calculation must include opportunity costs. The opportunity cost of tying up \$250,000 of capital, if the interest rate is 4%, amounts to \$10,000. In addition, if Felicity could earn \$55,000 in her best alternative job then an additional implicit cost of \$20,000 must be considered. When these two opportunity (or implicit) costs are added to the balance sheet, her profit is reduced to \$10,000. This is her economic profit. If Felicity's economic profit is representative of the retail clothing sector of the economy, then that profitability should attract new entrepreneurs. Our conclusion is that this sector of the economy *should experience new entrants and hence an outward shift of the supply curve*. In contrast, in the numerical example considered in Section 7.2, Felicity was experiencing losses (negative economic profits), and in the longer term she would have to consider leaving the business. If she and other suppliers exited, then the market supply curve would shift back to the left – representing a reduction in supply.

The critical point in this distinction between accounting and economic cost is that the decision to enter or leave a market in the longer term is based on what the entrepreneur can earn in the wider market place. That is, economic profits rather than accounting profits will determine the equilibrium number of firms in the long term. In terms of our cost curves, we will assume that the full economic costs are included in the various curves that we use. Consequently any profits (or losses) that arise are based upon the full economic costs of the firm's operation.

Economic (supernormal) profits are those profits above normal profits that induce firms to enter an industry.

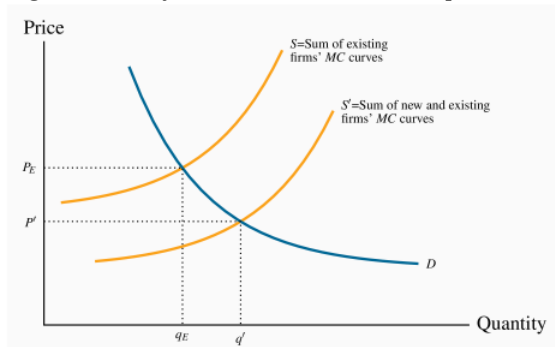
Let us return to our graphical analysis, and begin by supposing that the market equilibrium described in Figure 9.4 results in profits being made by some firms. Such an outcome is described in Figure 9.5, where the price exceeds the *ATC*. At the price P_E , a profit-making firm supplies the quantity q_E , as determined by its *MC* curve. On average, the cost of producing each unit of output, q_E , is defined by the point on the *ATC* at that output level, point k . Profit per unit is thus given by the value $(m-k)$ – the difference between revenue per unit and cost per unit. Total (economic) profit is therefore the area $P_E m k h$, which is quantity times profit per unit.

Figure 9.5 Short-run profits for the firm



At the price P_E , determined by the intersection of market demand and market supply, an individual firm produces the amount Q_E . The ATC of this output is k and therefore profit per unit is mk . Total profit is therefore $P_E mkh = 0q_E \times mk = TR - TC$.

Figure 9.6 Entry of firms due to economic profits



If economic profits result from the price P_E new firms enter the industry. This entry increases the market supply to S' and the equilibrium price falls to P' . Entry continues as long as economic profits are present. Eventually the price is driven to a level where only normal profits are made, and entry ceases.

While q_E represents an equilibrium for the firm, it is only a short-run, or temporary, equilibrium for the industry. The assumption of free entry and exit implies that the presence of economic profits will induce new entrepreneurs to enter and start producing. The impact of this dynamic is illustrated in Figure 9.6. An increased number of firms shifts supply rightwards to become S' , thereby increasing the amount supplied at any price. The impact on price of this supply shift is evident: With an unchanged demand, the equilibrium price must fall.

How far will the price fall, and how many new firms will enter this profitable industry? As long as economic profits exist new firms will enter and the resulting increase in supply will continue to drive the price downwards. But, once the price has been driven down to the minimum of the ATC of a representative firm, there is no longer an incentive for new entrepreneurs to enter. Therefore, the long-run industry equilibrium is where the market price equals the minimum point of a firm's ATC. This generates normal profits, and there is no incentive for firms to enter or exit.

A **long-run equilibrium** in a competitive industry requires a price equal to the minimum point of a firm's ATC. At this point, only normal profits exist, and there is no incentive for firms to enter or exit.

In developing this dynamic, we began with a situation in which economic profits were present. However, we could have equally started from a position of losses. With a market price between the minimum of the AVC and the minimum of the ATC in Figure 9.5, revenues per unit would exceed variable costs but not total costs per unit. When firms cannot cover their ATC in the long run, they will cease production. Such closures must reduce aggregate supply; consequently the market supply curve contracts, rather than expands as it did in Figure 9.6. The reduced supply drives up the price of the good. This process continues as long as firms are making losses. A final industry equilibrium is attained only when the price reaches a level where firms can make a normal profit. Again, this will be at the minimum of the typical firm's ATC.

Accordingly, the long-run equilibrium is the same, regardless of whether we begin from a position in which firms are incurring losses, or where they are making profits.

Application Box 9.2 Entry and exit: Oil rigs

Oil drilling is a competitive market. There are a large number of suppliers, information is ubiquitous, and entry and exit are relatively free.

In the years 2012 and 2013 the price of crude oil was around \$100 US per barrel. Towards the end of 2014 the price of oil began to drop on world markets, and by early 2015 it fluctuated around \$50. The response of drillers in the US was substantial and immediate. The number of active rigs declined dramatically. In terms of our economic model, certain suppliers exited; they moth-balled their rigs and waited for the price of oil to recover.

Another such cycle, even more pronounced, occurred in 2020. With the coronavirus pandemic, the demand for oil dropped and its price plummeted. Again, many firms shut down their rigs and had no choice but to sit out the price decline. A highly informative graphic is presented at <https://tradingeconomics.com/united-states/crude-oil-rigs>.

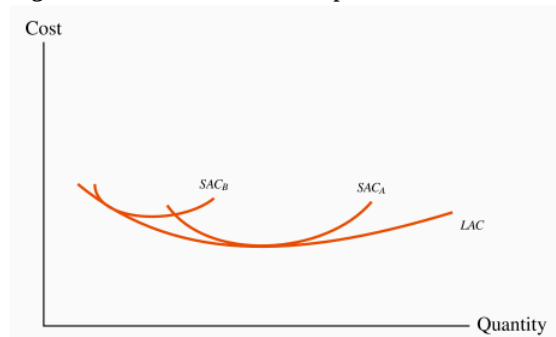
In addition to the decline in traditional oil recovery rigs, the number of operating shale crews declined by even greater amounts. Details at <https://www.forbes.com/sites/davidblackmon/2020/05/12/a-grim-earnings-season-for-the-us-shale-business/#6f55a95a1cf2>

This page titled [9.4: Dynamics- Entry and exit](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.5: Long-run industry supply

When aggregating the firm-level supply curves, as illustrated in Figure 9.3, we did not assume that all firms were identical. In that example, firm A has a cost structure with a lower AVC curve, since its supply curve starts at a lower dollar value. This indicates that firm A may have a larger plant size than firm B – one that puts A closer to the minimum efficient scale region of its long-run ATC curve.

Figure 9.7 Firms with different plant sizes



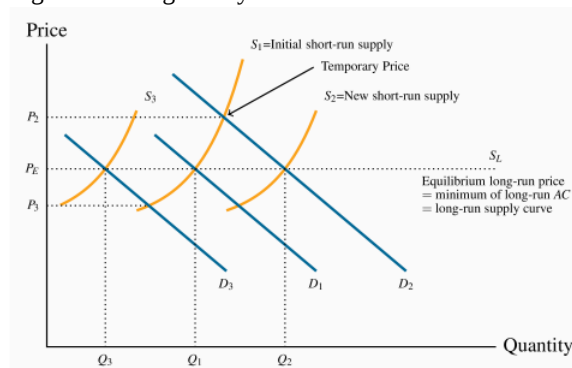
Firm B cannot compete with Firm A in the long run given that B has a less efficient plant size than firm A. The equilibrium long-run price equals the minimum of the LAC . At this price firm B must move to a more efficient plant size or make losses.

Can firm B survive with his current scale of operation in the long run? Our industry dynamics indicate that it cannot. The reason is that, provided *some* firms are making economic profits, new entrepreneurs will enter the industry and drive the price down to the minimum of the ATC curve of those firms who are operating with the lowest cost plant size. B-type firms will therefore be forced either to leave the industry or to adjust to the least-cost plant size—corresponding to the lowest point on its long-run ATC curve. Remember that the same technology is available to all firms; they each have the same long-run ATC curve, and may choose different scales of operation in the short run, as illustrated in Figure 9.7. But in the long run they must all produce using the minimum-cost plant size, or else they will be driven from the market.

This behaviour enables us to define a long-run industry supply. The long run involves the entry and exit of firms, and leads to a price corresponding to the minimum of the long-run ATC curve. Therefore, if the long-run equilibrium price corresponds to this minimum, the long-run supply curve of the industry is defined by a particular price value—it is horizontal at the price corresponding to the minimum of the $LATC$. More or less output is produced as a result of firms entering or leaving the industry, with those present always producing at the same unit cost in a long-run equilibrium.

Industry supply in the long run in perfect competition is horizontal at a price corresponding to the minimum of the representative firm's long-run ATC curve.

Figure 9.8 Long-run dynamics



The LR equilibrium price P_E is disturbed by a shift in demand from D_1 to D_2 . With a fixed number of firms, P_2 results. Profits accrue at this price and entry occurs. Therefore the SR supply shifts outwards until these profits are eroded and the new equilibrium output is Q_2 . If, instead, D falls to D_3 then firms exit because they make losses, S shifts back until the price is driven up sufficiently to restore normal profits. Different outputs are supplied in the long run at the same price P_E , therefore the long-run supply is horizontal at P_E .

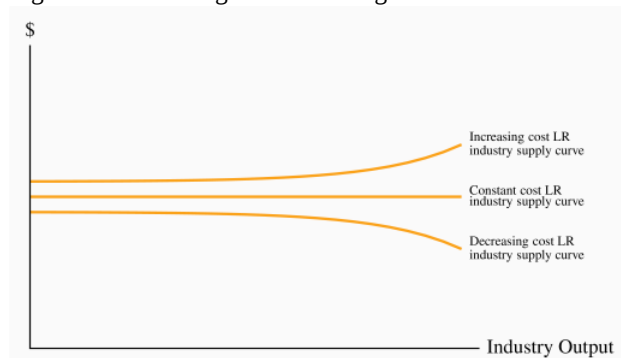
This industry's long-run supply curve, S_L , and a particular short-run supply are illustrated in Figure 9.8. Different points on S_L are attained when demand shifts. Suppose that, from an initial equilibrium Q_1 , defined by the intersection of D_1 and S_1 , demand increases from D_1 to D_2 because of a growth in income. With a fixed number of firms, the additional demand can be met only at a higher price (P_2), where each existing firm produces more using their existing plant size. The economic profits that result induce new operators to produce. This addition to the industry's production capacity shifts the short-run supply outwards and price declines until normal profits are once again being made. The new long-run equilibrium is at Q_2 , with more firms each producing at the minimum of their long-run ATC curve, P_E .

The same dynamic would describe the industry reaction to a decline in demand—price would fall, some firms would exit, and the resulting contraction in supply would force the price back up to the long-run equilibrium level. This is illustrated by a decline in demand from D_1 to D_3 .

Increasing and decreasing cost industries

While a horizontal long-run supply is the norm for perfect competition, in some industries costs increase with the scale of industry output; in others they decrease. This may be because all of the producers use a particular input that itself becomes more or less costly, depending upon the amount supplied.

Figure 9.9 Increasing and decreasing cost industries



When individual-supplier costs rise as the output of the industry increases we have an increasing cost supply curve for the industry in the long run. Conversely, when the costs of individual suppliers fall with the scale of the industry, we have a decreasing cost industry.

Decreasing cost sectors are those that benefit from a decline in the prices of their inputs as the size of their market expands. This is frequently because the suppliers of the inputs themselves can benefit from scale economies as a result of expansion in the market for the final good. A case in point has been the computer market, or the tablet market: As output in these markets has grown, the producers of videocards and random-access memory have benefited from scale economies and thus been able to sell these components at a lower price to the manufacturers of the final goods. An example of an increasing cost market is the market for landings and take-offs at airports. Airports are frequently limited in their ability to expand their size and build additional runways. In such markets, as use grows, planes about to land may have to adopt a circling holding pattern, while those departing encounter clearance delays. Such delays increase the time costs to passengers and the fuel and labour costs to the suppliers. Decreasing and increasing industry costs are reflected in the long-run industry supply curve by a downward-sloping segment or an upward sloping segment, as illustrated in Figure 9.9.

Increasing (decreasing) cost industry is one where costs rise (fall) for each firm because of the scale of industry operation.

This page titled [9.5: Long-run industry supply](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.6: Globalization and technological change

Globalization and technological change have had a profound impact on the way goods and services are produced and brought to market in the modern world. The cost structure of many firms has been reduced by *outsourcing to lower-wage economies*. Furthermore, the advent of the communications revolution has effectively *increased the minimum efficient scale for many industries*, as illustrated in Chapter 8 (Figure 8.7). Larger firms are less difficult to manage nowadays, and the *LAC* curve may not slope upwards until very high output levels are attained. The consequence is that some industries may not have sufficient "production space" to sustain a large number of firms. In order to reap the advantages of scale economies, firms become so large that they can supply a significant part of the market. They are no longer so small as to have no impact on the price.

Outsourcing and easier communications have in many cases simply eliminated many industries in the developed world. Garment making is an example. Some decades ago Quebec was Canada's main garment maker: Brokers dealt with 'cottage-type' garment assemblers outside Montreal and Quebec City. But ultimately the availability of cheaper labour in the developing world combined with efficient communications undercut the local manufacture. Most of Canada's garments are now imported. Other North American and European industries have been impacted in similar ways. Displaced labour has had to reskill, retool, reeducate itself, and either seek alternative employment in the manufacturing sector, or move to the service sector of the economy, or retire.

Globalization has had a third impact on the domestic economy, in so far as it *reduces the cost of components*. Even industries that continue to operate within national boundaries see a reduction in their cost structure on account of globalization's impact on input costs. This is particularly in evidence in the computing industry, where components are produced in numerous low-wage economies, imported to North America and assembled into computers domestically. Such components are termed *intermediate goods*.

This page titled [9.6: Globalization and technological change](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.7: Efficient resource allocation

Economists have a particular liking for competitive markets. The reason is not, as is frequently thought, that we love competitive battles; it really concerns resource allocation in the economy at large. In Chapter 5 we explained why markets are frequently an excellent vehicle for transporting the economy's resources to where they are most valued: A perfectly competitive marketplace in which there are no externalities results in resources being used up to the point where the demand and supply prices are equal. If demand is a measure of marginal benefit and supply is a measure of marginal cost, then a perfectly competitive market ensures that this condition will hold in equilibrium. *Perfect competition, therefore, results in resources being used efficiently.*

Our initial reaction to this perspective may be: If market equilibrium is such that the quantity supplied always equals the quantity demanded, is not every market efficient? The answer is no. As we shall see in the next chapter on monopoly, the monopolist's supply decision does not reflect the marginal cost of resources used in production, and therefore does not result in an efficient allocation in the economy.

This page titled [9.7: Efficient resource allocation](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.8: Key Terms

Perfect competition: an industry in which many suppliers, producing an identical product, face many buyers, and no one participant can influence the market.

Profit maximization is the goal of competitive suppliers – they seek to maximize the difference between revenues and costs.

Marginal revenue is the additional revenue accruing to the firm resulting from the sale of one more unit of output.

Shut-down price corresponds to the minimum value of the *AVC* curve.

Break-even price corresponds to the minimum of the *ATC* curve.

Short-run supply curve for perfect competitor: the portion of the *MC* curve above the minimum of the *AVC*.

Industry supply (short run) in perfect competition is the horizontal sum of all firms' supply curves.

Short-run equilibrium in perfect competition occurs when each firm maximizes profit by producing a quantity where $P=MC$.

Economic (supernormal) profits are those profits above normal profits that induce firms to enter an industry. Economic profits are based on the opportunity cost of the resources used in production.

Long-run equilibrium in a competitive industry requires a price equal to the minimum point of a firm's *ATC*. At this point, only normal profits exist, and there is no incentive for firms to enter or exit.

Industry supply in the long run in perfect competition is horizontal at a price corresponding to the minimum of the representative firm's long-run *ATC* curve.

Increasing (decreasing) cost industry is one where costs rise (fall) for each firm because of the scale of industry operation.

This page titled [9.8: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

9.9: Exercises for Chapter 9

EXERCISE 9.1

Wendy's Window Cleaning is a small local operation. Wendy presently cleans the outside windows in her neighbours' houses for \$36 per house. She does ten houses per day. She is incurring total costs of \$420, and of this amount \$100 is fixed. The cost per house is constant.

- What is the marginal cost associated with cleaning the windows of one house – we know it is constant?
- At a price of \$36, what is her break-even level of output (number of houses)?
- If the fixed cost is 'sunk' and she cannot increase her output in the short run, should she shut down?

EXERCISE 9.2

A manufacturer of vacuum cleaners incurs a constant variable cost of production equal to \$80. She can sell the appliances to a wholesaler for \$130. Her annual fixed costs are \$200,000. How many vacuums must she sell in order to cover her total costs?

EXERCISE 9.3

For the vacuum cleaner producer in Exercise 9.2:

- Draw the *MC* curve.
- Next, draw her *AFC* and her *AVC* curves.
- Finally, draw her *ATC* curve.
- In order for this cost structure to be compatible with a perfectly competitive industry, what must happen to her *MC* curve at some output level?

EXERCISE 9.4

Consider the supply curves of two firms in a competitive industry: $P=q_A$ and $P=2q_B$.

- On a diagram, draw these two supply curves, marking their intercepts and slopes numerically (remember that they are really *MC* curves).
- Now draw a supply curve that represents the combined supply of these two firms.

EXERCISE 9.5

Amanda's Apple Orchard Productions Limited produces 10,000 kilograms of apples per month. Her total production costs at this output level are \$8,000. Two of her many competitors have larger-scale operations and produce 12,000 and 15,000 kilos at total costs of \$9,500 and \$11,000 respectively. If this industry is competitive, on what segment of the *LAC* curve are these producers producing?

EXERCISE 9.6

Consider the data in the table below. *TC* is total cost, *TR* is total revenue, and *Q* is output.

Q	0	1	2	3	4	5	6	7	8	9	10
TC	10	18	24	31	39	48	58	69	82	100	120
TR	0	11	22	33	44	55	66	77	88	99	110

- Add some extra rows to the table and for each level of output calculate the *MR*, the *MC* and total profit.
- Next, compute *AFC*, *AVC*, and *ATC* for each output level, and draw these three cost curves on a diagram.
- What is the profit-maximizing output?
- How can you tell that this firm is in a competitive industry?

EXERCISE 9.7

Optional: The market demand and supply curves in a perfectly competitive industry are given by: $Q_d=30,000-600P$ and $Q_s=200P-2000$.

- a. Draw these functions on a diagram, and calculate the equilibrium price of output in this industry.
- b. Now assume that an additional firm is considering entering. This firm has a short-run MC curve defined by $MC=10+0.5q$, where q is the firm's output. If this firm enters the industry and it knows the equilibrium price in the industry, what output should it produce?

EXERCISE 9.8

Optional: Consider two firms in a perfectly competitive industry. They have the same MC curves and differ only in having higher and lower fixed costs. Suppose the ATC curves are of the form: $400/q+10+(1/4)q$ and $225/q+10+(1/4)q$. The MC for each is a straight line: $MC=10+(1/2)q$.

- a. In the first column of a spreadsheet enter quantity values of 1, 5, 10, 15, 20,..., 50. In the following columns compute the ATC curves for each quantity value.
- b. Compute the MC at each output in the next column, and plot all three curves.
- c. Compute the break-even price for each firm.
- d. Explain why both of these firms cannot continue to produce in the long run in a perfectly competitive market.

This page titled [9.9: Exercises for Chapter 9](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10: Monopoly

Chapter 10: Monopoly

In this chapter we will explore:

10.1	Why monopolies exist
10.2	How monopolists maximize profits
10.3	Long-run behaviour
10.4	Monopoly and market efficiency
10.5	Price discrimination
10.6	Cartels
10.7	Invention, innovation and rent seeking

10.1 Monopolies

In analyzing perfect competition we emphasized the difference between the industry and the individual supplier. The individual supplier is an atomistic unit with no market power. In contrast, a monopolist has a great deal of market power, for the simple reason that a monopolist is the sole supplier of a particular product and so really *is* the industry. The word monopoly, comes from the Greek words *monos*, meaning one, and *polein* meaning to sell. When there is just a single seller, our analysis need not distinguish between the industry and the individual firm. They are the same on the supply side.

Furthermore, the distinction between long run and short run is blurred, because a monopoly that continues to survive as a monopoly obviously sees no entry or exit. This is not to say that monopolized sectors of the economy do not evolve, they do. Sometimes they die, sometimes they evolve in a different role. For example, when digital cameras entered the market place in the eighties the Polaroid Land camera (which printed film straight out of the camera) 'died' because the demand side of the market lost interest. The Blackberry 'smart' phone had a virtual monopoly on this product into the new millennium until Apple and Nokia entered the market.

A **monopolist** is the sole supplier of an industry's output, and therefore the industry and the firm are one and the same.

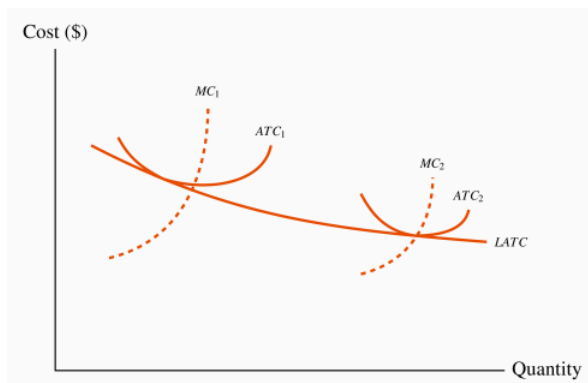
Monopolies can exist and exert their dominance in the market place for several reasons; scale economies, national policy, successful prevention of entry, research and development combined with patent protection.

Natural monopolies

Traditionally, monopolies were viewed as being 'natural' in some sectors of the economy. This means that scale economies define some industries' production and cost structures up to very high output levels, and that the whole market might be supplied at least cost by a single firm.

Consider the situation depicted in Figure 10.1. The long-run ATC curve declines indefinitely. There is no output level where average costs begin to increase. Imagine now having several firms, each producing with a plant size corresponding to the short-run average cost curve ATC_1 , or alternatively a single larger firm using a plant size denoted by ATC_2 . The small firms in this case cannot compete with the larger firm because the larger firm has lower production costs and can undercut the smaller firms, and *supply the complete market* in the process. Such a scenario is termed a natural monopoly.

Figure 10.1 A 'natural' monopolist



When LR average costs continue to decline at very high output, one large firm may be able to supply the industry at a lower unit cost than several smaller firms. With a plant size corresponding to ATC_2 , a single supplier can supply the whole market, whereas several smaller firms, each with plant size corresponding to ATC_1 , cannot compete with the larger firm on account of differential unit costs.

Natural monopoly: one where the ATC of producing any output declines with the scale of operation.

Electricity distribution in some of Canada's provinces is in the hands of a single supplier – *Hydro Quebec* or *Hydro One* in Ontario, for example. These distributors are natural monopolies in the sense described above: Unit distribution costs decline with size. In contrast, electricity production is not 'naturally' a monopoly. Other suppliers were once thought of as 'natural' monopolies also, but are no longer. *Bell Canada* was considered to be a natural monopoly in the era of land lines: It would not make economic sense to run several sets of phone lines to every residence. But that was before the arrival of cell phones, broadband and satellites. *Canada Post* was also thought to be a natural monopoly, until the advent of *FEDEX*, *UPS* and other couriers proved otherwise. Invention can compete away a 'natural' monopoly.

In reality there are very few *pure* monopolies. *Facebook*, *Microsoft*, *Amazon*, *Apple*, *Netflix* and *Google* may be extraordinarily dominant in their markets, but they are not the only suppliers of the services or products that they offer. There exist other products that are similar.

National and Provincial Policy

Government policy can foster monopolies. Some governments are, or once were, proud to have a 'national carrier' in the airline industry – *Air Canada* in Canada or *British Airways* in the UK. The mail service was viewed as a symbol of nationhood in Canada and the US: *Canada Post* and the *US Postal* system are national emblems that have historic significance. They were vehicles for integrating the provinces or states at various points in the federal lives of these countries.

In the modern era, most of Canada's provinces have decided to create a provincial monopoly crown corporation for the sale of cannabis. But competition abounds in the form of an illegal market.

The down side of such nationalist policies is that they can be costly to the taxpayer. Industries that are not subject to competition can become fat and uncompetitive: Managers have insufficient incentives to curtail costs; unions realize the government is committed to sustain the monopoly and push for higher wages than under a more competitive structure, and innovation may be less likely to occur.

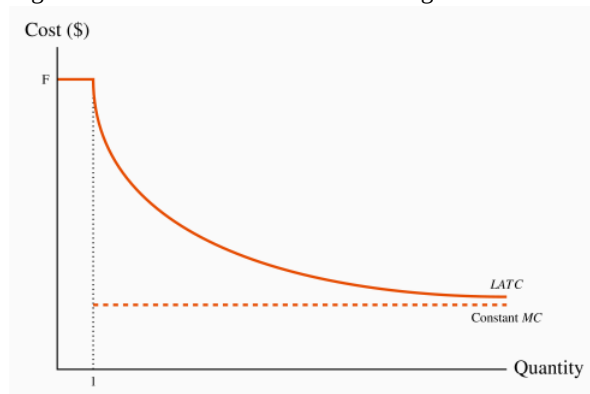
Maintaining barriers to entry

Monopolies can continue to survive if they are successful in preventing the entry of new firms and products. *Patents and copyrights* are one vehicle for preserving the sole-supplier role, and are certainly necessary to encourage firms to undertake the research and development (R&D) for new products.

Many corporations produce products that require a large up-front investment; this might be in the form of research and development, or the construction of costly production facilities. For example, *Boeing* or *Airbus* incurs billions of dollars in developing new aircraft; pharmaceuticals may have to invest a billion dollars to develop a new drug. However, once such an investment is complete, the cost of producing each unit of output may be relatively low. This is particularly true in pharmaceuticals. Such a phenomenon is displayed in Figure 10.2. In this case the average cost for a small number of units produced is high, but once the fixed cost is spread over an ever larger output, the average cost declines rapidly, and in the limit approaches the marginal cost.

These production structures are common in today's global economy, and they give rise to markets characterized either by a single supplier or a small number of suppliers.

Figure 10.2 Fixed cost and constant marginal cost



With a fixed cost of producing the first unit of output equal to F and a constant marginal cost thereafter, the long-run average total cost, $LATC$, declines indefinitely and becomes asymptotic to the marginal cost curve.

This figure is useful in understanding the role of patents. Suppose that *Pharma A* spends one billion dollars in developing a new drug and has constant unit production costs thereafter, while *Pharma B* avoids research and development and simply imitates *Pharma A*'s product. Clearly *Pharma B* would have a $LATC$ equal to its LMC , and would be able to undercut the initial developer of the drug. Such an outcome would discourage investment in new products and the economy at large would suffer as a consequence. Economies would be worse off if protection is not provided to the developers of new products because, if such protection is not offered, potential developers will not have the incentive to incur the up-front investment required.

While copyright and patent protection is legal, *predatory pricing* is an illegal form of entry barrier, and we explore it more fully in Chapter 14. An example would be where an existing firm that sells nationally may deliberately undercut the price of a small local entrant to the industry. Airlines with a national scope are frequently accused of posting low fares on flights in regional markets that a new carrier is trying to enter.

Political lobbying is another means of maintaining monopolistic power. For example, the *Canadian Wheat Board* had fought successfully for decades to prevent independent farmers from marketing wheat. This Board lost its monopoly status in August 2012, when the government of the day decided it was not beneficial to consumers or farmers in general. Numerous 'supply management' policies are in operation all across Canada. Agriculture is protected by production quotas. All maple syrup in Quebec must be marketed through a single monopoly supplier.

Critical networks also form a type of barrier, though not always a monopoly. *Microsoft's Office* package has an almost monopoly status in word processing and spreadsheet analysis for the reason that so many individuals and corporations use it. The fact that we know a business colleague will be able to edit our documents if written in *Word*, provides us with an incentive to use *Word*, even if we might prefer *Wordperfect* as a vehicle for composing documents. We develop the concept of strategic entry prevention further in Chapter 11.

10.2 Profit maximizing behaviour

We established in the previous chapter that, in deciding upon a profit-maximizing output, any firm should produce up to the point where the additional cost equals the additional revenue from a unit of output. What distinguishes the supply decision for a monopolist from the supply decision of the perfect competitor is that the monopolist faces a downward sloping demand. A monopolist is the sole supplier and therefore must meet the full market demand. This means that if more output is produced, the price must fall. We will illustrate the choice of a profit maximizing output using first a marginal-cost/marginal-revenue approach; then a supply/demand approach.

Marginal revenue and marginal cost

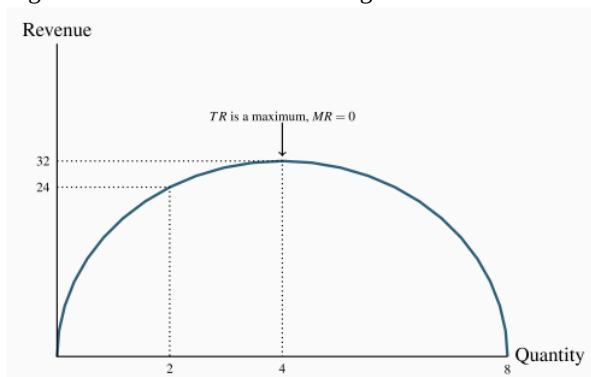
Table 10.1 displays price and quantity values for a demand curve in columns 1 and 2. Column 3 contains the sales revenue generated at each output. It is the product of price and quantity. Since the price denotes the revenue per unit, it is sometimes referred to as average revenue. The total revenue (TR) reaches a maximum at \$32, where 4 units of output are produced. A greater output necessitates a lower price on every unit sold, and in this case revenue falls if the fifth unit is brought to the market. Even

though the fifth unit sells for a positive price, the price on the other 4 units is now lower and the net effect is to reduce total revenue. This pattern reflects what we examined in Chapter 4: As price is lowered from the highest possible value of \$14 (where 1 unit is demanded) and the corresponding quantity increases, revenue rises, peaks, and ultimately falls as output increases. In Chapter 4 we explained that this maximum revenue point occurs where the price elasticity is unity (-1), at the midpoint of a linear demand curve.

Table 10.1 A profit maximizing monopolist

Quantity (Q)	Price (P)	Total revenue (TR)	Marginal revenue (MR)	Marginal cost (MC)	Total cost (TC)	Profit
0	16					
1	14	14	14	2	2	12
2	12	24	10	3	5	19
3	10	30	6	4	9	21
4	8	32	2	5	14	18
5	6	30	-2	6	20	10
6	4	24	-6	7	27	-3
7	2	14	-10	8	35	-21

Figure 10.3 Total revenue and marginal revenue



When the quantity sold increases total revenue/expenditure initially increases also. At a certain point, further sales require a price that not only increases quantity, but reduces revenue on units already being sold to such a degree that *TR* declines – where the demand elasticity equals -1 (the mid point of a linear demand curve). Here the midpoint occurs at $Q=4$. Where the *TR* is a maximum the $MR=0$.

Related to the total revenue function is the marginal revenue function. It is the addition to total revenue due to the sale of one more unit of the commodity.

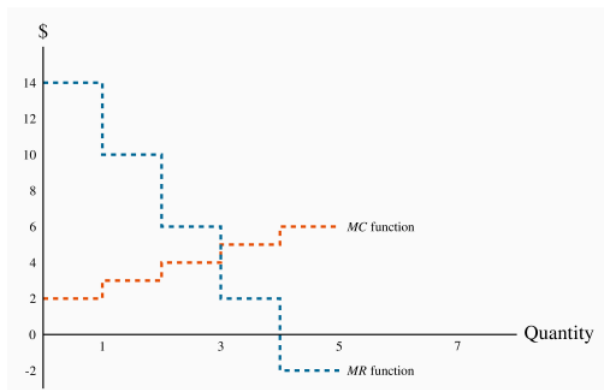
Marginal revenue is the change in total revenue due to selling one more unit of the good.

Average revenue is the price per unit sold.

The *MR* in this example is defined in the fourth column of Table 10.1. When the quantity sold increases from 1 unit to 2 units total revenue increases from \$14 to \$24. Therefore the marginal revenue associated with the second unit of output is \$10. When a third unit is sold *TR* increases to \$30 and therefore the *MR* of the third unit is \$6. As output increases the *MR* declines and eventually becomes negative – at the point where the *TR* is a maximum: If *TR* begins to decline then the additional revenue is by definition negative.

The *MR* function is plotted in Figure 10.4. It becomes negative when output increases from 4 to 5 units.

Figure 10.4 Monopolist's profit maximizing output



It is optimal for the monopolist to increase output as long as MR exceeds MC . In this case $MR > MC$ for units 1, 2 and 3. But for the fourth unit $MC > MR$ and therefore the monopolist would reduce total profit by producing it. He should produce only 3 units of output.

The optimal output

This producer has a marginal cost structure given in the fifth column of the table, and this too is plotted in Figure 10.4. Our profit maximizing rule from Chapter 8 states that it is optimal to produce a greater output as long as the additional revenue exceeds the additional cost of production on the next unit of output. In perfectly competitive markets the additional revenue is given by the fixed price for the individual producer, whereas for the monopolist the additional revenue is the marginal revenue. Consequently as long as MR exceeds MC for the next unit a greater output is profitable, but once MC exceeds MR the production of additional units should cease.

From Table 10.1 and Figure 10.4 it is clear that the optimal output is at 3 units. The third unit itself yields a profit of \$2, the difference between MR (\$6) and MC (\$4). A fourth unit however would reduce profit by \$3, because the MR (\$2) is less than the MC (\$5). What price should the producer charge? The price, as always, is given by the demand function. At a quantity sold of 3 units, the corresponding price is \$10, yielding total revenue of \$30.

Profit is the difference between total revenue and total cost. In Chapter 8 we computed total cost as the average cost times the number of units produced. It can also be computed as the sum of costs associated with each unit produced: The first unit costs \$2, the second \$3 and the third \$4. The total cost of producing 3 units is the sum of these dollar values: $\$9 = \$2 + \$3 + \4 . The profit-maximizing output therefore yields a profit of \$21 ($\$30 - \9).

Supply and demand

When illustrating market behaviour it is convenient to describe behaviour by simple linear supply and demand functions that are continuous, rather than the 'step' functions used in the preceding example. As explained in Chapter 5, in using continuous curves to represent a market we implicitly assume that a unit of output can be broken into subunits. In the example above we assumed that sales always involve one whole unit of the product being sold. In fact many goods can be sold in fractional units: Gasoline can be sold in fractions of a litre; fruits and vegetables can be sold in fractions of a kilogram, and so forth. Table 10.2 below furnishes the data for our analysis.

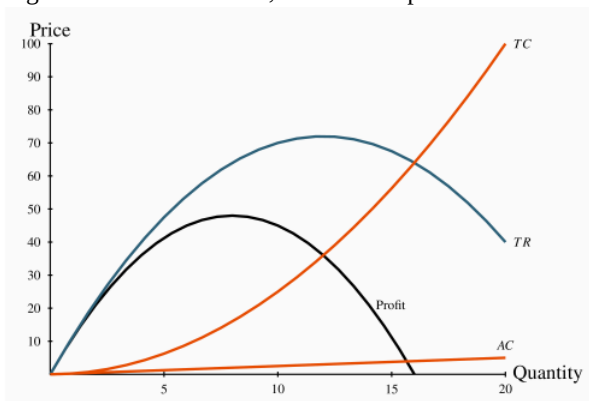
Table 10.2 Discrete quantities

Price	Quantity	Total	Total	Profit
	demand	revenue	cost	
12	0	0	0	0
11	2	22	1	21
10	4	40	4	36
9	6	54	9	45
8	8	64	16	48
7	10	70	25	45

6	12	72	36	36
5	14	70	49	21
4	16	64	64	0
3	18	54	81	-27
2	20	40	100	-60
1	22	22	121	-99
0	24	0	144	-144

The first two columns define the demand curve. Total revenue is the product of price and quantity and given in column 3. The cost data are given in column 4, and profit – the difference between total revenue and total cost is in the final column. Profit is maximized where the difference between revenue and cost is greatest; in this case where the output is 8 units. At lower or higher outputs profit is less. Figure 10.5 contains the curves defining total revenue (TR), total cost (TC) and profit. These functions can be obtained by mapping all of the revenue-quantity combinations, the cost-quantity combinations, and the profit-quantity combinations as a series of points, and joining these points to form the smooth functions displayed. The vertical axis is measured in dollars, the horizontal axis in units of output. Graphically, profit is maximized where the dollar difference between TR and TC is greatest; that is at the output where the vertical distance between the two curves is greatest. This difference, which is also defined by the profit curve, occurs at a value of 8 units, corresponding to the outcome in Table 10.2.

Figure 10.5 Total revenue, total cost & profit



At any quantity less than this output, profit would rise with additional output. This is because, from a less-than-optimal output, the additional revenue from increased sales exceeds the increased cost associated with producing those units: Stated differently, the marginal revenue would exceed the marginal cost. Conversely, outputs greater than the optimum result in a MR less than the associated MC . Accordingly, since outputs where $MR > MC$ are too low, and outputs where $MR < MC$ are too high, the optimum must be where the $MR = MC$. Hence, the equality between MR and MC is implied in this diagram at the output where the difference between TR and TC is greatest.

Note finally that total revenue is maximized where the TR curve reaches a peak. In this example that occurs at a value of 12 units of output. This is to be anticipated, as we learned in Chapter 4, because the midpoint of the demand schedule in Table 10.2 occurs at that value.

Figure 10.6 Market demand, the MR curve, and the monopolist's AC and MC curves

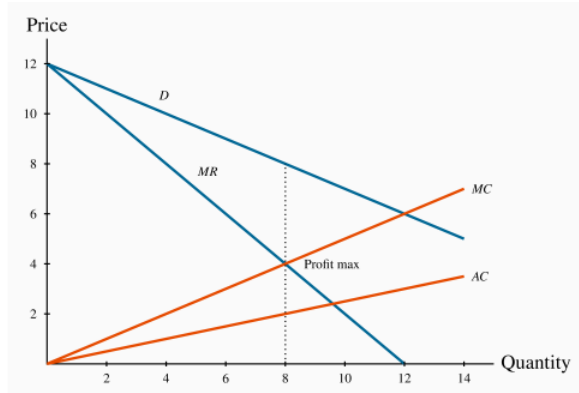


Figure 10.6 displays the demand curve for the market, the MR curve, and the monopolist's MC and AC curves. Consider first the marginal revenue curve. In contrast to the previous example, where only whole or integer units could be sold, in this example units can be sold in fractional amounts, and the MR curve must reflect this. To determine the position of the MR curve, note that with a straight-line demand curve total revenue is a maximum at the midpoint of the demand curve. Any increase in output results in reduced revenue: Stated differently, the marginal revenue becomes negative at that output. Up to that output the MR is positive, as illustrated in Figure 10.3. Accordingly, the MR curve must intersect the quantity axis midway between zero and the horizontal-axis intercept of the demand curve. Geometrically, since the MR intersects the quantity axis half way to the horizontal intercept of the demand curve, it must have a slope that is twice the slope of the demand curve.

By observing the data in columns 1 and 2 of the table, the demand curve intercepts are $\{ \$12, 24 \}$, and from above discussion the MR curve has intercepts $\{ \$12, 12 \}$. The AC is obtained by dividing TC by output in Table 10.2, and the MC can be also calculated as the change in total cost divided by the change in output from Table 10.2. The result of these calculations is displayed in Figure 10.6.

The profit maximizing output is 8 units, where $MC=MR$. The price at which 8 units can be sold is read from the demand curve¹, or the first column in Table 10.2. It is \$8. And, as expected, this price-quantity combination maximizes profit. Table 10.2 indicates that profit is maximized at \$48, at $q=8$.

Demand elasticity and marginal revenue

We have shown above that the MR curve cuts the horizontal axis at a quantity where the elasticity of demand is unity. We know from Chapter 4 that demand is elastic at points on the demand curve above this unit-elastic point. Furthermore, since the intersection of MR and MC must be at a positive dollar value (MC cannot be negative), then it must be the case that the *profit maximizing price for a monopolist always lies on the elastic segment of the demand curve*.

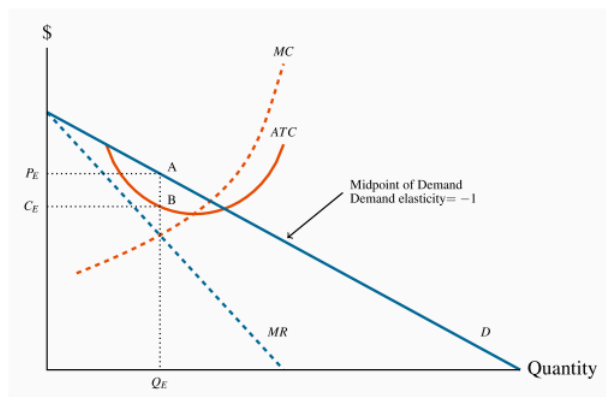
A general graphical representation

In Figure 10.7 we generalize the graphical representation of the monopoly profit maximizing output by allowing the MC and ATC curves to be nonlinear. The optimal output is at Q_E , where $MR=MC$, and the price P_E sustains that output. With the average cost known, profit per unit is AB , and therefore total profit is this margin multiplied by the number of units sold, Q_E .

Total profit is therefore $P_E ABC_E$.

Note that the monopolist may not always make a profit. Losses could result in Figure 10.7 if average costs were to rise so that the ATC were everywhere above the demand curve, or if the demand curve shifted down to being everywhere below the ATC curve. In the longer term the monopolist would have to either reduce costs or perhaps stimulate demand through advertising if she wanted to continue in operation.

Figure 10.7 The monopoly equilibrium



The profit maximizing output is Q_E , where $MC=MR$. This output can be sold at a price P_E . The cost per unit of Q_E is read from the ATC curve, and equals B . Per unit profit is therefore AB and total profit is P_EABC_E .

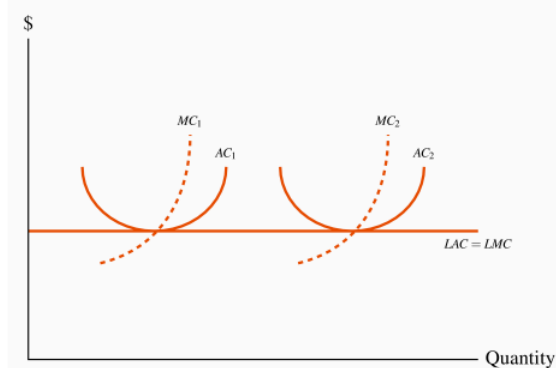
10.3 Long-run choices

Consider next the impact of a shift in demand upon the profit maximizing choice of this firm. A rightward shift in demand in Figure 10.7 also yields a new MR curve. The firm therefore chooses a new level of output, using the same profit maximizing rule: Set $MC=MR$. This output will be greater than the previous output, but again the price must be on an elastic portion of the new demand curve. If operating with the same plant size, the MC and ATC curves do not change and the new profit per unit is again read from the ATC curve.

By this stage the curious student will have asked: "What happens to plant size in the long run?" For example, is the monopolist in Figure 10.7 using the most appropriate plant size in the first place? Even if she is, should the monopolist consider adopting an expanded plant size in response to the shift in demand?

The answer is: In the long run the monopolist is free to choose whatever plant size is best. Her initial plant size might have been optimal for the demand she faced, but if it was, it is unlikely to be optimal for the larger scale of production associated with the demand shift. Accordingly, with the new demand curve, she must consider how much profit she could make using different plant sizes.

Figure 10.8 The monopolist's choice of plant size



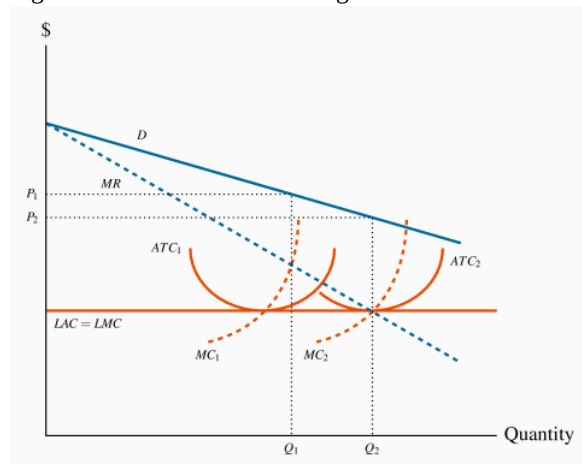
With constant returns to scale and constant prices per unit of labour and capital, a doubling of output involves exactly a doubling of costs. Thus, per unit costs, or average costs, are constant in the LR. Hence $LAC=LMC$, and each is constant.

To illustrate one possibility, we will think of this firm as having constant returns to scale at all output ranges, as displayed in Figure 10.8. (Our reasoning carries through if the LAC slopes downwards; the graph just becomes a little more complex.) The key characteristic of constant returns to scale is that a doubling of inputs leads to a doubling of output. Therefore, if the per-unit cost of inputs is fixed, a doubling of inputs (and therefore output) leads exactly to a doubling of costs. This implies that, when the firm varies its plant size *and* its labour use, the cost of producing each additional unit must be constant. The long-run marginal cost LMC is therefore constant and equals the ATC in the long run.

Figure 10.9 describes the market for this good. The optimal output and price are determined in the usual manner: Set $MC=MR$. If the monopolist has plant size corresponding to ATC_1 , the optimal output is Q_1 and should be sold at the price P_1 . The key issue now

is: Given the demand conditions, could the monopolist make more profit by choosing a plant size that differs from the one corresponding to ATC_1 ?

Figure 10.9 Plant size in the long run



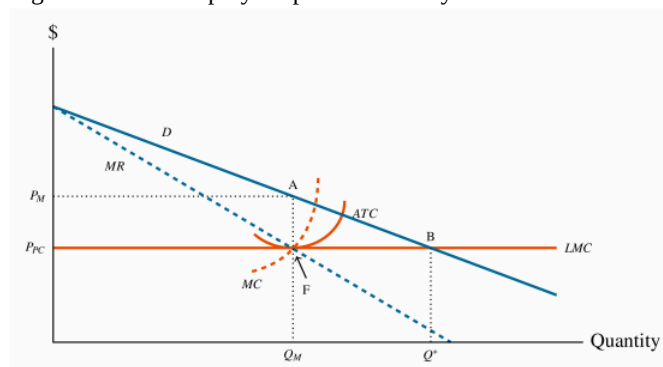
With demand conditions defined by D and MR , the optimal plant size is one corresponding to the point where $MR=MC$ in the long run. Therefore Q_2 is the optimal output and the optimal plant size corresponds to ATC_2 . If the current plant is defined by ATC_1 , then optimal SR production is Q_1 .

In this instance the answer is a clear 'yes'. Her LMC curve is horizontal and so, by increasing output from Q_1 to Q_2 she earns a profit on each additional unit in that range, because the MR curve lies above the LMC curve. In order to produce the output level Q_2 at least cost she must choose a plant size corresponding to AC_2 .

10.4 Output inefficiency

A characteristic of perfect competition is that it secures an efficient allocation of resources when there are no externalities in the market: Resources are used up to the point where their marginal cost equals their marginal value – as measured by the price that consumers are willing to pay. But a monopoly structure does not yield this output. Consider Figure 10.10.

Figure 10.10 Monopoly output inefficiency



A monopolist maximizes profit at Q_M . Here the value of marginal output exceeds cost. If output expands to Q^* a gain arises equal to the area ABF . This is the deadweight loss associated with the output Q_M rather than Q^* . If the monopolist's long-run MC is equivalent to a competitive industry's supply curve, then the deadweight loss is the cost of having a monopoly rather than a perfectly competitive market.

The monopolist's profit-maximizing output Q_M is where MC equals MR . This output is inefficient for the reason that we developed in Chapter 5: If output is increased beyond Q_M the additional benefit exceeds the additional cost of producing it. The additional benefit is measured by the willingness of buyers to pay – the market demand curve. The additional cost is the long-run MC curve under the assumption of constant returns to scale. Using the terminology from Chapter 5, there is a deadweight loss equal to the area ABF . This is termed allocative inefficiency.

Allocative inefficiency arises when resources are not appropriately allocated and result in deadweight losses .

Perfect competition versus monopoly

The area ABF can also be considered as the efficiency loss associated with having a monopoly rather than a perfectly competitive market structure. In perfect competition the supply curve is horizontal. This is achieved by having firms enter and exit when more or less must be produced. Accordingly, *if the perfectly competitive industry's supply curve approximates the monopolist's long-run marginal cost curve*², we can say that if the monopoly were turned into a competitive industry, output would increase from Q_M to Q^* . The deadweight loss is one measure of the superiority of the perfectly competitive structure over the monopoly structure.

Note that this critique of monopoly is not initially focused upon profit. While monopoly profits are what frequently irk the public, we have focused upon resource allocation inefficiencies. But in a real sense the two are related: Monopoly inefficiencies arise through output being restricted, and it is this output reduction – achieved by maintaining a higher than competitive price – that gives rise to those profits. Nonetheless, there is more than just a shift in purchasing power from the buyer to the seller. Deadweight losses arise because output is at a level lower than the point where the MC equals the value placed on the good; thus the economy is sacrificing the possibility of creating additional surplus.

Given that monopoly has this undesirable inefficiency, what measures should be taken, if any, to counter the inefficiency? We will see what Canada's Competition Act has to say in Chapter 14 and also examine what other measures are available to control monopolies.

10.5 Price discrimination

A common characteristic in the pricing of many goods is that different individuals pay different prices for goods or services that are essentially the same. Examples abound: Seniors get a reduced rate for coffee in *Burger King*; hair salons charge women more than they charge men; bank charges are frequently waived for juniors. Price discrimination involves charging different prices to different consumers in order to increase profit.

Price discrimination involves charging different prices to different consumers in order to increase profit.

A strict definition of discrimination involves different prices for *identical products*. We all know of a school friend who has been willing to take the midnight flight to make it home at school break at a price he can afford. In contrast, the business executive prefers the seven a.m. flight to arrive for a nine a.m. business meeting in the same city at several times the price. These are very mild forms of price discrimination, since a midnight flight (or a midday flight) is not a perfect substitute for an early morning flight. Price discrimination is practiced because buyers are willing to pay different amounts for a good or service, and the supplier may have a means of profiting from this. Consider the following example.

Family Flicks is the local movie theatre. It has two distinct groups of customers – those of prime age form one group; youth and seniors form the other. Family Flicks has done its market research and determined that each group accounts for 50 percent of the total market of 100 potential viewers per screening. It has also established that the prime-age group members are willing to pay \$12 to see a movie, while the seniors and youth are willing to pay just \$5. How should the tickets be priced?

Family Flicks has no variable costs, only fixed costs. It must pay a \$100 royalty to the movie maker each time it shows the current movie, and must pay a cashier and usher \$20 each. Total costs are therefore \$140, regardless of how many people show up – short-run MC is zero. On the pricing front, as illustrated in Table 10.3 below, if Family Flicks charges \$12 per ticket it will attract 50 viewers, generate \$600 in revenue and therefore make a profit of \$460.

Table 10.3 Price discrimination

	$P=\$5$	$P=\$12$	Twin price
No. of customers	100	50	
Total revenue	\$500	\$600	\$850
Total costs	\$140	\$140	\$140
Profit	\$360	\$460	\$710

In contrast, if it charges \$5 it can fill the theatre, because each of the prime-age individuals is willing to pay more than \$5, but the seniors and youth are now offered a price they too are willing to pay. However, the total revenue is now only \$500 ($100 \times \$5 = \500), and profits are reduced to \$360. It therefore decides to charge the high price and leave the theatre half-empty, because this strategy maximizes its profit.

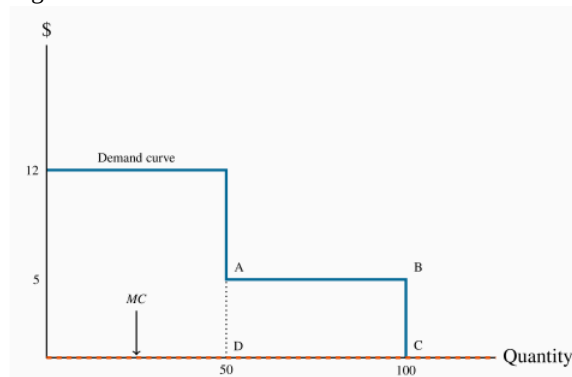
Suppose finally that the theatre is able to segregate its customers. It can ask the young and senior customers for identification upon entry, and in this way charge them a lower price, *while still maintaining the higher price to the prime-age customers*. If it can execute such a plan Family Flicks can now generate \$850 in revenue – \$600 from the prime-age group and \$250 from the youth and seniors groups. Profit soars to \$710.

There are two important conditions for this scheme to work:

1. The seller must be able to *segregate the market* at a reasonable cost. In the movie case this is achieved by asking for identification.
2. The second condition is that *resale must be impossible or impractical*. For example, we rule out the opportunity for young buyers to resell their tickets to the prime-age individuals. Sellers have many ways of achieving this – they can require immediate entry to the movie theatre upon ticket purchase, they can stamp the customer's hand, they can demand the showing of ID with the ticket when entering the theatre area.

Frequently we think of sellers who offer price reductions to specific groups as being generous. For example, hotels may levy only a nominal fee for the presence of a child, once the parents have paid a suitable rate for the room or suite in which a family stays. The hotel knows that if it charges too much for the child, it may lose the whole family as a paying unit. The coffee shop offering cheap coffee to seniors is interested in getting a price that will cover its variable cost and so contribute to its profit. It is unlikely to be motivated by philanthropy, or to be concerned with the financial circumstances of seniors.

Figure 10.11 Price discrimination at the movies



At $P=12$, 50 prime-age individuals demand movie tickets. At $P=5$, 50 more seniors and youths demand tickets. Since the MC is zero the efficient output is where the demand curve takes a zero value – where all 100 customers purchase tickets. Thus, any scheme that results in all 100 individuals buying ticket is efficient. Efficient output is at point C.

Price discrimination has a further interesting feature that is illustrated in Figure 10.11: It frequently *reduces the deadweight loss* associated with a monopoly seller!

In our Family Flicks example, the profit maximizing monopolist that did not, or could not, price discriminate *left 50 customers unsupplied who were willing to pay \$5 for a good that had a zero MC*. This is a deadweight loss of \$250 because 50 seniors and youth valued a commodity at \$5 that had a zero MC . Their demand was not met because, in the absence of an ability to discriminate between consumer groups, Family Flicks made more profit by satisfying the demand of the prime-age group alone. But in this example, by segregating its customers, the firm's profit maximization behaviour resulted in the DWL being eliminated, because it supplied the product to those additional 50 individuals. In this instance *price discrimination improves welfare*, because more of a good is supplied in a situation where market valuation exceeds marginal cost.

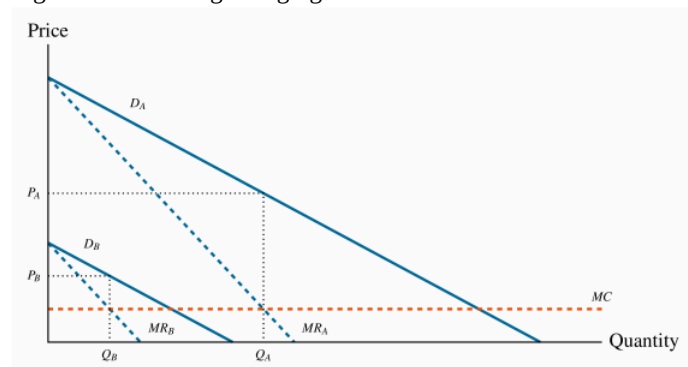
In the preceding example we simplified the demand side of the market by assuming that every individual in a given group was willing to pay the same price – either \$12 or \$5. More realistically each group can be defined by a downward-sloping demand curve, reflecting the variety of prices that buyers in a given market segment are willing to pay. It is valuable to extend the analysis to include this reality. For example, a supplier may face different demands from her domestic and foreign buyers, and if she can segment these markets she can price discriminate effectively.

Consider Figure 10.12 where two segmented demands are displayed, D_A and D_B , with their associated marginal revenue curves, MR_A and MR_B . We will assume that marginal costs are constant for the moment. It should be clear by this point that the profit maximizing solution for the monopoly supplier is to supply an amount to each market where the MC equals the MR in each market: Since the buyers in one market cannot resell to buyers in the other, the monopolist considers these as two different markets and

therefore maximizes profit by applying the standard rule. She will maximize profit in market A by supplying the quantity Q_A and in market B by supplying Q_B . The prices at which these quantities can be sold are P_A and P_B . These prices, unsurprisingly, are different – the objective of segmenting markets is to increase profit by treating the markets as distinct.

An example of this type of price discrimination is where pharmaceutical companies sell drugs to less developed economies at a lower price than to developed economies. The low price is sufficient to cover marginal cost and is therefore profitable - provided the high price market covers the fixed costs.

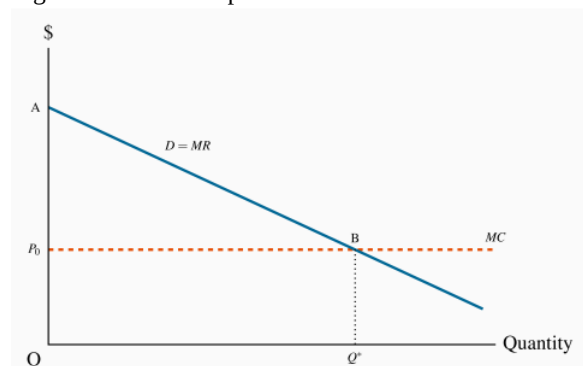
Figure 10.12 Pricing in segregated markets



With two separate markets defined by D_A and D_B , and their associated MR_A and MR_B , a profit maximizing strategy is to produce where $MC=MR_A=MR_B$, and discriminate between the two markets by charging prices P_A and P_B .

The preceding examples involved two separable groups of customers and are very real. This kind of group segregation is sometimes called *third degree price discrimination*. But it may be possible to segregate customers into several groups rather than just two. In the limit, if we could charge a different price to every consumer in a market, or for every unit sold, the revenue accruing to the monopolist would be the area under the demand curve up to the output sold. Though primarily of theoretical interest, this is illustrated in Figure 10.13. It is termed *perfect price discrimination*, and sometimes *first degree price discrimination*. Such discrimination is not so unrealistic: A tax accountant may charge different customers a different price for providing the same service; home renovators may try to charge as much as any client appears willing to pay.

Figure 10.13 Perfect price discrimination



A monopolist who can sell each unit at a different price maximizes profit by producing Q^* . With each consumer paying a different price the demand curve becomes the MR curve. The result is that the monopoly DWL is eliminated because the efficient output is produced, and the monopolist appropriates all the consumer surplus. Total revenue for the perfect price discriminator is $OABQ^*$.

Second degree price discrimination is based on a different concept of buyer identifiability. In the cases we have developed above, the seller is able to distinguish the buyers by *observing* a vital characteristic that signals their type. It is also possible that, while individuals might have defining traits which influence their demands, such traits might not be detectable by the supplier. Nonetheless, it is frequently possible for the supplier to offer different pricing options (corresponding to different uses of a product) that buyers would choose from, with the result that her profit would be greater than under a uniform price with no variation in the use of the service. Different cell phone 'plans', or different internet plans that users can choose from are examples of this second-degree discrimination.

10.6 Cartels: Acting like a monopolist

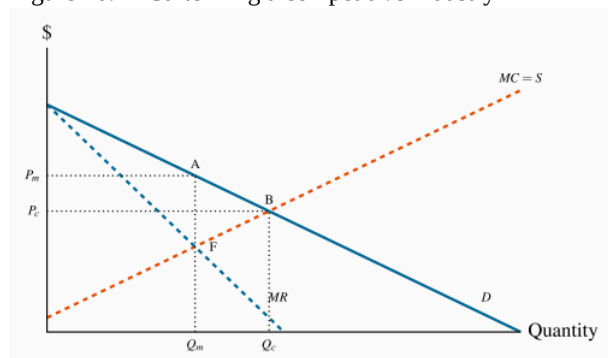
A cartel is a group of suppliers that colludes to operate like a monopolist. The cartel formed by the members of the Organization of Oil Exporting Countries (OPEC) is an example of a cartel that was successful in achieving its objectives for a long period. This cartel first flexed its muscles in 1973, by increasing the world price of oil from \$3 per barrel to \$10 per barrel. The result was to transfer billions of dollars from the energy-importing nations in Europe and North America to OPEC members – the demand for oil is relatively inelastic, hence an increase in price increases total expenditures.

A **cartel** is a group of suppliers that colludes to operate like a monopolist.

A second renowned cartel is managed by De Beers, which controls a large part of the world's diamond supply. In Canada, agricultural marketing boards are a means of restricting supply legally. Such cartels may have thousands of members. By limiting entry, through requiring a production 'quota', the incumbents can charge a higher price than if entry to the industry were free.

To illustrate the dynamics of cartels consider Figure 10.14. Several producers, with given production capacities, come together and agree to restrict output with a view to increasing price and therefore profit. This may be done with the agreement of the government, or it may be done secretly, and possibly against the law. Each firm has a MC curve, and the industry supply is defined as the sum of these marginal cost curves, as illustrated in Figure 9.3. The resulting cartel is effectively one in which there is a single supplier with many different plants – a multi-plant monopolist. To maximize profits this organization will choose an output level Q_m where the MR equals the MC . In contrast, if these firms act competitively the output chosen will be Q_c . The competitive output yields no supernormal profit, whereas the monopoly/cartel output does.

Figure 10.14 Cartelizing a competitive industry



A cartel is formed when individual suppliers come together and act like a monopolist in order to increase profit. If MC is the joint supply curve of the cartel, profits are maximized at the output Q_m , where $MC=MR$. In contrast, if these firms operate competitively output increases to Q_c .

The cartel results in a deadweight loss equal to the area ABE , just as in the standard monopoly model.

Cartel instability

Some cartels are unstable in the long run. In the first instance, the degree of instability depends on the authority that the governing body of the cartel can exercise over its members, and upon the degree of information it has on the operations of its members. If a cartel is simply an arrangement among producers to limit output, each individual member of the cartel has an incentive to increase its output, because the monopoly price that the cartel attempts to sustain exceeds the cost of producing a marginal unit of output. In Figure 10.14 each firm has a MC of output equal to $\$F$ when the group collectively produces the output Q_m . Yet any firm that brings output to market, *beyond its agreed production limit*, at the price P_m will make a profit of AF on that additional output – *provided the other members of the cartel agree to restrict their output*. Since each firm faces the same incentive to increase output, it is difficult to restrain all members from doing so.

Individual members are more likely to abide by the cartel rules if the organization can sanction them for breaking the supply-restriction agreement. Alternatively, if the actions of individual members are not observable by the organization, then the incentive to break ranks may be too strong for the cartel to sustain its monopoly power.

We will see in Chapter 14 that Canada's Competition Act forbids the formation of cartels, as it forbids many other anti-competitive practices. At the same time, our governments frequently are the driving force in the formation of domestic cartels.

In the second instance, cartels may be undermined eventually by the emergence of new products and new technologies. OPEC has lost much of its power in the modern era because of technological developments in oil recovery. Canada's 'tar sands' yield oil, as a

result of technological developments that enabled producers to separate the oil from the earth it is mixed with. Fracking technologies are another means of extracting oil that is discovered in small pockets and encased in rock. The supply coming from these new technologies has limited the ability of the old OPEC cartel to increase prices through supply restriction.

Application Box 10.1 The taxi cartel

The new sharing economy has brought competition to some traditional cartels. City taxis are an example of such a formation: Traditionally, entry has been restricted to drivers who hold a permit (medallion), and fares are higher as a consequence of the resulting reduced supply. A secondary market then develops for these medallions, in which the city may offer new medallions through auction, or existing owners may exit and sell their medallions. Restricted entry has characterized most of Canada's major cities. Depending on the strictness of the entry process, medallions are worth correspondingly more. By 2012, medallions were selling in New York and Boston for a price in the neighborhood of one million dollars.

But ride-sharing start-up companies changed all of that. As Western examples, Uber and Lyft developed smart-phone apps that link demanders for rides with drivers, who may, or may not be, part of the traditional taxi companies. Such start-ups have succeeded in taking a significant part of the taxi business away from the traditional operators. As a result, the price of taxi medallions on the open market has plunged. From trading in the range of \$1m. in New York in 2012, medallions are being offered in 2019 at about one fifth of that price. In Toronto, some medallions were traded in the range of \$300,000 in 2012, but are on offer in 2019 for prices in the range of \$30,000.

Not surprisingly, the traditional taxi companies charge that ride-hailing operators are violating the accepted rules governing the taxi business, and have launched legal suits against them and against local governments, and lobbied governments to keep them out of their cities.

In the new 'sharing economy', of which ride hailing companies are an example, participants operate with less traditional capital, and the communications revolution has been critical to their success. Home owners can use an online site to rent a spare bedroom in their house to visitors to their city (*Airbnb*), and thus compete with hotels. The main capital in this business is in the form of the information technology that links potential buyers to potential sellers.

Information on medallion prices in Canada can be found by, for example, searching at <http://www.kijiji.ca>

10.7 Invention, innovation and rent seeking

Invention and innovation are critical aspects of the modern economy. In some sectors of the economy, firms that cannot invent or innovate are liable to die. Invention is a genuine discovery, whereas innovation is the introduction of a new product or process.

Invention is the discovery of a new product or process through research.

Product innovation refers to new or better goods or services.

Process innovation refers to new or better production or supply.

To this point we have said little that is good about monopolies. However, the economist Joseph Schumpeter argued that, while monopoly leads to resource misallocation in the economy, this cost might be offset by the greater tendency for monopoly firms to invent and innovate. This is because such firms have more profit and therefore more resources with which to fund R&D and may therefore be more innovative than competitive firms. If this were true then, taking a long-run dynamic view of the marketplace, monopolies could have lower costs and more advanced products than competitive firms and thus benefit the consumer.

While this argument has some logical appeal, it falls short on several counts. First, even if large firms carry out more research than competitive firms, there is no guarantee that the ensuing benefits carry over to the consumer. Second, the results of such research may be used to prevent entry into the industry in question. Firms may register their inventions and gain use protection before a competitor can come up with the same or a similar invention. *Apple* and *Samsung* each own tens of thousands of patents. Third, the empirical evidence on the location of most R&D is inconclusive: A sector with several large firms, rather than one with a single or very many firms, may be best. For example, if *Apple* did not have *Samsung* as a competitor, or vice versa, would the pace of innovation be as strong?

Fourth, much research has a 'public good' aspect to it. Research carried out at universities and government-funded laboratories is sometimes referred to as basic research: It explores the principles underlying chemistry, social relations, engineering forces, microbiology, etc., and has multiple applications in the commercial world. If disseminated, this research is like a public good – its fruits can be used in many different applications, and its use in one area does not preclude its use in others. Consequently, rather

than protecting monopolies on the promise of more R&D, a superior government policy might be to invest directly in research and make the fruits of the research publicly available.

Modern economies have patent laws, which grant inventors a legal monopoly on use for a fixed period of time – perhaps fifteen years. By preventing imitation, patent laws raise the incentive to conduct R&D but do not establish a monopoly in the long run. Over the life of a patent the inventor charges a higher price than would exist if his invention were not protected; this both yields greater profits and provides the research incentive. When the patent expires, competition from other producers leads to higher output and lower prices for the product. Generic drugs are a good example of this phenomenon.

Patent laws grant inventors a legal monopoly on use for a fixed period of time.

The power of globalization once again is very relevant in patents. Not all countries have patent laws that are as strong as those in North America and Europe. The BRIC economies (Brazil, Russia, India and China) form an emerging power block. But their legal systems and enforcement systems are less well-developed than in Europe or North America. The absence of a strong and transparent legal structure inhibits research and development, because their fruits may be appropriated by competitors.

Rent seeking

Citizens are frequently appalled when they read of lobbying activities in their nation's capital. Every capital city in the world has an army of lobbyists, seeking to influence legislators and regulators. Such individuals are in the business of rent seeking, whose goal is to direct profit to particular groups, and protect that profit from the forces of competition. In Virginia and Kentucky we find that state taxes on cigarettes are the lowest in the US – because the tobacco leaf is grown in these states, and the tobacco industry makes major contributions to the campaigns of some political representatives.

Rent-seeking carries a resource cost: Imagine that we could outlaw the lobbying business and put these lobbyists to work producing goods and services in the economy instead. Their purpose is to maintain as much quasi-monopoly power in the hands of their clients as possible, and to ensure that the fruits of this effort go to those same clients. If this practice could be curtailed then the time and resources involved could be redirected to other productive ends.

Rent seeking is an activity that uses productive resources to redistribute rather than create output and value.

Industries in which rent seeking is most prevalent tend to be those in which the potential for economic profits is greatest – monopolies or near-monopolies. These, therefore, are the industries that allocate resources to the preservation of their protected status. We do not observe laundromat owners or shoe-repair businesses lobbying in Ottawa.

Conclusion

We have now examined two extreme types of market structure – perfect competition and monopoly. While many sectors of the economy operate in a way that is close to the competitive paradigm, very few are pure monopolies in that they have no close substitute products. Even firms like *Microsoft*, or *De Beers*, that supply a huge percentage of the world market for their product would deny that they are monopolies and would argue that they are subject to strong competitive pressures from smaller or 'fringe' producers. As a result we must look upon the monopoly paradigm as a useful way of analyzing markets, rather than being an exact description of the world. Accordingly, our next task is to examine how sectors with a few, several or multiple suppliers act when pursuing the objective of profit maximization. Many different market structures define the real economy, and we will concentrate on a limited number of the more important structures in the next chapter.

Key Terms

Monopolist: is the sole supplier of an industry's output, and therefore the industry and the firm are one and the same.

Natural monopoly: one where the ATC of producing any output declines with the scale of operation.

Marginal revenue is the change in total revenue due to selling one more unit of the good.

Average revenue is the price per unit sold.

Allocative inefficiency arises when resources are not appropriately allocated and result in deadweight losses.

Price discrimination involves charging different prices to different consumers in order to increase profit.

A **cartel** is a group of suppliers that colludes to operate like a monopolist.

Rent seeking is an activity that uses productive resources to redistribute rather than create output and value.

Invention is the discovery of a new product or process through research.

Product innovation refers to new or better products or services.

Process innovation refers to new or better production or supply.

Patent laws grant inventors a legal monopoly on use for a fixed period of time.

Exercises for Chapter 10

EXERCISE 10.1

Consider a monopolist with demand curve defined by $P=100-2Q$. The MR curve is $MR=100-4Q$ and the marginal cost is $MC=10+Q$. The demand intercepts are $\{ \$100, 50 \}$, the MR intercepts are $\{ \$100, 25 \}$.

1. Develop a diagram that illustrates this market, using either graph paper or an Excel spreadsheet, for values of output $Q = 1 \dots 25$.
2. Identify visually the profit-maximizing price and output combination.
3. *Optional:* Compute the profit maximizing price and output combination.

EXERCISE 10.2

Consider a monopolist who wants to maximize revenue rather than profit. She has the demand curve $P=72-Q$, with marginal revenue $MR=72-2Q$, and $MC=12$. The demand intercepts are $\{ \$72, 72 \}$, the MR intercepts are $\{ \$72, 36 \}$.

1. Graph the three functions, using either graph paper or an Excel spreadsheet.
2. Calculate the price she should charge in order to maximize revenue. [*Hint:* Where the $MR=0$.]
3. Compute the total revenue she will obtain using this strategy.

EXERCISE 10.3

Suppose that the monopoly in Exercise 10.2 has a large number of plants. Consider what could happen if each of these plants became a separate firm, and acted competitively. In this perfectly competitive world you can assume that the MC curve of the monopolist becomes the industry supply curve.

1. Illustrate graphically the output that would be produced in the industry?
2. What price would be charged in the marketplace?
3. *Optional:* Compute the gain to the economy in dollar terms as a result of the DWL being eliminated [*Hint:* It resembles the area ABF in Figure 10.14].

EXERCISE 10.4

In the text example in Table 10.1, compute the profit that the monopolist would make if he were able to price discriminate, by selling each unit at the demand price in the market.

EXERCISE 10.5

A monopolist is able to discriminate perfectly among his consumers – by charging a different price to each one. The market demand curve facing him is given by $P=72-Q$. His marginal cost is given by $MC=24$ and marginal revenue is $MR=72-2Q$.

1. In a diagram, illustrate the profit-maximizing equilibrium, where discrimination is not practiced. The demand intercepts are $\{ \$72, 72 \}$, the MR intercepts are $\{ \$72, 36 \}$.
2. Illustrate the equilibrium output if he discriminates perfectly.
3. *Optional:* If he has no fixed cost beyond the marginal production cost of \$24 per unit, calculate his profit in each pricing scenario.

EXERCISE 10.6

A monopolist faces two distinct markets A and B for her product, and she is able to insure that resale is not possible. The demand curves in these markets are given by $P_A=20-(1/4)Q_A$ and $P_B=14-(1/4)Q_B$. The marginal cost is constant: $MC=4$. There are no fixed costs.

1. Graph these two markets and illustrate the profit maximizing price and quantity in each market. [You will need to insert the MR curves to determine the optimal output.] The demand intercepts in A are $\{ \$20, 80 \}$, and in B are $\{ \$14, 56 \}$.

2. In which market will the monopolist charge a higher price?

EXERCISE 10.7

A concert organizer is preparing for the arrival of the Grateful Living band in his small town. He knows he has two types of concert goers: One group of 40 people, each willing to spend \$60 on the concert, and another group of 70 people, each willing to spend \$40. His total costs are purely fixed at \$3,500.

1. Draw the market demand curve faced by this monopolist.
2. Draw the MR and MC curves.
3. With two-price discrimination what will be the monopolist's profit?
4. If he must charge a single price for all tickets can he make a profit?

EXERCISE 10.8

Optional: A monopolist faces a demand curve $P=64-2Q$ and $MR=64-4Q$. His marginal cost is $MC=16$.

1. Graph the three functions and compute the profit maximizing output and price.
 2. Compute the efficient level of output (where $MC=demand$), and compute the DWL associated with producing the profit maximizing output rather than the efficient output.
1. It is not difficult to show that the demand curve corresponding to the data in the table is given by the expression $P=12-0.5q$. Since the MR curve has twice the slope of the demand curve, it is given by $MR=12-q$. The data indicate that the MC curve can be written as $MC=0.5q$ and the average cost curve by $AC=0.25q$. Setting $MR=MC$ yields $q=8$. At this output the demand curve implies the price is \$8. Profit is $(P-AC) \times q = (\$8 - \$2) \times 8 = \$48$.
 2. We can think of such a transformation coming from a single supplier taking over a number of small suppliers, and the monopolist would thus be a multi-plant firm.

This page titled [10: Monopoly](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.1: Monopolies

In analyzing perfect competition we emphasized the difference between the industry and the individual supplier. The individual supplier is an atomistic unit with no market power. In contrast, a monopolist has a great deal of market power, for the simple reason that a monopolist is the sole supplier of a particular product and so really *is* the industry. The word monopoly, comes from the Greek words *monos*, meaning one, and *polein* meaning to sell. When there is just a single seller, our analysis need not distinguish between the industry and the individual firm. They are the same on the supply side.

Furthermore, the distinction between long run and short run is blurred, because a monopoly that continues to survive as a monopoly obviously sees no entry or exit. This is not to say that monopolized sectors of the economy do not evolve, they do. Sometimes they die, sometimes they evolve in a different role. For example, when digital cameras entered the market place in the eighties the Polaroid Land camera (which printed film straight out of the camera) 'died' because the demand side of the market lost interest. The Blackberry 'smart' phone had a virtual monopoly on this product into the new millennium until Apple and Nokia entered the market.

A **monopolist** is the sole supplier of an industry's output, and therefore the industry and the firm are one and the same.

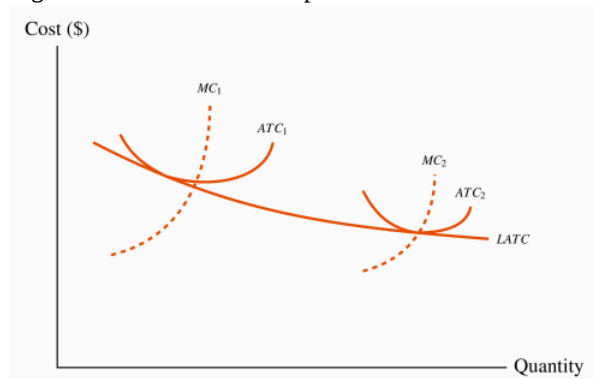
Monopolies can exist and exert their dominance in the market place for several reasons; scale economies, national policy, successful prevention of entry, research and development combined with patent protection.

Natural monopolies

Traditionally, monopolies were viewed as being 'natural' in some sectors of the economy. This means that scale economies define some industries' production and cost structures up to very high output levels, and that the whole market might be supplied at least cost by a single firm.

Consider the situation depicted in Figure 10.1. The long-run *ATC* curve declines indefinitely. There is no output level where average costs begin to increase. Imagine now having several firms, each producing with a plant size corresponding to the short-run average cost curve ATC_1 , or alternatively a single larger firm using a plant size denoted by ATC_2 . The small firms in this case cannot compete with the larger firm because the larger firm has lower production costs and can undercut the smaller firms, and *supply the complete market* in the process. Such a scenario is termed a natural monopoly.

Figure 10.1 A 'natural' monopolist



When LR average costs continue to decline at very high output, one large firm may be able to supply the industry at a lower unit cost than several smaller firms. With a plant size corresponding to ATC_2 , a single supplier can supply the whole market, whereas several smaller firms, each with plant size corresponding to ATC_1 , cannot compete with the larger firm on account of differential unit costs.

Natural monopoly: one where the *ATC* of producing any output declines with the scale of operation.

Electricity distribution in some of Canada's provinces is in the hands of a single supplier – *Hydro Quebec* or *Hydro One* in Ontario, for example. These distributors are natural monopolies in the sense described above: Unit distribution costs decline with size. In contrast, electricity production is not 'naturally' a monopoly. Other suppliers were once thought of as 'natural' monopolies also, but are no longer. *Bell Canada* was considered to be a natural monopoly in the era of land lines: It would not make economic sense to run several sets of phone lines to every residence. But that was before the arrival of cell phones, broadband and satellites. *Canada*

Post was also thought to be a natural monopoly, until the advent of *FEDEX*, *UPS* and other couriers proved otherwise. Invention can compete away a 'natural' monopoly.

In reality there are very few *pure* monopolies. *Facebook*, *Microsoft*, *Amazon*, *Apple*, *Netflix* and *Google* may be extraordinarily dominant in their markets, but they are not the only suppliers of the services or products that they offer. There exist other products that are similar.

National and Provincial Policy

Government policy can foster monopolies. Some governments are, or once were, proud to have a 'national carrier' in the airline industry – *Air Canada* in Canada or *British Airways* in the UK. The mail service was viewed as a symbol of nationhood in Canada and the US: *Canada Post* and the *US Postal* system are national emblems that have historic significance. They were vehicles for integrating the provinces or states at various points in the federal lives of these countries.

In the modern era, most of Canada's provinces have decided to create a provincial monopoly crown corporation for the sale of cannabis. But competition abounds in the form of an illegal market.

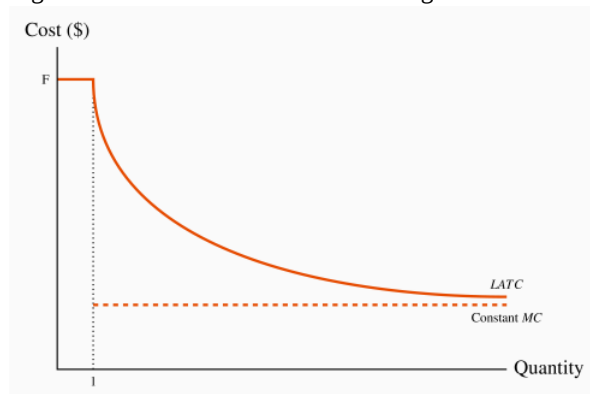
The down side of such nationalist policies is that they can be costly to the taxpayer. Industries that are not subject to competition can become fat and uncompetitive: Managers have insufficient incentives to curtail costs; unions realize the government is committed to sustain the monopoly and push for higher wages than under a more competitive structure, and innovation may be less likely to occur.

Maintaining barriers to entry

Monopolies can continue to survive if they are successful in preventing the entry of new firms and products. *Patents* and *copyrights* are one vehicle for preserving the sole-supplier role, and are certainly necessary to encourage firms to undertake the research and development (R&D) for new products.

Many corporations produce products that require a large up-front investment; this might be in the form of research and development, or the construction of costly production facilities. For example, *Boeing* or *Airbus* incurs billions of dollars in developing new aircraft; pharmaceuticals may have to invest a billion dollars to develop a new drug. However, once such an investment is complete, the cost of producing each unit of output may be relatively low. This is particularly true in pharmaceuticals. Such a phenomenon is displayed in Figure 10.2. In this case the average cost for a small number of units produced is high, but once the fixed cost is spread over an ever larger output, the average cost declines rapidly, and in the limit approaches the marginal cost. These production structures are common in today's global economy, and they give rise to markets characterized either by a single supplier or a small number of suppliers.

Figure 10.2 Fixed cost and constant marginal cost



With a fixed cost of producing the first unit of output equal to F and a constant marginal cost thereafter, the long-run average total cost, $LATC$, declines indefinitely and becomes asymptotic to the marginal cost curve.

This figure is useful in understanding the role of patents. Suppose that *Pharma A* spends one billion dollars in developing a new drug and has constant unit production costs thereafter, while *Pharma B* avoids research and development and simply imitates *Pharma A*'s product. Clearly *Pharma B* would have a $LATC$ equal to its LMC , and would be able to undercut the initial developer of the drug. Such an outcome would discourage investment in new products and the economy at large would suffer as a

consequence. Economies would be worse off if protection is not provided to the developers of new products because, if such protection is not offered, potential developers will not have the incentive to incur the up-front investment required.

While copyright and patent protection is legal, *predatory pricing* is an illegal form of entry barrier, and we explore it more fully in Chapter 14. An example would be where an existing firm that sells nationally may deliberately undercut the price of a small local entrant to the industry. Airlines with a national scope are frequently accused of posting low fares on flights in regional markets that a new carrier is trying to enter.

Political lobbying is another means of maintaining monopolistic power. For example, the *Canadian Wheat Board* had fought successfully for decades to prevent independent farmers from marketing wheat. This Board lost its monopoly status in August 2012, when the government of the day decided it was not beneficial to consumers or farmers in general. Numerous 'supply management' policies are in operation all across Canada. Agriculture is protected by production quotas. All maple syrup in Quebec must be marketed through a single monopoly supplier.

Critical networks also form a type of barrier, though not always a monopoly. *Microsoft's Office* package has an almost monopoly status in word processing and spreadsheet analysis for the reason that so many individuals and corporations use it. The fact that we know a business colleague will be able to edit our documents if written in *Word*, provides us with an incentive to use *Word*, even if we might prefer *Wordperfect* as a vehicle for composing documents. We develop the concept of strategic entry prevention further in Chapter 11.

This page titled [10.1: Monopolies](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.2: Profit maximizing behaviour

We established in the previous chapter that, in deciding upon a profit-maximizing output, any firm should produce up to the point where the additional cost equals the additional revenue from a unit of output. What distinguishes the supply decision for a monopolist from the supply decision of the perfect competitor is that the monopolist faces a downward sloping demand. A monopolist is the sole supplier and therefore must meet the full market demand. This means that if more output is produced, the price must fall. We will illustrate the choice of a profit maximizing output using first a marginal-cost/marginal-revenue approach; then a supply/demand approach.

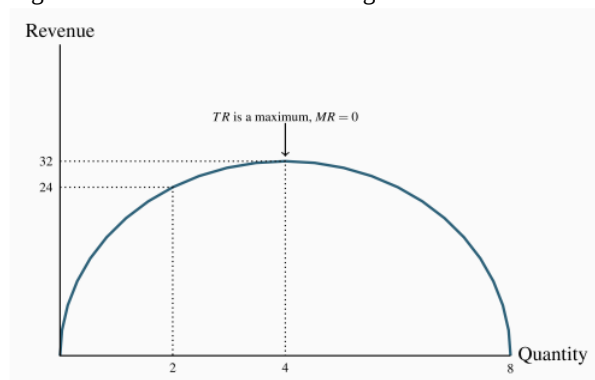
Marginal revenue and marginal cost

Table 10.1 displays price and quantity values for a demand curve in columns 1 and 2. Column 3 contains the sales revenue generated at each output. It is the product of price and quantity. Since the price denotes the revenue per unit, it is sometimes referred to as average revenue. The total revenue (TR) reaches a maximum at \$32, where 4 units of output are produced. A greater output necessitates a lower price on every unit sold, and in this case revenue falls if the fifth unit is brought to the market. Even though the fifth unit sells for a positive price, the price on the other 4 units is now lower and the net effect is to reduce total revenue. This pattern reflects what we examined in Chapter 4: As price is lowered from the highest possible value of \$14 (where 1 unit is demanded) and the corresponding quantity increases, revenue rises, peaks, and ultimately falls as output increases. In Chapter 4 we explained that this maximum revenue point occurs where the price elasticity is unity (-1), at the midpoint of a linear demand curve.

Table 10.1 A profit maximizing monopolist

Quantity (Q)	Price (P)	Total revenue (TR)	Marginal revenue (MR)	Marginal cost (MC)	Total cost (TC)	Profit
0	16					
1	14	14	14	2	2	12
2	12	24	10	3	5	19
3	10	30	6	4	9	21
4	8	32	2	5	14	18
5	6	30	-2	6	20	10
6	4	24	-6	7	27	-3
7	2	14	-10	8	35	-21

Figure 10.3 Total revenue and marginal revenue



When the quantity sold increases total revenue/expenditure initially increases also. At a certain point, further sales require a price that not only increases quantity, but reduces revenue on units already being sold to such a degree that TR declines – where the demand elasticity equals -1 (the mid point of a linear demand curve). Here the midpoint occurs at $Q=4$. Where the TR is a maximum the $MR=0$.

Related to the total revenue function is the marginal revenue function. It is the addition to total revenue due to the sale of one more unit of the commodity.

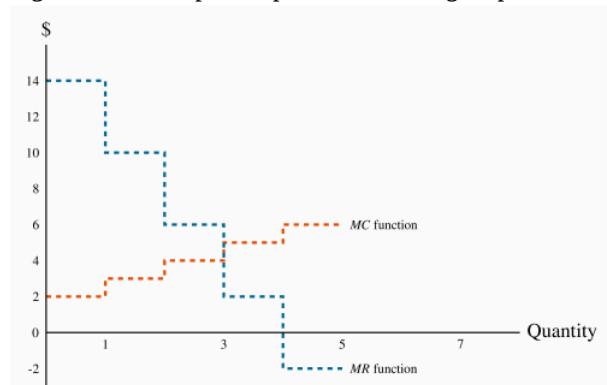
Marginal revenue is the change in total revenue due to selling one more unit of the good.

Average revenue is the price per unit sold.

The *MR* in this example is defined in the fourth column of Table 10.1. When the quantity sold increases from 1 unit to 2 units total revenue increases from \$14 to \$24. Therefore the marginal revenue associated with the second unit of output is \$10. When a third unit is sold *TR* increases to \$30 and therefore the *MR* of the third unit is \$6. As output increases the *MR* declines and eventually becomes negative – at the point where the *TR* is a maximum: If *TR* begins to decline then the additional revenue is by definition negative.

The *MR* function is plotted in Figure 10.4. It becomes negative when output increases from 4 to 5 units.

Figure 10.4 Monopolist's profit maximizing output



It is optimal for the monopolist to increase output as long as *MR* exceeds *MC*. In this case $MR > MC$ for units 1, 2 and 3. But for the fourth unit $MC > MR$ and therefore the monopolist would reduce total profit by producing it. He should produce only 3 units of output.

The optimal output

This producer has a marginal cost structure given in the fifth column of the table, and this too is plotted in Figure 10.4. Our profit maximizing rule from Chapter 8 states that it is optimal to produce a greater output as long as the additional revenue exceeds the additional cost of production on the next unit of output. In perfectly competitive markets the additional revenue is given by the fixed price for the individual producer, whereas for the monopolist the additional revenue is the marginal revenue. Consequently as long as *MR* exceeds *MC* for the next unit a greater output is profitable, but once *MC* exceeds *MR* the production of additional units should cease.

From Table 10.1 and Figure 10.4 it is clear that the optimal output is at 3 units. The third unit itself yields a profit of \$2, the difference between *MR* (\$6) and *MC* (\$4). A fourth unit however would reduce profit by \$3, because the *MR* (\$2) is less than the *MC* (\$5). What price should the producer charge? The price, as always, is given by the demand function. At a quantity sold of 3 units, the corresponding price is \$10, yielding total revenue of \$30.

Profit is the difference between total revenue and total cost. In Chapter 8 we computed total cost as the average cost times the number of units produced. It can also be computed as the sum of costs associated with each unit produced: The first unit costs \$2, the second \$3 and the third \$4. The total cost of producing 3 units is the sum of these dollar values: $\$9 = \$2 + \$3 + \4 . The profit-maximizing output therefore yields a profit of \$21 ($\$30 - \9).

Supply and demand

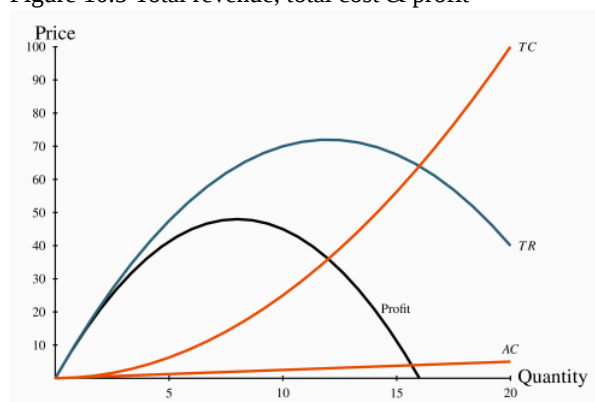
When illustrating market behaviour it is convenient to describe behaviour by simple linear supply and demand functions that are continuous, rather than the 'step' functions used in the preceding example. As explained in Chapter 5, in using continuous curves to represent a market we implicitly assume that a unit of output can be broken into subunits. In the example above we assumed that sales always involve one whole unit of the product being sold. In fact many goods can be sold in fractional units: Gasoline can be sold in fractions of a litre; fruits and vegetables can be sold in fractions of a kilogram, and so forth. Table 10.2 below furnishes the data for our analysis.

Table 10.2 Discrete quantities

Price	Quantity	Total	Total	Profit
	demanded	revenue	cost	
12	0	0	0	0
11	2	22	1	21
10	4	40	4	36
9	6	54	9	45
8	8	64	16	48
7	10	70	25	45
6	12	72	36	36
5	14	70	49	21
4	16	64	64	0
3	18	54	81	-27
2	20	40	100	-60
1	22	22	121	-99
0	24	0	144	-144

The first two columns define the demand curve. Total revenue is the product of price and quantity and given in column 3. The cost data are given in column 4, and profit – the difference between total revenue and total cost is in the final column. Profit is maximized where the difference between revenue and cost is greatest; in this case where the output is 8 units. At lower or higher outputs profit is less. Figure 10.5 contains the curves defining total revenue (TR), total cost (TC) and profit. These functions can be obtained by mapping all of the revenue-quantity combinations, the cost-quantity combinations, and the profit-quantity combinations as a series of points, and joining these points to form the smooth functions displayed. The vertical axis is measured in dollars, the horizontal axis in units of output. Graphically, profit is maximized where the dollar difference between TR and TC is greatest; that is at the output where the vertical distance between the two curves is greatest. This difference, which is also defined by the profit curve, occurs at a value of 8 units, corresponding to the outcome in Table 10.2.

Figure 10.5 Total revenue, total cost & profit



At any quantity less than this output, profit would rise with additional output. This is because, from a less-than-optimal output, the additional revenue from increased sales exceeds the increased cost associated with producing those units: Stated differently, the marginal revenue would exceed the marginal cost. Conversely, outputs greater than the optimum result in a MR less than the associated MC . Accordingly, since outputs where $MR > MC$ are too low, and outputs where $MR < MC$ are too high, the optimum must be where the $MR = MC$. Hence, the equality between MR and MC is implied in this diagram at the output where the difference between TR and TC is greatest.

Note finally that total revenue is maximized where the TR curve reaches a peak. In this example that occurs at a value of 12 units of output. This is to be anticipated, as we learned in Chapter 4, because the midpoint of the demand schedule in Table 10.2 occurs at that value.

Figure 10.6 Market demand, the MR curve, and the monopolist's AC and MC curves

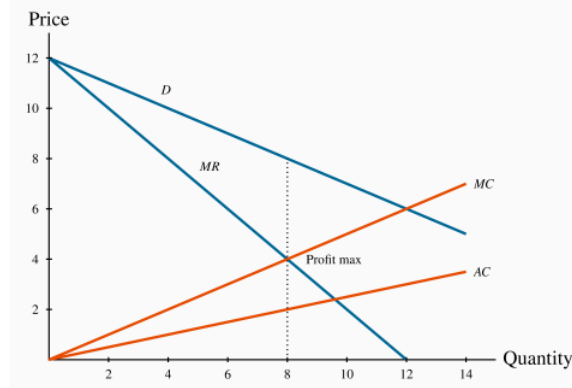


Figure 10.6 displays the demand curve for the market, the MR curve, and the monopolist's MC and AC curves. Consider first the marginal revenue curve. In contrast to the previous example, where only whole or integer units could be sold, in this example units can be sold in fractional amounts, and the MR curve must reflect this. To determine the position of the MR curve, note that with a straight-line demand curve total revenue is a maximum at the midpoint of the demand curve. Any increase in output results in reduced revenue: Stated differently, the marginal revenue becomes negative at that output. Up to that output the MR is positive, as illustrated in Figure 10.3. Accordingly, the MR curve must intersect the quantity axis midway between zero and the horizontal-axis intercept of the demand curve. Geometrically, since the MR intersects the quantity axis half way to the horizontal intercept of the demand curve, it must have a slope that is twice the slope of the demand curve.

By observing the data in columns 1 and 2 of the table, the demand curve intercepts are $\{ \$12, 24 \}$, and from above discussion the MR curve has intercepts $\{ \$12, 12 \}$. The AC is obtained by dividing TC by output in Table 10.2, and the MC can be also calculated as the change in total cost divided by the change in output from Table 10.2. The result of these calculations is displayed in Figure 10.6.

The profit maximizing output is 8 units, where $MC=MR$. The price at which 8 units can be sold is read from the demand curve¹, or the first column in Table 10.2. It is \$8. And, as expected, this price-quantity combination maximizes profit. Table 10.2 indicates that profit is maximized at \$48, at $q=8$.

Demand elasticity and marginal revenue

We have shown above that the MR curve cuts the horizontal axis at a quantity where the elasticity of demand is unity. We know from Chapter 4 that demand is elastic at points on the demand curve above this unit-elastic point. Furthermore, since the intersection of MR and MC must be at a positive dollar value (MC cannot be negative), then it must be the case that the *profit maximizing price for a monopolist always lies on the elastic segment of the demand curve*.

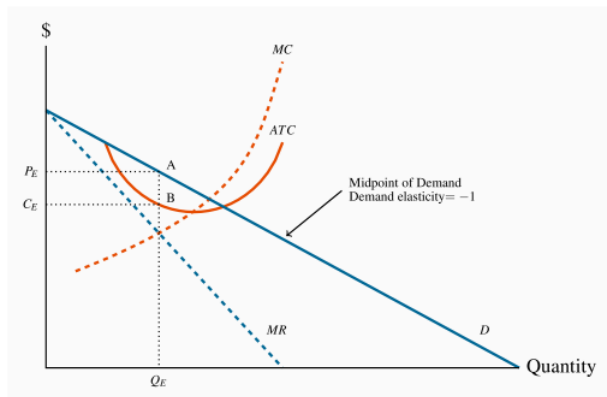
A general graphical representation

In Figure 10.7 we generalize the graphical representation of the monopoly profit maximizing output by allowing the MC and ATC curves to be nonlinear. The optimal output is at Q_E , where $MR=MC$, and the price P_E sustains that output. With the average cost known, profit per unit is AB , and therefore total profit is this margin multiplied by the number of units sold, Q_E .

Total profit is therefore $P_E ABC_E$.

Note that the monopolist may not always make a profit. Losses could result in Figure 10.7 if average costs were to rise so that the ATC were everywhere above the demand curve, or if the demand curve shifted down to being everywhere below the ATC curve. In the longer term the monopolist would have to either reduce costs or perhaps stimulate demand through advertising if she wanted to continue in operation.

Figure 10.7 The monopoly equilibrium



The profit maximizing output is Q_E , where $MC=MR$. This output can be sold at a price P_E . The cost per unit of Q_E is read from the ATC curve, and equals B . Per unit profit is therefore AB and total profit is P_EABC_E .

This page titled [10.2: Profit maximizing behaviour](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

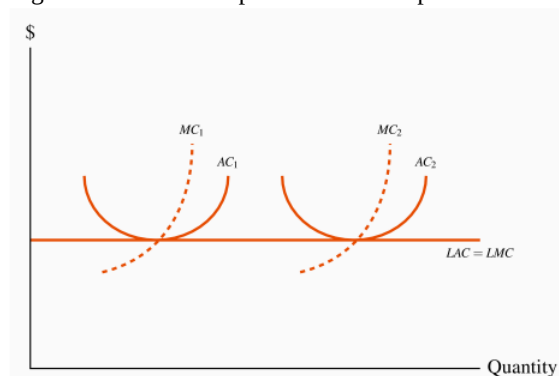
10.3: Long-run choices

Consider next the impact of a shift in demand upon the profit maximizing choice of this firm. A rightward shift in demand in Figure 10.7 also yields a new MR curve. The firm therefore chooses a new level of output, using the same profit maximizing rule: Set $MC=MR$. This output will be greater than the previous output, but again the price must be on an elastic portion of the new demand curve. If operating with the same plant size, the MC and ATC curves do not change and the new profit per unit is again read from the ATC curve.

By this stage the curious student will have asked: "What happens to plant size in the long run?" For example, is the monopolist in Figure 10.7 using the most appropriate plant size in the first place? Even if she is, should the monopolist consider adopting an expanded plant size in response to the shift in demand?

The answer is: In the long run the monopolist is free to choose whatever plant size is best. Her initial plant size might have been optimal for the demand she faced, but if it was, it is unlikely to be optimal for the larger scale of production associated with the demand shift. Accordingly, with the new demand curve, she must consider how much profit she could make using different plant sizes.

Figure 10.8 The monopolist's choice of plant size

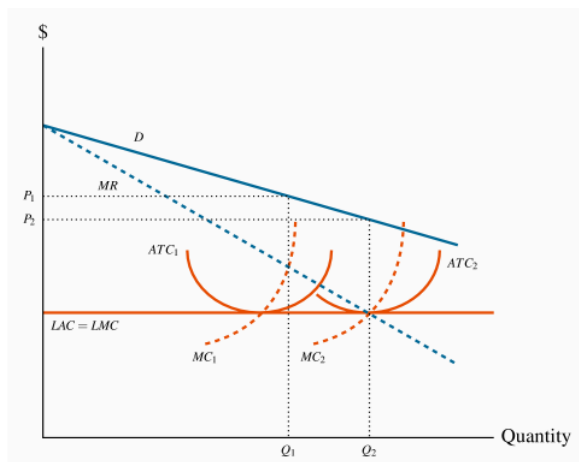


With constant returns to scale and constant prices per unit of labour and capital, a doubling of output involves exactly a doubling of costs. Thus, per unit costs, or average costs, are constant in the LR. Hence $LAC=LMC$, and each is constant.

To illustrate one possibility, we will think of this firm as having constant returns to scale at all output ranges, as displayed in Figure 10.8. (Our reasoning carries through if the LAC slopes downwards; the graph just becomes a little more complex.) The key characteristic of constant returns to scale is that a doubling of inputs leads to a doubling of output. Therefore, if the per-unit cost of inputs is fixed, a doubling of inputs (and therefore output) leads exactly to a doubling of costs. This implies that, when the firm varies its plant size *and* its labour use, the cost of producing each additional unit must be constant. The long-run marginal cost LMC is therefore constant and equals the ATC in the long run.

Figure 10.9 describes the market for this good. The optimal output and price are determined in the usual manner: Set $MC=MR$. If the monopolist has plant size corresponding to ATC_1 , the optimal output is Q_1 and should be sold at the price P_1 . The key issue now is: Given the demand conditions, could the monopolist make more profit by choosing a plant size that differs from the one corresponding to ATC_1 ?

Figure 10.9 Plant size in the long run



With demand conditions defined by D and MR , the optimal plant size is one corresponding to the point where $MR=MC$ in the long run. Therefore Q_2 is the optimal output and the optimal plant size corresponds to ATC_2 . If the current plant is defined by ATC_1 , then optimal SR production is Q_1 .

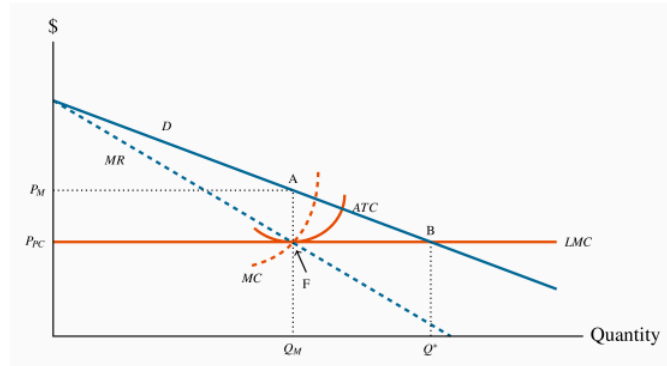
In this instance the answer is a clear 'yes'. Her LMC curve is horizontal and so, by increasing output from Q_1 to Q_2 she earns a profit on each additional unit in that range, because the MR curve lies above the LMC curve. In order to produce the output level Q_2 at least cost she must choose a plant size corresponding to AC_2 .

This page titled [10.3: Long-run choices](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis](#) and [Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.4: Output inefficiency

A characteristic of perfect competition is that it secures an efficient allocation of resources when there are no externalities in the market: Resources are used up to the point where their marginal cost equals their marginal value – as measured by the price that consumers are willing to pay. But a monopoly structure does not yield this output. Consider Figure 10.10.

Figure 10.10 Monopoly output inefficiency



A monopolist maximizes profit at Q_M . Here the value of marginal output exceeds cost. If output expands to Q^* a gain arises equal to the area ABF. This is the deadweight loss associated with the output Q_M rather than Q^* . If the monopolist's long-run MC is equivalent to a competitive industry's supply curve, then the deadweight loss is the cost of having a monopoly rather than a perfectly competitive market.

The monopolist's profit-maximizing output Q_M is where MC equals MR . This output is inefficient for the reason that we developed in Chapter 5: If output is increased beyond Q_M the additional benefit exceeds the additional cost of producing it. The additional benefit is measured by the willingness of buyers to pay – the market demand curve. The additional cost is the long-run MC curve under the assumption of constant returns to scale. Using the terminology from Chapter 5, there is a deadweight loss equal to the area ABF. This is termed allocative inefficiency.

Allocative inefficiency arises when resources are not appropriately allocated and result in deadweight losses .

Perfect competition versus monopoly

The area ABF can also be considered as the efficiency loss associated with having a monopoly rather than a perfectly competitive market structure. In perfect competition the supply curve is horizontal. This is achieved by having firms enter and exit when more or less must be produced. Accordingly, *if the perfectly competitive industry's supply curve approximates the monopolist's long-run marginal cost curve*², we can say that if the monopoly were turned into a competitive industry, output would increase from Q_M to Q^* . The deadweight loss is one measure of the superiority of the perfectly competitive structure over the monopoly structure.

Note that this critique of monopoly is not initially focused upon profit. While monopoly profits are what frequently irk the public, we have focused upon resource allocation inefficiencies. But in a real sense the two are related: Monopoly inefficiencies arise through output being restricted, and it is this output reduction – achieved by maintaining a higher than competitive price – that gives rise to those profits. Nonetheless, there is more than just a shift in purchasing power from the buyer to the seller. Deadweight losses arise because output is at a level lower than the point where the MC equals the value placed on the good; thus the economy is sacrificing the possibility of creating additional surplus.

Given that monopoly has this undesirable inefficiency, what measures should be taken, if any, to counter the inefficiency? We will see what Canada's Competition Act has to say in Chapter 14 and also examine what other measures are available to control monopolies.

This page titled [10.4: Output inefficiency](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.5: Price discrimination

A common characteristic in the pricing of many goods is that different individuals pay different prices for goods or services that are essentially the same. Examples abound: Seniors get a reduced rate for coffee in *Burger King*; hair salons charge women more than they charge men; bank charges are frequently waived for juniors. Price discrimination involves charging different prices to different consumers in order to increase profit.

Price discrimination involves charging different prices to different consumers in order to increase profit.

A strict definition of discrimination involves different prices for *identical products*. We all know of a school friend who has been willing to take the midnight flight to make it home at school break at a price he can afford. In contrast, the business executive prefers the seven a.m. flight to arrive for a nine a.m. business meeting in the same city at several times the price. These are very mild forms of price discrimination, since a midnight flight (or a midday flight) is not a perfect substitute for an early morning flight. Price discrimination is practiced because buyers are willing to pay different amounts for a good or service, and the supplier may have a means of profiting from this. Consider the following example.

Family Flicks is the local movie theatre. It has two distinct groups of customers – those of prime age form one group; youth and seniors form the other. Family Flicks has done its market research and determined that each group accounts for 50 percent of the total market of 100 potential viewers per screening. It has also established that the prime-age group members are willing to pay \$12 to see a movie, while the seniors and youth are willing to pay just \$5. How should the tickets be priced?

Family Flicks has no variable costs, only fixed costs. It must pay a \$100 royalty to the movie maker each time it shows the current movie, and must pay a cashier and usher \$20 each. Total costs are therefore \$140, regardless of how many people show up – short-run *MC* is zero. On the pricing front, as illustrated in Table 10.3 below, if Family Flicks charges \$12 per ticket it will attract 50 viewers, generate \$600 in revenue and therefore make a profit of \$460.

Table 10.3 Price discrimination

	<i>P</i> =\$5	<i>P</i> =\$12	Twin price
No. of customers	100	50	
Total revenue	\$500	\$600	\$850
Total costs	\$140	\$140	\$140
Profit	\$360	\$460	\$710

In contrast, if it charges \$5 it can fill the theatre, because each of the prime-age individuals is willing to pay more than \$5, but the seniors and youth are now offered a price they too are willing to pay. However, the total revenue is now only \$500 ($100 \times \$5 = \500), and profits are reduced to \$360. It therefore decides to charge the high price and leave the theatre half-empty, because this strategy maximizes its profit.

Suppose finally that the theatre is able to segregate its customers. It can ask the young and senior customers for identification upon entry, and in this way charge them a lower price, *while still maintaining the higher price to the prime-age customers*. If it can execute such a plan Family Flicks can now generate \$850 in revenue – \$600 from the prime-age group and \$250 from the youth and seniors groups. Profit soars to \$710.

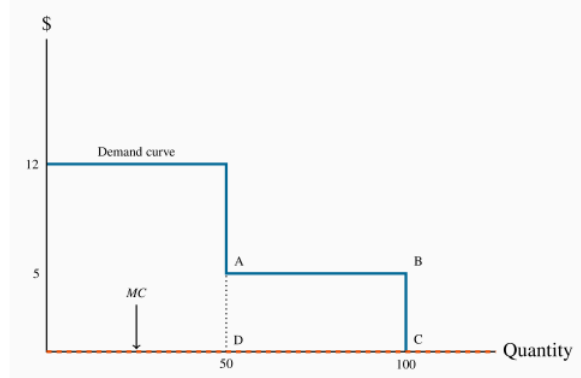
There are two important conditions for this scheme to work:

1. The seller must be able to *segregate the market* at a reasonable cost. In the movie case this is achieved by asking for identification.
2. The second condition is that *resale must be impossible or impractical*. For example, we rule out the opportunity for young buyers to resell their tickets to the prime-age individuals. Sellers have many ways of achieving this – they can require immediate entry to the movie theatre upon ticket purchase, they can stamp the customer's hand, they can demand the showing of ID with the ticket when entering the theatre area.

Frequently we think of sellers who offer price reductions to specific groups as being generous. For example, hotels may levy only a nominal fee for the presence of a child, once the parents have paid a suitable rate for the room or suite in which a family stays. The hotel knows that if it charges too much for the child, it may lose the whole family as a paying unit. The coffee shop offering cheap

coffee to seniors is interested in getting a price that will cover its variable cost and so contribute to its profit. It is unlikely to be motivated by philanthropy, or to be concerned with the financial circumstances of seniors.

Figure 10.11 Price discrimination at the movies



At $P=12$, 50 prime-age individuals demand movie tickets. At $P=5$, 50 more seniors and youths demand tickets. Since the MC is zero the efficient output is where the demand curve takes a zero value – where all 100 customers purchase tickets. Thus, any scheme that results in all 100 individuals buying ticket is efficient. Efficient output is at point C.

Price discrimination has a further interesting feature that is illustrated in Figure 10.11: It frequently *reduces the deadweight loss* associated with a monopoly seller!

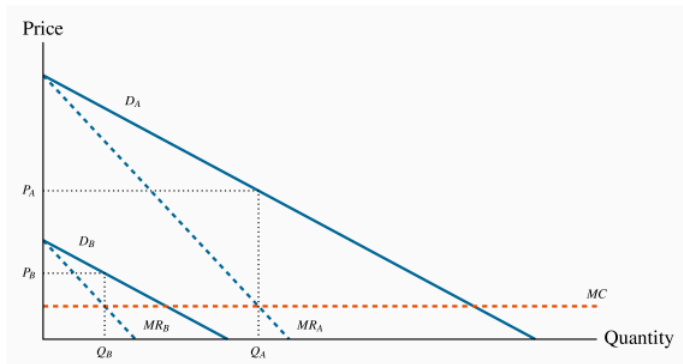
In our Family Flicks example, the profit maximizing monopolist that did not, or could not, price discriminate *left 50 customers unsupplied who were willing to pay \$5 for a good that had a zero MC*. This is a deadweight loss of \$250 because 50 seniors and youth valued a commodity at \$5 that had a zero MC . Their demand was not met because, in the absence of an ability to discriminate between consumer groups, Family Flicks made more profit by satisfying the demand of the prime-age group alone. But in this example, by segregating its customers, the firm's profit maximization behaviour resulted in the DWL being eliminated, because it supplied the product to those additional 50 individuals. In this instance *price discrimination improves welfare*, because more of a good is supplied in a situation where market valuation exceeds marginal cost.

In the preceding example we simplified the demand side of the market by assuming that every individual in a given group was willing to pay the same price – either \$12 or \$5. More realistically each group can be defined by a downward-sloping demand curve, reflecting the variety of prices that buyers in a given market segment are willing to pay. It is valuable to extend the analysis to include this reality. For example, a supplier may face different demands from her domestic and foreign buyers, and if she can segment these markets she can price discriminate effectively.

Consider Figure 10.12 where two segmented demands are displayed, D_A and D_B , with their associated marginal revenue curves, MR_A and MR_B . We will assume that marginal costs are constant for the moment. It should be clear by this point that the profit maximizing solution for the monopoly supplier is to supply an amount to each market where the MC equals the MR in each market: Since the buyers in one market cannot resell to buyers in the other, the monopolist considers these as two different markets and therefore maximizes profit by applying the standard rule. She will maximize profit in market A by supplying the quantity Q_A and in market B by supplying Q_B . The prices at which these quantities can be sold are P_A and P_B . These prices, unsurprisingly, are different – the objective of segmenting markets is to increase profit by treating the markets as distinct.

An example of this type of price discrimination is where pharmaceutical companies sell drugs to less developed economies at a lower price than to developed economies. The low price is sufficient to cover marginal cost and is therefore profitable - provided the high price market covers the fixed costs.

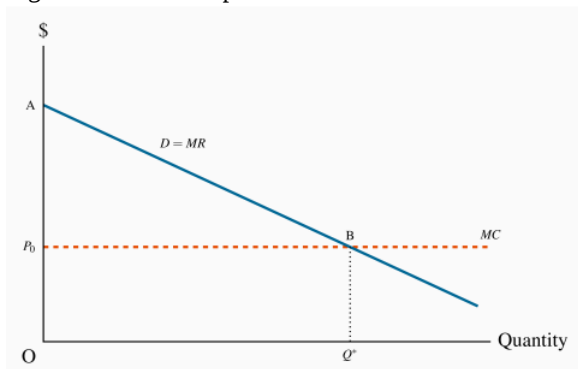
Figure 10.12 Pricing in segregated markets



With two separate markets defined by D_A and D_B , and their associated MR_A and MR_B , a profit maximizing strategy is to produce where $MC=MR_A=MR_B$, and *discriminate* between the two markets by charging prices P_A and P_B .

The preceding examples involved two separable groups of customers and are very real. This kind of group segregation is sometimes called *third degree price discrimination*. But it may be possible to segregate customers into several groups rather than just two. In the limit, if we could charge a different price to every consumer in a market, or for every unit sold, the revenue accruing to the monopolist would be the area under the demand curve up to the output sold. Though primarily of theoretical interest, this is illustrated in Figure 10.13. It is termed *perfect price discrimination*, and sometimes *first degree price discrimination*. Such discrimination is not so unrealistic: A tax accountant may charge different customers a different price for providing the same service; home renovators may try to charge as much as any client appears willing to pay.

Figure 10.13 Perfect price discrimination



A monopolist who can sell each unit at a different price maximizes profit by producing Q^* . With each consumer paying a different price the demand curve becomes the MR curve. The result is that the monopoly DWL is eliminated because the efficient output is produced, and the monopolist appropriates all the consumer surplus. Total revenue for the perfect price discriminator is $OABQ^*$.

Second degree price discrimination is based on a different concept of buyer identifiability. In the cases we have developed above, the seller is able to distinguish the buyers by *observing* a vital characteristic that signals their type. It is also possible that, while individuals might have defining traits which influence their demands, such traits might not be detectable by the supplier. Nonetheless, it is frequently possible for the supplier to offer different pricing options (corresponding to different uses of a product) that buyers would choose from, with the result that her profit would be greater than under a uniform price with no variation in the use of the service. Different cell phone 'plans', or different internet plans that users can choose from are examples of this second-degree discrimination.

This page titled 10.5: Price discrimination is shared under a CC BY-NC-SA 4.0 license and was authored, remixed, and/or curated by Douglas Curtis and Ian Irvine (Lyryx) via source content that was edited to the style and standards of the LibreTexts platform.

10.6: Cartels- Acting like a monopolist

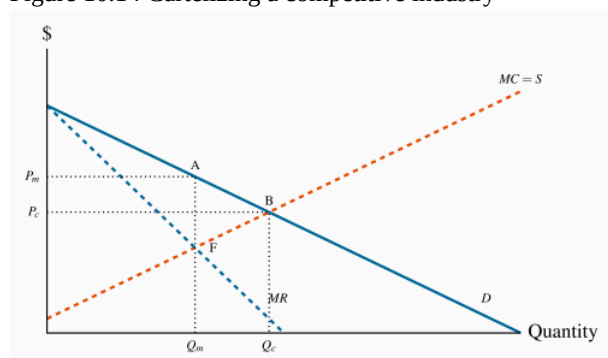
A cartel is a group of suppliers that colludes to operate like a monopolist. The cartel formed by the members of the Organization of Oil Exporting Countries (OPEC) is an example of a cartel that was successful in achieving its objectives for a long period. This cartel first flexed its muscles in 1973, by increasing the world price of oil from \$3 per barrel to \$10 per barrel. The result was to transfer billions of dollars from the energy-importing nations in Europe and North America to OPEC members – the demand for oil is relatively inelastic, hence an increase in price increases total expenditures.

A **cartel** is a group of suppliers that colludes to operate like a monopolist.

A second renowned cartel is managed by De Beers, which controls a large part of the world's diamond supply. In Canada, agricultural marketing boards are a means of restricting supply legally. Such cartels may have thousands of members. By limiting entry, through requiring a production 'quota', the incumbents can charge a higher price than if entry to the industry were free.

To illustrate the dynamics of cartels consider Figure 10.14. Several producers, with given production capacities, come together and agree to restrict output with a view to increasing price and therefore profit. This may be done with the agreement of the government, or it may be done secretly, and possibly against the law. Each firm has a MC curve, and the industry supply is defined as the sum of these marginal cost curves, as illustrated in Figure 9.3. The resulting cartel is effectively one in which there is a single supplier with many different plants – a multi-plant monopolist. To maximize profits this organization will choose an output level Q_m where the MR equals the MC . In contrast, if these firms act competitively the output chosen will be Q_c . The competitive output yields no supernormal profit, whereas the monopoly/cartel output does.

Figure 10.14 Cartelizing a competitive industry



A cartel is formed when individual suppliers come together and act like a monopolist in order to increase profit. If MC is the joint supply curve of the cartel, profits are maximized at the output Q_m , where $MC=MR$. In contrast, if these firms operate competitively output increases to Q_c .

The cartel results in a deadweight loss equal to the area ABF , just as in the standard monopoly model.

Cartel instability

Some cartels are unstable in the long run. In the first instance, the degree of instability depends on the authority that the governing body of the cartel can exercise over its members, and upon the degree of information it has on the operations of its members. If a cartel is simply an arrangement among producers to limit output, each individual member of the cartel has an incentive to increase its output, because the monopoly price that the cartel attempts to sustain exceeds the cost of producing a marginal unit of output. In Figure 10.14 each firm has a MC of output equal to $\$F$ when the group collectively produces the output Q_m . Yet any firm that brings output to market, *beyond its agreed production limit*, at the price P_m will make a profit of AF on that additional output – *provided the other members of the cartel agree to restrict their output*. Since each firm faces the same incentive to increase output, it is difficult to restrain all members from doing so.

Individual members are more likely to abide by the cartel rules if the organization can sanction them for breaking the supply-restriction agreement. Alternatively, if the actions of individual members are not observable by the organization, then the incentive to break ranks may be too strong for the cartel to sustain its monopoly power.

We will see in Chapter 14 that Canada's Competition Act forbids the formation of cartels, as it forbids many other anti-competitive practices. At the same time, our governments frequently are the driving force in the formation of domestic cartels.

In the second instance, cartels may be undermined eventually by the emergence of new products and new technologies. OPEC has lost much of its power in the modern era because of technological developments in oil recovery. Canada's 'tar sands' yield oil, as a result of technological developments that enabled producers to separate the oil from the earth it is mixed with. Fracking technologies are another means of extracting oil that is discovered in small pockets and encased in rock. The supply coming from these new technologies has limited the ability of the old OPEC cartel to increase prices through supply restriction.

Application Box 10.1 The taxi cartel

The new sharing economy has brought competition to some traditional cartels. City taxis are an example of such a formation: Traditionally, entry has been restricted to drivers who hold a permit (medallion), and fares are higher as a consequence of the resulting reduced supply. A secondary market then develops for these medallions, in which the city may offer new medallions through auction, or existing owners may exit and sell their medallions. Restricted entry has characterized most of Canada's major cities. Depending on the strictness of the entry process, medallions are worth correspondingly more. By 2012, medallions were selling in New York and Boston for a price in the neighborhood of one million dollars.

But ride-sharing start-up companies changed all of that. As Western examples, Uber and Lyft developed smart-phone apps that link demanders for rides with drivers, who may, or may not be, part of the traditional taxi companies. Such start-ups have succeeded in taking a significant part of the taxi business away from the traditional operators. As a result, the price of taxi medallions on the open market has plunged. From trading in the range of \$1m. in New York in 2012, medallions are being offered in 2019 at about one fifth of that price. In Toronto, some medallions were traded in the range of \$300,000 in 2012, but are on offer in 2019 for prices in the range of \$30,000.

Not surprisingly, the traditional taxi companies charge that ride-hailing operators are violating the accepted rules governing the taxi business, and have launched legal suits against them and against local governments, and lobbied governments to keep them out of their cities.

In the new 'sharing economy', of which ride hailing companies are an example, participants operate with less traditional capital, and the communications revolution has been critical to their success. Home owners can use an online site to rent a spare bedroom in their house to visitors to their city (*Airbnb*), and thus compete with hotels. The main capital in this business is in the form of the information technology that links potential buyers to potential sellers.

Information on medallion prices in Canada can be found by, for example, searching at <http://www.kijiji.ca>

This page titled [10.6: Cartels- Acting like a monopolist](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine](#) (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.7: Invention, innovation and rent seeking

Invention and innovation are critical aspects of the modern economy. In some sectors of the economy, firms that cannot invent or innovate are liable to die. Invention is a genuine discovery, whereas innovation is the introduction of a new product or process.

Invention is the discovery of a new product or process through research.

Product innovation refers to new or better goods or services.

Process innovation refers to new or better production or supply.

To this point we have said little that is good about monopolies. However, the economist Joseph Schumpeter argued that, while monopoly leads to resource misallocation in the economy, this cost might be offset by the greater tendency for monopoly firms to invent and innovate. This is because such firms have more profit and therefore more resources with which to fund R&D and may therefore be more innovative than competitive firms. If this were true then, taking a long-run dynamic view of the marketplace, monopolies could have lower costs and more advanced products than competitive firms and thus benefit the consumer.

While this argument has some logical appeal, it falls short on several counts. First, even if large firms carry out more research than competitive firms, there is no guarantee that the ensuing benefits carry over to the consumer. Second, the results of such research may be used to prevent entry into the industry in question. Firms may register their inventions and gain use protection before a competitor can come up with the same or a similar invention. *Apple* and *Samsung* each own tens of thousands of patents. Third, the empirical evidence on the location of most R&D is inconclusive: A sector with several large firms, rather than one with a single or very many firms, may be best. For example, if *Apple* did not have *Samsung* as a competitor, or vice versa, would the pace of innovation be as strong?

Fourth, much research has a 'public good' aspect to it. Research carried out at universities and government-funded laboratories is sometimes referred to as basic research: It explores the principles underlying chemistry, social relations, engineering forces, microbiology, etc., and has multiple applications in the commercial world. If disseminated, this research is like a public good – its fruits can be used in many different applications, and its use in one area does not preclude its use in others. Consequently, rather than protecting monopolies on the promise of more R&D, a superior government policy might be to invest directly in research and make the fruits of the research publicly available.

Modern economies have patent laws, which grant inventors a legal monopoly on use for a fixed period of time – perhaps fifteen years. By preventing imitation, patent laws raise the incentive to conduct R&D but do not establish a monopoly in the long run. Over the life of a patent the inventor charges a higher price than would exist if his invention were not protected; this both yields greater profits and provides the research incentive. When the patent expires, competition from other producers leads to higher output and lower prices for the product. Generic drugs are a good example of this phenomenon.

Patent laws grant inventors a legal monopoly on use for a fixed period of time.

The power of globalization once again is very relevant in patents. Not all countries have patent laws that are as strong as those in North America and Europe. The BRIC economies (Brazil, Russia, India and China) form an emerging power block. But their legal systems and enforcement systems are less well-developed than in Europe or North America. The absence of a strong and transparent legal structure inhibits research and development, because their fruits may be appropriated by competitors.

Rent seeking

Citizens are frequently appalled when they read of lobbying activities in their nation's capital. Every capital city in the world has an army of lobbyists, seeking to influence legislators and regulators. Such individuals are in the business of rent seeking, whose goal is to direct profit to particular groups, and protect that profit from the forces of competition. In Virginia and Kentucky we find that state taxes on cigarettes are the lowest in the US – because the tobacco leaf is grown in these states, and the tobacco industry makes major contributions to the campaigns of some political representatives.

Rent-seeking carries a resource cost: Imagine that we could outlaw the lobbying business and put these lobbyists to work producing goods and services in the economy instead. Their purpose is to maintain as much quasi-monopoly power in the hands of their clients as possible, and to ensure that the fruits of this effort go to those same clients. If this practice could be curtailed then the time and resources involved could be redirected to other productive ends.

Rent seeking is an activity that uses productive resources to redistribute rather than create output and value.

Industries in which rent seeking is most prevalent tend to be those in which the potential for economic profits is greatest – monopolies or near-monopolies. These, therefore, are the industries that allocate resources to the preservation of their protected status. We do not observe laundromat owners or shoe-repair businesses lobbying in Ottawa.

This page titled [10.7: Invention, innovation and rent seeking](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.8: Conclusion

We have now examined two extreme types of market structure – perfect competition and monopoly. While many sectors of the economy operate in a way that is close to the competitive paradigm, very few are pure monopolies in that they have no close substitute products. Even firms like *Microsoft*, or *De Beers*, that supply a huge percentage of the world market for their product would deny that they are monopolies and would argue that they are subject to strong competitive pressures from smaller or 'fringe' producers. As a result we must look upon the monopoly paradigm as a useful way of analyzing markets, rather than being an exact description of the world. Accordingly, our next task is to examine how sectors with a few, several or multiple suppliers act when pursuing the objective of profit maximization. Many different market structures define the real economy, and we will concentrate on a limited number of the more important structures in the next chapter.

This page titled [10.8: Conclusion](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.9: Key Terms

Monopolist: is the sole supplier of an industry's output, and therefore the industry and the firm are one and the same.

Natural monopoly: one where the *ATC* of producing any output declines with the scale of operation.

Marginal revenue is the change in total revenue due to selling one more unit of the good.

Average revenue is the price per unit sold.

Allocative inefficiency arises when resources are not appropriately allocated and result in deadweight losses.

Price discrimination involves charging different prices to different consumers in order to increase profit.

A **cartel** is a group of suppliers that colludes to operate like a monopolist.

Rent seeking is an activity that uses productive resources to redistribute rather than create output and value.

Invention is the discovery of a new product or process through research.

Product innovation refers to new or better products or services.

Process innovation refers to new or better production or supply.

Patent laws grant inventors a legal monopoly on use for a fixed period of time.

This page titled [10.9: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

10.10: Exercises for Chapter 10

EXERCISE 10.1

Consider a monopolist with demand curve defined by $P=100-2Q$. The MR curve is $MR=100-4Q$ and the marginal cost is $MC=10+Q$. The demand intercepts are $\{ \$100, 50 \}$, the MR intercepts are $\{ \$100, 25 \}$.

- Develop a diagram that illustrates this market, using either graph paper or an Excel spreadsheet, for values of output $Q = 1 \dots 25$.
- Identify visually the profit-maximizing price and output combination.
- Optional:* Compute the profit maximizing price and output combination.

EXERCISE 10.2

Consider a monopolist who wants to maximize revenue rather than profit. She has the demand curve $P=72-Q$, with marginal revenue $MR=72-2Q$, and $MC=12$. The demand intercepts are $\{ \$72, 72 \}$, the MR intercepts are $\{ \$72, 36 \}$.

- Graph the three functions, using either graph paper or an Excel spreadsheet.
- Calculate the price she should charge in order to maximize revenue. [*Hint:* Where the $MR=0$.]
- Compute the total revenue she will obtain using this strategy.

EXERCISE 10.3

Suppose that the monopoly in Exercise 10.2 has a large number of plants. Consider what could happen if each of these plants became a separate firm, and acted competitively. In this perfectly competitive world you can assume that the MC curve of the monopolist becomes the industry supply curve.

- Illustrate graphically the output that would be produced in the industry?
- What price would be charged in the marketplace?
- Optional:* Compute the gain to the economy in dollar terms as a result of the DWL being eliminated [*Hint:* It resembles the area ABF in Figure 10.14].

EXERCISE 10.4

In the text example in Table 10.1, compute the profit that the monopolist would make if he were able to price discriminate, by selling each unit at the demand price in the market.

EXERCISE 10.5

A monopolist is able to discriminate perfectly among his consumers – by charging a different price to each one. The market demand curve facing him is given by $P=72-Q$. His marginal cost is given by $MC=24$ and marginal revenue is $MR=72-2Q$.

- In a diagram, illustrate the profit-maximizing equilibrium, where discrimination is not practiced. The demand intercepts are $\{ \$72, 72 \}$, the MR intercepts are $\{ \$72, 36 \}$.
- Illustrate the equilibrium output if he discriminates perfectly.
- Optional:* If he has no fixed cost beyond the marginal production cost of \$24 per unit, calculate his profit in each pricing scenario.

EXERCISE 10.6

A monopolist faces two distinct markets A and B for her product, and she is able to insure that resale is not possible. The demand curves in these markets are given by $P_A=20-(1/4)Q_A$ and $P_B=14-(1/4)Q_B$. The marginal cost is constant: $MC=4$. There are no fixed costs.

- Graph these two markets and illustrate the profit maximizing price and quantity in each market. [You will need to insert the MR curves to determine the optimal output.] The demand intercepts in A are $\{ \$20, 80 \}$, and in B are $\{ \$14, 56 \}$.
- In which market will the monopolist charge a higher price?

EXERCISE 10.7

A concert organizer is preparing for the arrival of the Grateful Living band in his small town. He knows he has two types of concert goers: One group of 40 people, each willing to spend \$60 on the concert, and another group of 70 people, each willing to spend

\$40. His total costs are purely fixed at \$3,500.

- a. Draw the market demand curve faced by this monopolist.
- b. Draw the MR and MC curves.
- c. With two-price discrimination what will be the monopolist's profit?
- d. If he must charge a single price for all tickets can he make a profit?

EXERCISE 10.8

Optional: A monopolist faces a demand curve $P=64-2Q$ and $MR=64-4Q$. His marginal cost is $MC=16$.

- a. Graph the three functions and compute the profit maximizing output and price.
 - b. Compute the efficient level of output (where $MC=demand$), and compute the DWL associated with producing the profit maximizing output rather than the efficient output.
1. It is not difficult to show that the demand curve corresponding to the data in the table is given by the expression $P=12-0.5q$. Since the MR curve has twice the slope of the demand curve, it is given by $MR=12-q$. The data indicate that the MC curve can be written as $MC=0.5q$ and the average cost curve by $AC=0.25q$. Setting $MR=MC$ yields $q=8$. At this output the demand curve implies the price is \$8. Profit is $(P-AC) \times q = (\$8 - \$2) \times 8 = \$48$.
 2. We can think of such a transformation coming from a single supplier taking over a number of small suppliers, and the monopolist would thus be a multi-plant firm.

This page titled [10.10: Exercises for Chapter 10](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11: Imperfect competition

Chapter 11: Imperfect competition

In this chapter we will explore:

11.1	The principle ideas
11.2	Imperfect competitors
11.3	Imperfect competitors: measures of structure and market power
11.4	Imperfect competition: monopolistic competition
11.5	Imperfect competition: economies of scope and platforms
11.6	Strategic behaviour: oligopoly and games
11.7	Strategic behaviour: duopoly and Cournot games
11.8	Strategic behaviour: entry, exit and potential competition
11.9	Matching markets: design

11.1 The principle ideas

The preceding chapters have explored extreme forms of supply: The monopolist is the sole supplier and possesses as much market power as possible. In contrast, the perfect competitor is small and has no market power whatsoever. He simply accepts the price for his product that is determined in the market by the forces of supply and demand. These are very useful paradigms to explore, but the real world for the most part lies between these extremes. We observe that there are a handful of dominant brewers in Canada who supply more than three quarters of the market, and they are accompanied by numerous micro brewers that form the fringe of the brewing business. We have a small number of air carriers and one of them controls half of the national market. The communications market has just three major suppliers; the Canadian Football League has nine teams and there are just a handful of major hardware/builders' suppliers stores nationally. At the other end of the spectrum we have countless restaurants and fitness centres, but they do not supply exactly the same product to the marketplaces for 'food' or 'health', and so these markets are not perfectly competitive, despite the enormous number of participants.

In this chapter we will explore three broad topics: First is the relationship between firm behaviour and firm size relative to the whole sector. This comes broadly under the heading of *imperfect competition* and covers a variety of market forms. Second, we will explore the principle modern ideas in *strategic behavior*. In a sense all decisions in microeconomics have an element of strategy to them - economic agents aim to attain certain goals and they adopt specific maximizing strategies to attain them. But in this chapter we explore a more specific concept of strategic behavior - one that focuses upon direct interactions between a small number of players in the market place. Third, we explore the principle characteristics of what are termed *matching' markets*. These are markets where transactions take place without money and involve matching heterogeneous suppliers with heterogeneous buyers.

11.2 Imperfect competitors

Imperfect competitors can be defined by the number of firms in their sector, or the share of total sales going to a small number of suppliers. They can also be defined in terms of the characteristics of the demand curves they all face. A perfect competitor faces a perfectly elastic demand at the existing market price, and this is the only market structure to have this characteristic. In all other market structures suppliers effectively face a downward-sloping demand. This means that they have some influence on the price of the good, and also that if they change the price they charge, they can expect demand to reflect this in a predictable manner. So, in theory, we can classify all market structures apart from perfect competition as being imperfectly competitive. In practice we use the term to denote firms that fall between the extremes of perfect competition and monopoly.

Imperfectly competitive firms face a downward-sloping demand curve, and their output price reflects the quantity sold.

The demand curve for the firm and industry coincide for the monopolist, but not for other imperfectly competitive firms. It is convenient to categorize the producing sectors of the economy as either having a relatively small number of participants, or having

a large number. The former market structures are called oligopolistic, and the latter are called monopolistically competitive. The word *oligopoly* comes from the Greek word *oligos* meaning few, and *polein* meaning to sell.

Oligopoly defines a market with a small number of suppliers.

Monopolistic competition defines a market with many sellers of products that have similar characteristics. Monopolistically competitive firms can exert only a small influence on the whole market.

The home appliance industry is an oligopoly. The prices of *KitchenAid* appliances depend not only on their own output and sales, but also on the prices of *Whirlpool*, *Maytag* and *Bosch*. If a firm has just two main producers it is called a duopoly. *Canadian National* and *Canadian Pacific* are the only two major rail freight carriers in Canada; they thus form a duopoly. In contrast, the local Italian restaurant is a monopolistic competitor. Its output is a package of distinctive menu choices, personal service, and convenience for local customers. It can charge a different price than the out-of-neighbourhood restaurant, but if its prices are too high local diners may travel elsewhere for their food experience, or switch to a different cuisine locally. Many markets are defined by producers who supply similar but not identical products. Canada's universities all provide degrees, but they differ one from another in their programs, their balance of in-class and on-line courses, their student activities, whether they are science based or liberal arts based, whether they have cooperative programs or not, and so forth. While universities are not in the business of making profit, they certainly wish to attract students, and one way of doing this is to differentiate themselves from other institutions. The profit-oriented world of commerce likewise seeks to increase its market share by distinguishing its product line.

Duopoly defines a market or sector with just two firms.

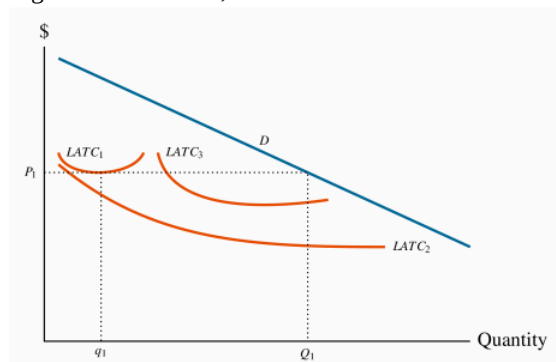
These distinctions are not completely airtight. For example, if a sole domestic producer is subject to international competition it cannot act in the way we described in the previous chapter – it has potential, or actual, competition. *Bombardier* may be Canada's sole rail car manufacturer, but it is not a monopolist, even in Canada. It could best be described as being part of an international oligopoly in rail-car manufacture. Likewise, it is frequently difficult to delineate the boundary of a given market. For example, is *Canada Post* a monopoly in mail delivery, or an oligopolist in hard-copy communication? We can never fully remove these ambiguities.

The role of cost structures

A critical determinant of market structure is the way in which demand and cost interact to determine the likely number of market participants in a given sector or market. Structure also evolves over the long run: Time is required for entry and exit.

Figure 11.1 shows the demand curve D for the output of an industry in the long run. Suppose, initially, that all firms and potential entrants face the long-run average cost curve $LATC_1$. At the price P_1 , free entry and exit means that each firm produces q_1 . With the demand curve D , industry output is Q_1 . The number of firms in the industry is $N_1 (=Q_1/q_1)$. If q_1 , the minimum average cost output on $LATC_1$, is small relative to D , then N_1 is large. This outcome might be perfect competition (N virtually infinite), or monopolistic competition (N large) with slightly differentiated products produced by each firm.

Figure 11.1 Demand, costs and market structure



With a cost structure defined by $LATC_1$ this market has space for many firms – perfect or monopolistic competition, each producing approximately q_1 . If costs correspond to $LATC_2$, where scale economies are substantial, there may be space for just one producer. The intermediate case, $LATC_3$, can give rise to oligopoly, with each firm producing more than q_1 but less than a monopolist. These curves encounter their MES at very different output levels.

Instead, suppose that the production structure in the industry is such that the long-run average cost curve is $LATC_2$. Here, scale economies are vast, relative to the market size. At the lowest point on this cost curve, output is large relative to the demand curve

D. If this one firm were to act like a monopolist it would produce an output where $MR=MC$ in the long run and set a price such that the chosen output is sold. Given the scale economies, there may be no scope for another firm to enter this market, because such a firm would have to produce a very high output to compete with the existing producer. This situation is what we previously called a "natural" monopolist.

Finally, the cost structure might involve curves of the type $LATC_3$, which would give rise to the possibility of several producers, rather than one or very many. This results in oligopoly.

It is clear that one *crucial determinant of market structure is minimum efficient scale relative to the size of the total market* as shown by the demand curve. The larger the minimum efficient scale relative to market size, the smaller is the number of producers in the industry.

11.3 Imperfect competitors: measures of structure and market power

Sectors of the economy do not fit neatly into the limited number of categories described above. The best we can say in most cases is that they resemble more closely one type of market than another. Consider the example of Canada's brewing sector: It has two large brewers in *Molson-Coors* and *Labatt*, a couple of intermediate sized firms such as *Sleeman*, and an uncountable number of small boutique brew pubs. While such a large number of brewers satisfy one requirement for perfect competition, it would not be true to say that the biggest brewers wield no market power; and this is the most critical element in defining market structure.

By the same token, we could not define this market as a duopoly: Even though there are just two major participants, there are countless others who, together, are important.

One way of defining what a particular structure most closely resembles is to examine the percentage of sales in the market that is attributable to a small number of firms. For example: What share is attributable to the largest three or four firms? The larger the share, the more concentrated the market power. Such a statistic is called a concentration ratio. The N -firm concentration ratio is the sales share of the largest N firms in that sector of the economy.

The **N -firm concentration ratio** is the sales share of the largest N firms in that sector of the economy.

Table 11.1 Concentration in Canadian food processing 2011

Sector	% of shipments
Sugar	98
Breakfast cereal	96
Canning	60
Meat processing	23

Source: "Four Firm Concentration Ratios (CR4s) for selected food processing sectors," adapted from Statistics Canada publication Measuring industry concentration in Canada's food processing sectors, Agriculture and Rural Working Paper series no. 70, Catalogue 21-601, <http://www.statcan.gc.ca/pub/21-601-m/21-601-m2004070-eng.pdf>.

Table 11.1 contains information on the 4-firm concentration ratio for several sectors of the Canadian economy. It indicates that, at one extreme, sectors such as breakfast cereals and sugars have a high degree of concentration, whereas meat processing has much less. A high degree of concentration suggests market power, and possibly economies of scale.

11.4 Imperfect competition: monopolistic competition

Monopolistic competition presumes a large number of quite small producers or suppliers, each of whom may have a slightly differentiated product. The competition element of this name signifies that there are many participants, while the monopoly component signifies that each supplier faces a downward-sloping demand. In concrete terms, your local coffee shop that serves "fair trade" coffee has a product that differs slightly from that of neighbouring shops that sell the traditional product. They coexist in the same sector, and probably charge different prices: The fair trade supplier likely charges a higher price, but knows nonetheless that too large a difference between her price and the prices of her competitors will see some of her clientele migrate to those lower-priced establishments. That is to say, she faces a downward-sloping demand curve.

The competition part of the name also indicates that there is *free entry and exit*. There are no barriers to entry. As a consequence, we know at the outset that only normal profits will exist in a long-run equilibrium. Economic profits will be competed away by entry, just as losses will erode due to exit.

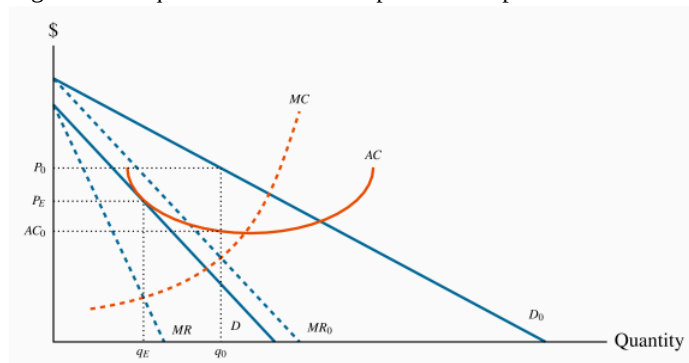
As a general rule then, each firm can influence its market share to some extent by changing its price. Its demand curve is not horizontal because different firms' products are only limited substitutes. A lower price level may draw some new customers away from competitors, but convenience or taste will prevent most patrons from deserting their local businesses. In concrete terms: A pasta special at the local Italian restaurant that reduces the price below the corresponding price at the competing local Thai restaurant will indeed draw clients away from the latter, but the foods are sufficiently different that only some customers will leave the Thai restaurant. The differentiated menus mean that many customers will continue to pay the higher price.

A **differentiated product** is one that differs slightly from other products in the same market.

Given that there are very many firms, the theory also envisages limits to scale economies. Firms are small and, with many competitors, individual firms do not compete strategically with *particular* rivals. Because the various products offered are slightly differentiated, we avoid graphics with a *market* demand, because this would imply that a uniform product is being considered. At the same time the market is a well-defined concept—it might be composed of all those restaurants within a reasonable distance, for example, even though each one is slightly different from the others. The market share of each firm depends on the price that it charges *and* on the number of competing firms. For a given number of suppliers, a shift in industry demand also shifts the demand facing each firm. Likewise, the presence of more firms in the industry reduces the demand facing each one.

Equilibrium is illustrated in Figure 11.2. Here D_0 is the initial demand facing a representative firm, and MR_0 is the corresponding marginal revenue curve. Profit is maximized where $MC=MR$, and the price P_0 is obtained from the demand curve corresponding to the output q_0 . Total profit is the product of output times the difference between price and average cost, which equals $q_0 \times (P_0 - AC_0)$.

Figure 11.2 Equilibrium for a monopolistic competitor



Profits exist at the initial equilibrium (q_0, P_0) . Hence, new firms enter and reduce the share of the total market faced by each firm, thereby shifting back their demand curve. A final equilibrium is reached where economic profits are eliminated: At $AC=P_E$ and $MR=MC$.

With free entry, such profits attract new firms. The increased number of firms reduces the share of the market that any one firm can claim. That is, the firm's demand curve shifts inwards when entry occurs. As long as (economic) profits exist, this process continues. For entry to cease, average cost must equal price. A final equilibrium is illustrated by the combination (P_E, q_E) , where the demand has shifted inward to D .

At this long-run equilibrium, two conditions must hold: First, the optimal pricing rule must be satisfied—that is $MC=MR$; second it must be the case that only normal profits are made at the final equilibrium. Economic profits are competed away as a result of free entry. Graphically this implies that ATC must equal price at the output where $MC=MR$. In turn this implies that the ATC is tangent to the demand curve where $P=ATC$. While this could be proven mathematically, it is easy to intuit why this tangency must exist: If ATC merely intersected the demand curve at the output where $MC=MR$, we could find some other output where the demand price would be above ATC , suggesting that profits could be made at such an output. Clearly that could not represent an equilibrium.

The **monopolistically competitive equilibrium** in the long run requires the firm's demand curve to be tangent to the ATC curve at the output where $MR=MC$.

11.5 Imperfect competition: economies of scope and platforms

The communications revolution has impacted market structure in modern economies profoundly: it has facilitated economies of scope, meaning that firms may yield more collective profit if merged than if operating independently.

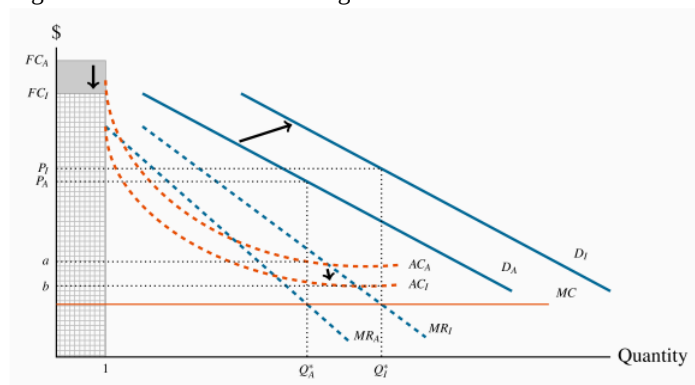
Economies of Scope

Imagine an aspiring entrepreneur who envisages a revolution of the traditional taxi sector of the economy. He decides to develop a smartphone application that will match independent income-seeking vehicle owners (drivers) with individuals seeking transport (passengers) from point A to point B. We know how this adventure evolves. In one case it takes the form of the corporation Uber, in another the corporation Lyft, and others worldwide.

These corporations have grown in leaps and bounds and have taken business from the conventional taxi corporations. As of 2019 they cannot turn a profit, yet the stock market continues to bet upon future success: investors believe that when these corporations evolve into fully integrated multi-product suppliers, both costs will decline and demand will increase for each component of the business. In the case of transportation companies, they aim to become a 'one-stop-shop' for mobility services. Uber is not only a ride-hailing service, it also transports meals through its Uber-eats platform, and is developing the electric scooter and electric bike markets in addition. In some local markets it is linked to public transport services. All of this is being achieved through a single smartphone application. The objective is to simplify movement for persons, by providing multiple options on a variety of transport modes, accessed through a single portal.

This phenomenon is described in Figure 11.3. The subscripts *A* and *I* represent market conditions when the service supplier is operating Alone or in an Integrated corporation. The initial equilibrium is defined by the *A* demand and cost conditions. The profit maximizing output occurs when $MC_A = MR_A$, leading to a price P_A and a quantity Q_A^* . Each unit of the good yields a profit margin of $(P_A - a)$.

Figure 11.3 Summa's ride hailing service



Demand for a particular product increases when the autonomous supplier (*A*) merges with another firm to become an integrated firm (*I*), because customers switch to firms that offer several different services from the same platform: the demand curve shifts outward, from D_A to D_I . With integration, the fixed costs fall and average costs fall, even with marginal costs constant. Output and profit increase, and concentration in the marketplace rises.

This firm now merges with another transportation corporation - perhaps a food delivery service, perhaps an electric bike service. Since each firm has a similar type of fixed cost, these costs can be reduced by the merger. In technical terms, the merged firms, or merged operations, share a common hardware-cum-software platform. Each firm will therefore incur lower average costs, even if marginal costs remain unchanged: the *AC* curve declines to AC_I . In addition to the decline in average costs, each firm sees an increase in its customer base, because transportation service buyers find it preferable to choose their mode of transport through a single portal rather than through several different modes of access. This is represented by an outward shift in the demand curve for vehicle rides to D_I .

The new profit maximizing equilibrium occurs at Q_I^* . Total profit necessarily increases both because average costs have fallen and the number of buyers willing to buy at any price has risen. The analytics in this figure also describe the benefits accruing to the other firm or firms in the merger.

A **platform** describes a technology that is common to more than one product in a multi-product organization.

We conclude from this analysis that, if scope economies are substantial, it may be difficult for stand-alone firms specializing in just one component of the transportation services sector to remain profitable. It may also be impossible to define a conventional equilibrium in this kind of marketplace. This is because some conglomerate firms may have different component producers in their suite of firms. For example, Lyft may not have a food delivery service, but it may have a limousine or bus service. What is critical

for an equilibrium is that firms of a particular type, whether they are part of a conglomerate or not, be able to compete with corresponding firms. This means that their cost structure must be similar.

As a further example: Amazon initially was primarily an on-line book seller. But it expanded to include the sale of other products. And once it became a 'market for everything' the demand side of the market exploded in parallel with the product line, because it becomes easy to shop for 'anything' or even different objects on a single site. Only Walmart, in North America, comes close to being able to compete with Amazon.

11.6 Strategic behaviour: Oligopoly and games

Under perfect competition or monopolistic competition, there are so many firms in the industry that each one can ignore the immediate effect of its own actions on particular rivals. However, in an oligopolistic industry *each firm must consider how its actions affect the decisions of its relatively few competitors*. Each firm must guess how its rivals will react. Before discussing what constitutes an intelligent guess, we investigate whether they are likely to collude or compete. Collusion is a means of reducing competition with a view to increasing profit.

Collusion is an explicit or implicit agreement to avoid competition with a view to increasing profit.

A particular form of collusion occurs when firms co-operate to form a cartel, as we saw in the last chapter. Collusion is more difficult if there are many firms in the industry, if the product is not standardized, or if demand and cost conditions are changing rapidly. In the absence of collusion, each firm's demand curve depends upon how competitors react: If Air Canada contemplates offering customers a seat sale on a particular route, how will West Jet react? Will it, too, make the same offer to buyers? If Air Canada thinks about West Jet's likely reaction, will it go ahead with the contemplated promotion? A conjecture is a belief that one firm forms about the strategic reaction of another competing firm.

A **conjecture** is a belief that one firm forms about the strategic reaction of another competing firm.

Good poker players will attempt to anticipate their opponents' moves or reactions. Oligopolists are like poker players, in that they try to anticipate their rivals' moves. To study interdependent decision making, we use game theory. A game is a situation in which contestants plan strategically to maximize their payoffs, taking account of rivals' behaviour.

A **game** is a situation in which contestants plan strategically to maximize their payoffs, taking account of rivals' behaviour.

The *players* in the game try to maximize their own *payoffs*. In an oligopoly, the firms are the players and their payoffs are their profits. Each player must choose a strategy, which is a plan describing how a player moves or acts in different situations.

A **strategy** is a game plan describing how a player acts, or moves, in each possible situation.

Equilibrium outcomes

How do we arrive at an equilibrium in these games? Let us begin by defining a commonly used concept of equilibrium. A Nash equilibrium is one in which each player chooses the best strategy, given the strategies chosen by the other players, and there is no incentive to move or change choice.

A **Nash equilibrium** is one in which each player chooses the best strategy, given the strategies chosen by the other player, and there is no incentive for any player to move.

In such an equilibrium, no player wants to change strategy, since the other players' strategies were already figured into determining each player's own best strategy. This concept and theory are attributable to the Princeton mathematician John Nash, who was popularized by the Hollywood movie version of his life, *A Beautiful Mind*.

In most games, each player's best strategy depends on the strategies chosen by their opponents. Occasionally, a player's best strategy is independent of those chosen by rivals. Such a strategy is called a dominant strategy.

A **dominant strategy** is a player's best strategy, independent of the strategies adopted by rivals.

We now illustrate these concepts with the help of two different games. These games differ in their outcomes and strategies. Table 11.2 contains the domestic happiness game¹. Will and Kate are attempting to live in harmony, and their happiness depends upon each of them carrying out domestic chores such as shopping, cleaning and cooking. The first element in each pair defines Will's outcome, the second Kate's outcome. If both contribute to domestic life they each receive a happiness or utility level of 5 units. If one contributes and the other does not the happiness levels are 2 for the contributor and 6 for the non-contributor, or 'free-rider'. If neither contributes happiness levels are 3 each. When each follows the same strategy the payoffs are on the diagonal, when they

follow different strategies the payoffs are on the off-diagonal. Since the elements of the table define the payoffs resulting from various choices, this type of matrix is called a payoff matrix.

A **payoff matrix** defines the rewards to each player resulting from particular choices.

So how is the game likely to unfold? In response to Will's choice of a contribute strategy, Kate's utility maximizing choice involves lazing: She gets 6 units by not contributing as opposed to 5 by contributing. Instead, if Will decides to be lazy what is in Kate's best interest? Clearly it is to be lazy also because that strategy yields 3 units of happiness compared to 2 units if she contributes. In sum, Kate's best strategy is to be lazy, regardless of Will's behaviour. So the strategy of not contributing is a *dominant strategy*, in this particular game.

Will also has a dominant strategy – identical to Kate's. This is not surprising since the payoffs are symmetric in the table. Hence, since each has a dominant strategy of not contributing the Nash equilibrium is in the bottom right cell, where each receives a payoff of 3 units. Interestingly, this equilibrium is not the one that yields maximum combined happiness.

Table 11.2 A game with dominant strategies

		Kate's choice	
		Contribute	Laze
Will's choice	Contribute	5,5	2,6
	Laze	6,2	3,3

The first element in each cell denotes the payoff or utility to Will; the second element the utility to Kate.

The reason that the equilibrium yields less utility for each player in this game is that the game is competitive: Each player tends to their own interest and seeks the best outcome conditional on the choice of the other player. This is evident from the (5,5) combination. From this position Kate would do better to defect to the Laze strategy, because her utility would increase².

To summarize: This game has a unique equilibrium and each player has a dominant strategy. But let us change the payoffs just slightly to the values in Table 11.3. The off-diagonal elements have changed. The contributor now gets no utility as a result of his or her contributions: Even though the household is a better place, he or she may be so annoyed with the other person that no utility flows to the contributor.

Table 11.3 A game without dominant strategies

		Kate's choice	
		Contribute	Laze
Will's choice	Contribute	5,5	0,4
	Laze	4,0	3,3

The first element in each cell denotes the payoff or utility to Will; the second element the utility to Kate.

What are the optimal choices here? Starting again from Will choosing to contribute, what is Kate's best strategy? It is to contribute: She gets 5 units from contributing and 4 from lazing, hence she is better contributing. But what is her best strategy if Will decides to laze? It is to laze, because that yields her 3 units as opposed to 0 by contributing. This set of payoffs therefore *contains no dominant strategy* for either player.

As a result of there being no dominant strategy, there arises the possibility of more than one equilibrium outcome. In fact there are two equilibria in this game now: If the players find themselves both contributing and obtaining a utility level of (5,5) it would not be sensible for either one to defect to a laze option. For example, if Kate decided to laze she would obtain a payoff of 4 utils rather than the 5 she enjoys at the (5,5) equilibrium. By the same reasoning, if they find themselves at the (laze, laze) combination there is no incentive to move to a contribute strategy.

Once again, it is to be emphasized that the twin equilibria emerge in a competitive environment. If this game involved cooperation or collusion the players should be able to reach the (5,5) equilibrium rather than the (3,3) equilibrium. But in the competitive environment we cannot say *ex ante* which equilibrium will be attained.

Repeated games

This game illustrates the tension between collusion and competition. While we have developed the game in the context of the household, it can equally be interpreted in the context of a profit maximizing game between two market competitors. Suppose the numbers define profit levels rather than utility as in Table 11.4. The 'contribute' option can be interpreted as 'cooperate' or 'collude', as we described for a cartel in the previous chapter. They collude by agreeing to restrict output, sell that restricted output at a higher price, and in turn make a greater total profit which they split between themselves. The combined best profit outcome (5,5) arises when each firm restricts its output.

Table 11.4 Collusion possibilities

		Firm K's profit	
		Low output	High output
Firm W's profit	Low output	5,5	2,6
	High output	6,2	3,3

The first element in each cell denotes the profit to Firm W; the second element the profit to Firm K.

But again there arises an incentive to defect: If Firm W agrees to maintain a high price and restrict output, then Firm K has an incentive to renege and increase output, hoping to improve its profit through the willingness of Firm W to restrict output. Since the game is symmetric, each firm has an incentive to renege. Each firm has a dominant strategy – high output, and there is a unique equilibrium (3,3).

Obviously there arises the question of whether these firms can find an operating mechanism that would ensure they each generate a profit of 5 units rather than 3 units, while remaining purely self-interested. This question brings us to the realm of repeated games. For example, suppose that firms make strategic choices each quarter of the year. If firm K had 'cheated' on the collusive strategy it had agreed with firm W in the previous quarter, what would happen in the following quarter? Would firms devise a strategy so that cheating would not be in the interest of either one, or would the competitive game just disintegrate into an unpredictable pattern? These are interesting questions and have provoked a great deal of thought among game theorists. But they are beyond our scope at the present time.

A repeated game is one that is repeated in successive time periods and where the knowledge that the game will be repeated influences the choices and outcomes in earlier periods.

We now examine what might happen in one-shot games of the type we have been examining, but in the context of many possible choices. In particular, instead of assuming that each firm can choose a high or low output, how would the outcome of the game be determined if each firm can choose an output that can lie anywhere between a high and low output? In terms of the demand curve for the market, this means that the firms can choose some output and price that is consistent with demand conditions: There may be an infinite number of choices. This framing of a game enables us to explore new concepts in strategic behavior.

11.7 Strategic behaviour: Duopoly and Cournot games

The duopoly model that we frequently use in economics to analyze competition between a small number of competitors is fashioned after the ideas of French economist Augustin Cournot. Consequently it has come to be known as the Cournot duopoly model. While the maximizing behaviour that is incorporated in this model can apply to a situation with several firms rather than two, we will develop the model with two firms. This differs slightly from the preceding section, where each firm has simply a choice between a high or low output.

The critical element of the Cournot approach is that the firms each determine their optimal strategy – one that maximizes profit – by reacting optimally to their opponent's strategy, which in this case involves their choice of output.

Cournot behaviour involves each firm reacting optimally in their choice of output to their competitors' output decisions.

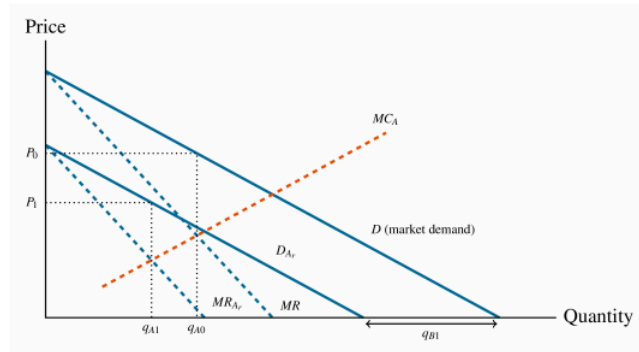
A central element here is the reaction function of each firm, which defines the optimal output choice conditional upon their opponent's choice.

Reaction functions define the optimal choice of output conditional upon a rival's output choice.

We can develop an optimal strategy with the help of Figure 11.4. D is the market demand, and two firms supply this market. If B supplies a zero output, then A would face the whole demand, and would maximize profit where $MC=MR$. Let this output be defined

by q_{A0} . We transfer this output combination to Figure 11.5, where the output of each firm is on one of the axes—A on the vertical axis and B on the horizontal. This particular combination of zero output for B and q_{A0} for A is represented on the vertical axis as the point q_{A0} .

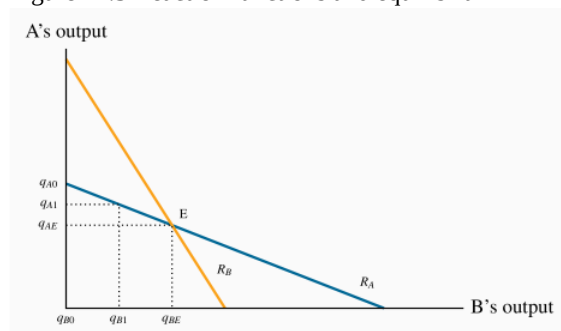
Figure 11.4 Duopoly behaviour



When one firm, B, chooses a specific output, e.g. q_{B1} , then A's residual demand D_{A_r} is the difference between the market demand and q_{B1} . A's profit is maximized at q_{A1} —where $MC = MR_{A_r}$. This is an optimal reaction by A to B's choice. For all possible choices by B, A can form a similar optimal response. The combination of these responses forms A's reaction function.

Instead, suppose that B produces a quantity q_{B1} in Figure 11.4. This reduces the demand curve facing A correspondingly from D to D_{A_r} , which we call A's *residual demand*. When subject to such a choice by B, firm A maximizes profit by producing where $MR_{A_r} = MC$, where MR_{A_r} is the marginal revenue corresponding to the residual demand D_{A_r} . The optimum for A is now q_{A1} , and this pair of outputs is represented by the combination (q_{A1}, q_{B1}) in Figure 11.5.

Figure 11.5 Reaction functions and equilibrium



The reaction function for A (R_A) defines the optimal output response for A to any output choice by B. The reaction function for B is defined similarly. The equilibrium occurs at the intersection of R_A and R_B . Any other combination will induce one firm to change its output, and therefore could not be an equilibrium.

Firm A forms a similar optimal response for every possible output level that B could choose, and these responses define A's reaction function. The reaction function illustrated for A in Figure 11.5 is thus the locus of all optimal response outputs on the part of A. The downward-sloping function makes sense: The more B produces, the smaller is the residual market for A, and therefore the less A will produce.

But A is just one of the players in the game. If B acts in the same optimizing fashion, B too can formulate a series of optimal reactions to A's output choices. The combination of such choices would yield a reaction function for B. This is plotted as R_B in Figure 11.5.

An equilibrium is defined by the intersection of the two reaction functions, in this case by the point E. At this output level *each firm is making an optimal decision, conditional upon the choice of its opponent*. Consequently, neither firm has an incentive to change its output; therefore it can be called the Nash equilibrium.

Any other combination of outputs on either reaction function would lead one of the players to change its output choice, and therefore could not constitute an equilibrium. To see this, suppose that B produces an output greater than q_{BE} ; how will A react? A's reaction function indicates that it should choose a quantity to supply less than q_{AE} . If so, how will B respond in turn to that optimal choice? It responds with a quantity read from its reaction function, and this will be less than the amount chosen at the previous stage. By tracing out such a sequence of reactions it is clear that the output of each firm will move to the equilibrium q_E .

Application Box 11.1 Cournot: Fixed costs and brand

Why do we observe so many industries on the national, and even international, stages with only a handful of firms? For example, Intel produces more than half of the world's computer chips, and AMD produces a significant part of the remainder. Why are there only two major commercial aircraft producers in world aviation – Boeing and Airbus? Why are there only a handful of major North American suppliers in pharmaceuticals, automobile tires, soda pop, internet search engines and wireless telecommunications?

The answer lies primarily in the nature of modern product development. Product development (fixed) costs, coupled with a relatively small marginal cost of production, leads to markets where there is enough space for only a few players. The development cost for a new cell phone, or a new aircraft, or a new computer-operating system may run into billions, while the cost of producing each unit may in fact be constant. The enormous development cost associated with many products explains not only why there may be a small number of firms in the domestic market for the product, but also why the number of firms in some sectors is small worldwide.

The Cournot model yields an outcome that lies between monopoly (or collusion/cartel) and competitive market models. It does not necessarily assume that the firms are identical in terms of their cost structure, although the lower-cost producer will end up with a larger share of the market.

The next question that arises is whether this duopoly market will be sustained as a duopoly, or if entry may take place. In particular, if economic profits accrue to the participants will such profits be competed away by the arrival of new producers, or might there be barriers of either a 'natural' or 'constructed' type that operate against new entrants?

11.8 Strategic behaviour: Entry, exit & potential competition

At this point we inquire about the potential entry and impact of new firms – firms who might enter the industry if conditions were sufficiently enticing, meaning the presence of economic profits. One way of examining entry in this oligopolistic world is to envisage potential entry barriers as being either intended or unintended, though the difference between the two can be blurred. Broadly, an unintended or 'natural' barrier is one related to scale economies and the size of the market. An intended barrier involves a strategic decision on the part of the firm to prevent entry.

Unintended entry barriers

Oligopolists tend to have substantial fixed costs, accompanied by declining average costs up to high output levels. Such a cost structure 'naturally' gives rise to a supply side with a small number of suppliers. For examples, given demand and cost structures, could Vancouver support two professional soccer teams; could Calgary support two professional hockey teams; could Montreal sustain two professional football teams? The answer to each of these questions is likely 'no'. Because given the cost structure of these markets, it would not be possible to induce twice as many spectators without reducing the price per game ticket to such a degree that revenue would be insufficient to cover costs. (We will neglect for the moment that the governing bodies of these sports also have the power to limit entry.) Fixed costs include stadium costs, staff payrolls and player payrolls. In fact most costs in these markets are relatively fixed. Market size relative to fixed and variable costs is not large enough to sustain two teams in most cities. Exceptions in reality are huge urban areas such as New York and Los Angeles.

Accordingly, it is possible that the existing team, or teams, may earn economic profit from their present operation; but such profit does not entice further entry, because the market structure is such that the entry of an additional team could lead to each team making losses.

Intended entry barriers

Patent Law

This is one form of protection for incumbent firms. Research and development is required for the development of many products in the modern era. Pharmaceuticals are an example. If innovations were not protected, firms and individuals would not be incentivized to devote their energies and resources to developing new drugs. Society would be poorer as a result. Patent protection is obviously a legal form of protection. At the same time, patent protection can be excessive. If patents provide immunity from replication or copying for an excessive period of time - for longer than required to recoup R & D costs - then social welfare declines because monopoly profits are being generated as a result of output restriction at too high a price.

Advertising

Advertising is a second form of entry deterrence. In this instance firms attempt to market their product as being distinctive and even enviable. For example, *Coca-Cola* and *PepsiCo* invest hundreds of millions annually to project their products in this light. They

sponsor sports, artistic and cultural events. Entry into the cola business is not impossible, but brand image is so strong for these firms that potential competitors would have a very low probability of entering this sector profitably. Likewise, in the 'energy-drinks' market, *Red Bull* spends hundreds of millions of dollars per annum on Formula One racing, kite surfing contests, mountain biking events and other extreme sports. In doing this it is reinforcing its brand image and distinguishing its product from Pepsi or Coca-Cola. This form of advertising is one of product differentiation and enables the manufacturer to maintain a higher price for its products by convincing its buyers that there are no close substitutes.

Predatory pricing

This form of pricing constitutes an illegal form of entry deterrence. It involves an incumbent charging an artificially low price for its product in the event of entry of a new competitor. This is done with a view to making it impossible for the entrant to earn a profit. Given that incumbents have generally greater resources than entrants, they can survive a battle of losses for a more prolonged period, thus ultimately driving out the entrant.

An iconic example of predatory pricing is that of Amazon deciding to take on a startup called Quidsi that operated the website *diapers.com*.³ The latter was proving to be a big hit with consumers in 2009 and Amazon decided that it was eating into Amazon profits on household and baby products. Amazon reacted by cutting its own prices dramatically, to the point where it was ready to lose a huge amount of money in order to grind Quidsi into the ground. The ultimate outcome was that Quidsi capitulated and sold to Amazon.

Whether this was a legal tactic or not we do not know, but it underlines the importance of war chests.

Maintaining a war chest

Many large corporations maintain a mountain of cash. This might seem like an odd thing to do when it could be paying that cash out to owners in the form of dividends. But there are at least two reasons for not doing this. First, personal taxes on dividends are frequently higher than taxes on capital gains; accordingly if a corporation can transform its cash into capital gain by making judicious investments, that strategy ultimately yields a higher post-tax return to the stock holders. A second reason is that a cash war chest serves as a credible threat to competitors of the type described involving Amazon and Quidsi above.

Network externalities

These externalities arise when the existing number of buyers itself influences the total demand for a product. *Facebook* is now a classic example. An individual contemplating joining a social network has an incentive to join one where she has many existing 'friends'. Not everyone views the *Microsoft* operating system (OS) as the best. Many prefer a simpler system such as *Linux* that also happens to be free. However, the fact that almost every new computer (that is not *Apple*) coming onto the market place uses Microsoft OS, there is an incentive for users to continue to use it because it is so easy to find a technician to repair a breakdown.

Transition costs and loyalty cards

Transition costs can be erected by firms who do not wish to lose their customer base. Cell-phone plans are a good example. Contract-termination costs are one obstacle to moving to a new supplier. Some carriers grant special low rates to users communicating with other users within the same network, or offer special rates for a block of users (perhaps within a family). Tim Hortons and other coffee chains offer loyalty cards that give one free cup of coffee for every eight purchased. These suppliers are not furnishing love to their caffeine consumers, they are providing their consumers with an incentive not to switch to a competing supplier. Air miles rewards operate on a similar principle. So too do loyalty cards for hotel chains.

How do competitors respond to these loyalty programs? Usually by offering their own. Hilton and Marriot each compete by offering a free night after a given points threshold is reached.

Over-investment

An *over-investment* strategy means that an existing supplier generates additional production capacity through investment in new plant or capital. This is costly to the incumbent and is intended as a signal to any potential entrant that this capacity could be brought on-line immediately should a potential competitor contemplate entry. For example, a ski-resort owner may invest in a new chair-lift, even if she does not use it frequently. The existence of the additional capacity may scare potential entrants. A key component of this strategy is that the incumbent firm invests ahead of time – and inflicts a cost on itself. The incumbent does not simply say "I will build another chair-lift if you decide to develop a nearby mountain into a ski hill." That policy does not carry the same degree of credibility as actually incurring the cost of construction ahead of time. However, such a strategy may not always be feasible: It might be just too costly to pre-empt entry by putting spare capacity in place. Spare capacity is not so different from

brand development through advertising; both are types of sunk cost. The threats associated with the incumbent's behaviour become a credible threat because the incumbent incurs costs up front.

A **credible threat** is one that is effective in deterring specific behaviours; a competitor must believe that the threat will be implemented if the competitor behaves in a certain way.

Lobbying

In our chapter on monopoly we stressed the role of political/lobbying activity. Large firms invariably employ public relations firms, and maintain their own public relations departments. The role of these units is not simply to portray a positive image of the corporation to the public; it is to maintain and increase whatever market power such firms already possess. It is as much in the interest of an oligopolistic firm as a monopolist to prevent entry and preserve supernormal profits.

In analyzing perfect competition, we saw that free entry is critical to maintaining normal profits. Lobbying is designed to obstruct entry, and it is also designed to facilitate mergers and acquisitions. The economist Thomas Philippon has written about the increasing concentration of economic power in recent decades in the hands of a small number of corporations in many sectors of the North American economy. He argues that this concentration of power contributes to making the distribution of income more favorable to corporate interests and less favorable to workers. In his recent book ("The Great Reversal: How America Gave up Free Markets"), he shows that, contrary to traditional beliefs, Europe is now much more competitive than the US in most sectors of the economy. More broadband suppliers result in rates in Europe that are about half of US rates. Whereas in the US four airlines control 80% of the market, In Europe they control 40%. If scale economies were the prime determinant of corporate concentration we should not expect such large differences. Likewise, if globalization and technological change were the main determinants of corporate concentration, we should expect experiences in Europe and North America to be similar. But they are not. Hence, it is reasonable to conclude that entry barriers in North America are more effective, or that regulatory forces are stronger in Europe.

11.9 Matching markets: design

Markets are institutions that facilitate the exchange of goods and services. They act as clearing houses. The normal medium of exchange is money in some form. But many markets deal in exchanges that do not involve money and frequently involve matching: Graduating medical students are normally matched with hospitals in order that graduates complete their residency requirement; in many jurisdictions in the US applicants for places in public schools that form a pool within a given school-board must go through an application process that sorts the applicants into the different schools within the board; patients in need of a new kidney must be matched with kidney donors.

These markets are clearinghouses and have characteristics that distinguish them from traditional currency-based markets that we have considered to this point.

- The good or service being traded is generally *heterogeneous*. For example, patients in search of a kidney donor must be medically compatible with the eventual donor if the organ transplant is not to be rejected. Hospitals may seek residents in particular areas of health, and they must find residents who are, likewise, seeking such placements. Students applying to public schools may be facing a choice between schools that focus upon science or upon the arts. Variety is key.
- Frequently the idea of a market that is mediated by money is *repugnant*. For example, the only economy in the world that permits the sale of human organs is Iran. Elsewhere the idea of a monetary payment for a kidney is unacceptable. A market in which potential suppliers of kidneys registered their reservation prices and demanders registered their willingness to pay is incompatible with our social mores. Consequently, potential living donors or actual deceased donors must be directly matched with a patient in need. While some individuals believe that a market in kidneys would do more good than harm, because a monetary payment might incentivize the availability of many more organs and therefore save many more lives, virtually every society considers the downside to such a trading system to outweigh the benefits.
- Modern matching markets are more frequently *electronically mediated*, and the communications revolution has led to an increase in the efficiency of these markets.

The Economics prize in memory of Alfred Nobel was awarded to Alvin Roth and Lloyd Shapley in 2011 in recognition of their contributions to designing markets that function efficiently in the matching of demanders and suppliers of the goods and services. What do we mean by an efficient mechanism? One way is to define it is similar to how we described the market for apartments in Chapter 5: following an equilibrium in the market, is it possible to improve the wellbeing of one participant without reducing the wellbeing of another? We showed in that example that the market performed efficiently: a different set of renters getting the apartments would reduce total surplus in the system.

Consider a system in which medical graduates are matched with hospitals, and the decision process results in the potential for improvement: Christina obtains a residency in the local University Hospital while Ulrich obtains a residency at the Childrens' Hospital. But Christina would have preferred the Childrens' and Ulrich would have preferred the University. The matching algorithm here was not efficient because, at the end of the allocation process, there is scope for gains for each individual. Alvin Roth devised a matching mechanism that surmounts this type of inefficiency. He called it the deferred acceptance algorithm.

Roth also worked on the matching of kidney donors to individuals in need of a kidney. The fundamental challenge in this area is that a patient in need of a kidney may have a family member, say a sibling, who is willing to donate a kidney, but the siblings are not genetically compatible. The patient's immune system may attack the implantation of a 'foreign' organ. One solution to this incompatibility is to find matching pairs of donors that come from a wider choice set. Two families in each of which there is patient and a donor may be able to cross-donate: donor in Family A can donate to patient in Family B, and donor in Family B can donate to patient in Family A, in the sense that the donor organs will not be rejected by recipients' immune systems. Hence if many patient-donor families register in a clearinghouse, a computer algorithm can search for matching pairs. Surgical operations may be performed simultaneously in order to prevent one donor from backing out following his sibling's receipt of a kidney.

A more recent development concerns 'chains'. In this case a good Samaritan ('unaligned donor') offers a kidney while seeking nothing in return. The algorithm then seeks a match for the good Samaritan's kidney among all of the recipient-donor couples registered in the data bank. Having found (at least) one, the algorithm seeks a recipient for the kidney that will come from the first recipient's donor partner. And so on. It turns out that an algorithm which seeks to maximize the potential number of participating pairs is fraught with technical and ethical challenges: should a young patient, who could benefit from the organ for a whole lifetime, get priority over an older patient, who will benefit for fewer years of life, even if the older patient is in greater danger of dying in the absence of a transplant? This is an ethical problem.

Examples where these algorithms have achieved more than a dozen linked transplants are easy to find on an internet search - they are called *chains*, for the obvious reason.

Consider the following efficiency aspect of the exchange. Suppose a patient has two siblings, each of whom is willing to donate (though only one of the two actually will); should such a patient get priority in the computer algorithm over a patient who has just a single sibling willing to donate? The answer may be yes; the dual donor patient should get priority because if his two siblings have different blood types, this greater variety on the supply side increases the chances for matching in the system as a whole and is therefore beneficial. If a higher priority were not given to the dual-donor patient, there would be an incentive for him to name just one potential donor, and that would impact the efficiency of the whole matching algorithm.

It is not always recognized that the discipline of Economics explores social problems of the nature we have described here, despite the fact that the discipline has developed the analytical tools to address them.

Conclusion

Monopoly and perfect competition are interesting paradigms; but few markets resemble them in the real world. In this chapter we addressed some of the complexities that define the economy we inhabit: It is characterized by strategic planning, entry deterrence, differentiated products and so forth.

Entry and exit are critical to competitive markets. Frequently entry is blocked because of scale economies – an example of a natural or unintended entry barrier. In addition, incumbents can formulate numerous strategies to limit entry.

Firms act strategically – particularly when there are just a few participants in the market. Before acting, firms make conjectures about how their competitors will react, and incorporate such reactions into their own planning. Competition between suppliers can frequently be analyzed in terms of a game, and such games usually have an equilibrium outcome. The Cournot duopoly model that we developed is a game between two competitors in which an equilibrium market output is determined from a pair of reaction functions.

Scale economies are critical. Large development costs or setup costs may mean that the market can generally support just a limited number of producers. In turn this implies that potential new (small-scale) firms cannot benefit from the scale economies and will not survive competition from large-scale suppliers.

Product differentiation is critical. If small differences exist between products produced in markets where there is free entry we get a monopolistically competitive structure. In these markets long-run profits are 'normal' and firms operate with some excess capacity. It is not possible to act strategically in this kind of market.

The modern economy also has sectors that have successfully erected barriers. These barriers lead to fewer competitors than could efficiently supply the market. Ultimately the owners of capital are the beneficiaries of these barriers and consumers suffer from higher prices.

Key Terms

Imperfectly competitive firms face a downward-sloping demand curve, and their output price reflects the quantity sold.

Oligopoly defines an industry with a small number of suppliers.

Monopolistic competition defines a market with many sellers of products that have similar characteristics. Monopolistically competitive firms can exert only a small influence on the whole market.

Duopoly defines a market or sector with just two firms.

Concentration ratio: N -firm concentration ratio is the sales share of the largest N firms in that sector of the economy.

Differentiated product is one that differs slightly from other products in the same market.

The **monopolistically competitive equilibrium** in the long run requires the firm's demand curve to be tangent to the ATC curve at the output where $MR=MC$.

Collusion is an explicit or implicit agreement to avoid competition with a view to increasing profit.

Conjecture: a belief that one firm forms about the strategic reaction of another competing firm.

Game: a situation in which contestants plan strategically to maximize their profits, taking account of rivals' behaviour.

Strategy: a game plan describing how a player acts, or moves, in each possible situation.

Nash equilibrium: one in which each player chooses the best strategy, given the strategies chosen by the other player, and there is no incentive for any player to move.

Dominant strategy: a player's best strategy, whatever the strategies adopted by rivals.

Payoff matrix: defines the rewards to each player resulting from particular choices.

Credible threat: one that, after the fact, is still optimal to implement.

Cournot behaviour involves each firm reacting optimally in their choice of output to their competitors' decisions.

Reaction functions define the optimal choice of output conditional upon a rival's output choice.

Exercises for Chapter 11

EXERCISE 11.1

Imagine that the biggest four firms in each of the sectors listed below produce the amounts defined in each cell. Compute the three-firm and four-firm concentration ratios for each sector, and rank the sectors by degree of industry concentration.

Sector	Firm 1	Firm 2	Firm 3	Firm 4	Total market
Shoes	60	45	20	12	920
Chemicals	120	80	36	24	480
Beer	45	40	3	2	110
Tobacco	206	84	30	5	342

EXERCISE 11.2

You own a company in a monopolistically competitive market. Your marginal cost of production is \$12 per unit. There are no fixed costs. The demand for your own product is given by the equation $P=48-(1/2)Q$.

1. Plot the demand curve, the marginal revenue curve, and the marginal cost curve.
2. Compute the profit-maximizing output and price combination.
3. Compute total revenue and total profit [*Hint:* Remember $AC=MC$ here].

4. In this monopolistically competitive industry, can these profits continue indefinitely?

EXERCISE 11.3

Two firms in a particular industry face a market demand curve given by the equation $P=100-(1/3)Q$. The marginal cost is \$40 per unit and the marginal revenue is $MR=100-(2/3)Q$. The quantity intercepts for demand and MR are 300 and 150.

1. Draw the demand curve and MR curve to scale on a diagram. Then insert the MC curve.
2. If these firms got together to form a cartel, what output would they produce and what price would they charge?
3. Assuming they each produce half of the total what is their individual profit?

EXERCISE 11.4

The classic game theory problem is the "prisoners' dilemma." In this game, two criminals are apprehended, but the police have only got circumstantial evidence to prosecute them for a small crime, without having the evidence to prosecute them for the major crime of which they are suspected. The interrogators then pose incentives to the crooks-incentives to talk. The crooks are put in separate jail cells and have the option to confess or deny. Their payoff depends upon what course of action each adopts. The payoff matrix is given below. The first element in each box is the payoff (years in jail) to the player in the left column, and the second element is the payoff to the player in the top row.

		B's strategy	
		Confess	Deny
A's strategy	Confess	6,6	0,10
	Deny	10,0	1,1

1. Does a "dominant strategy" present itself for each or both of the crooks?
2. What is the Nash equilibrium to this game?
3. Is the Nash equilibrium unique?
4. Was it important for the police to place the crooks in separate cells?

EXERCISE 11.5

Taylor-made and Titleist are considering a production strategy for their new golf drivers. If they each produce a small output, they can price the product higher and make more profit than if they each produce a large output. Their payoff/profit matrix is given below.

		Taylor-made strategy	
		Low output	High output
Titleist strategy	Low output	50,50	20,70
	High output	70,20	40,40

1. Does either player have a dominant strategy here?
2. What is the Nash equilibrium to the game?
3. Do you think that a cartel arrangement would be sustainable?

EXERCISE 11.6

Ronnie's Wraps is the only supplier of sandwich food and makes a healthy profit. It currently charges a high price and makes a profit of six units. However, Flash Salads is considering entering the same market. The payoff matrix below defines the profit outcomes for different possibilities. The first entry in each cell is the payoff/profit to Flash Salads and the second to Ronnie's Wraps.

		Ronnie's Wraps

		High price	Low price
Flash Salads	Enter the market	2,3	-1,1
	Stay out of market	0,6	0,4

1. If Ronnie's Wraps threatens to lower its price in response to the entry of a new competitor, should Flash Salads stay away or enter?
2. Explain the importance of threat credibility here.

EXERCISE 11.7

Optional: Consider the market demand curve for appliances: $P=3,200-(1/4)Q$. There are no fixed production costs, and the marginal cost of each appliance is $MC = \$400$. As usual, the MR curve has a slope that is twice as great as the slope of the demand curve.

1. Illustrate this market geometrically.
2. Determine the output that will be produced in a 'perfectly competitive' market structure where no profits accrue in equilibrium.
3. If this market is supplied by a monopolist, illustrate the choice of output.

EXERCISE 11.8

Optional: Consider the outputs you have obtained in Exercise 11.7.

1. Can you figure out how many firms would produce at the perfectly competitive output? If not, can you think of a reason?
2. If, in contrast, each firm in that market had to cover some fixed costs, in addition to the variable costs defined by the MC value, would that put a limit on the number of firms that could produce in this market?
1. This presentation is inspired by a similar problem in Ted Bergstrom and Hal Varian's book "*Workouts*".
2. This game is sometimes called the "prisoners' dilemma" game because it can be constructed to reflect a scenario in which two prisoners, under isolated questioning, each confess to a crime and each ends up worse off than if neither confessed.
3. <https://slate.com/technology/2013/10/amazon-book-how-jeff-bezos-went-thermonuclear-on-diapers-com.html>

This page titled [11: Imperfect competition](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.1: The principle ideas

The preceding chapters have explored extreme forms of supply: The monopolist is the sole supplier and possesses as much market power as possible. In contrast, the perfect competitor is small and has no market power whatsoever. He simply accepts the price for his product that is determined in the market by the forces of supply and demand. These are very useful paradigms to explore, but the real world for the most part lies between these extremes. We observe that there are a handful of dominant brewers in Canada who supply more than three quarters of the market, and they are accompanied by numerous micro brewers that form the fringe of the brewing business. We have a small number of air carriers and one of them controls half of the national market. The communications market has just three major suppliers; the Canadian Football League has nine teams and there are just a handful of major hardware/builders' suppliers stores nationally. At the other end of the spectrum we have countless restaurants and fitness centres, but they do not supply exactly the same product to the marketplaces for 'food' or 'health', and so these markets are not perfectly competitive, despite the enormous number of participants.

In this chapter we will explore three broad topics: First is the relationship between firm behaviour and firm size relative to the whole sector. This comes broadly under the heading of *imperfect competition* and covers a variety of market forms. Second, we will explore the principle modern ideas in *strategic behavior*. In a sense all decisions in microeconomics have an element of strategy to them - economic agents aim to attain certain goals and they adopt specific maximizing strategies to attain them. But in this chapter we explore a more specific concept of strategic behavior - one that focuses upon direct interactions between a small number of players in the market place. Third, we explore the principle characteristics of what are termed *matching' markets*. These are markets where transactions take place without money and involve matching heterogeneous suppliers with heterogeneous buyers.

This page titled [11.1: The principle ideas](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.2: Imperfect competitors

Imperfect competitors can be defined by the number of firms in their sector, or the share of total sales going to a small number of suppliers. They can also be defined in terms of the characteristics of the demand curves they all face. A perfect competitor faces a perfectly elastic demand at the existing market price, and this is the only market structure to have this characteristic. In all other market structures suppliers effectively face a downward-sloping demand. This means that they have some influence on the price of the good, and also that if they change the price they charge, they can expect demand to reflect this in a predictable manner. So, in theory, we can classify all market structures apart from perfect competition as being imperfectly competitive. In practice we use the term to denote firms that fall between the extremes of perfect competition and monopoly.

Imperfectly competitive firms face a downward-sloping demand curve, and their output price reflects the quantity sold.

The demand curve for the firm and industry coincide for the monopolist, but not for other imperfectly competitive firms. It is convenient to categorize the producing sectors of the economy as either having a relatively small number of participants, or having a large number. The former market structures are called oligopolistic, and the latter are called monopolistically competitive. The word *oligopoly* comes from the Greek word *oligos* meaning few, and *polein* meaning to sell.

Oligopoly defines a market with a small number of suppliers.

Monopolistic competition defines a market with many sellers of products that have similar characteristics. Monopolistically competitive firms can exert only a small influence on the whole market.

The home appliance industry is an oligopoly. The prices of *KitchenAid* appliances depend not only on their own output and sales, but also on the prices of *Whirlpool*, *Maytag* and *Bosch*. If a firm has just two main producers it is called a duopoly. *Canadian National* and *Canadian Pacific* are the only two major rail freight carriers in Canada; they thus form a duopoly. In contrast, the local Italian restaurant is a monopolistic competitor. Its output is a package of distinctive menu choices, personal service, and convenience for local customers. It can charge a different price than the out-of-neighbourhood restaurant, but if its prices are too high local diners may travel elsewhere for their food experience, or switch to a different cuisine locally. Many markets are defined by producers who supply similar but not identical products. Canada's universities all provide degrees, but they differ one from another in their programs, their balance of in-class and on-line courses, their student activities, whether they are science based or liberal arts based, whether they have cooperative programs or not, and so forth. While universities are not in the business of making profit, they certainly wish to attract students, and one way of doing this is to differentiate themselves from other institutions. The profit-oriented world of commerce likewise seeks to increase its market share by distinguishing its product line.

Duopoly defines a market or sector with just two firms.

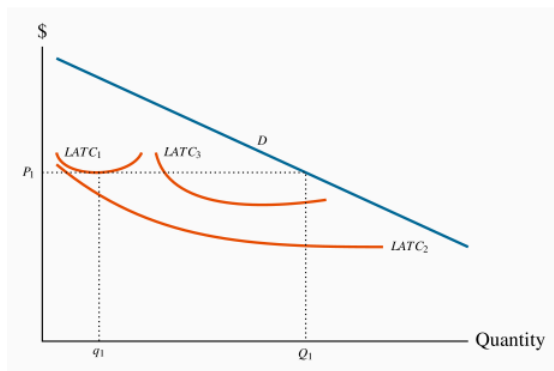
These distinctions are not completely airtight. For example, if a sole domestic producer is subject to international competition it cannot act in the way we described in the previous chapter – it has potential, or actual, competition. *Bombardier* may be Canada's sole rail car manufacturer, but it is not a monopolist, even in Canada. It could best be described as being part of an international oligopoly in rail-car manufacture. Likewise, it is frequently difficult to delineate the boundary of a given market. For example, is *Canada Post* a monopoly in mail delivery, or an oligopolist in hard-copy communication? We can never fully remove these ambiguities.

The role of cost structures

A critical determinant of market structure is the way in which demand and cost interact to determine the likely number of market participants in a given sector or market. Structure also evolves over the long run: Time is required for entry and exit.

Figure 11.1 shows the demand curve D for the output of an industry in the long run. Suppose, initially, that all firms and potential entrants face the long-run average cost curve $LATC_1$. At the price P_1 , free entry and exit means that each firm produces q_1 . With the demand curve D , industry output is Q_1 . The number of firms in the industry is $N_1 (=Q_1/q_1)$. If q_1 , the minimum average cost output on $LATC_1$, is small relative to D , then N_1 is large. This outcome might be perfect competition (N virtually infinite), or monopolistic competition (N large) with slightly differentiated products produced by each firm.

Figure 11.1 Demand, costs and market structure



With a cost structure defined by $LATC_1$ this market has space for many firms – perfect or monopolistic competition, each producing approximately q_1 . If costs correspond to $LATC_2$, where scale economies are substantial, there may be space for just one producer. The intermediate case, $LATC_3$, can give rise to oligopoly, with each firm producing more than q_1 but less than a monopolist. These curves encounter their MES at very different output levels.

Instead, suppose that the production structure in the industry is such that the long-run average cost curve is $LATC_2$. Here, scale economies are vast, relative to the market size. At the lowest point on this cost curve, output is large relative to the demand curve D . If this one firm were to act like a monopolist it would produce an output where $MR=MC$ in the long run and set a price such that the chosen output is sold. Given the scale economies, there may be no scope for another firm to enter this market, because such a firm would have to produce a very high output to compete with the existing producer. This situation is what we previously called a "natural" monopolist.

Finally, the cost structure might involve curves of the type $LATC_3$, which would give rise to the possibility of several producers, rather than one or very many. This results in oligopoly.

It is clear that one *crucial determinant of market structure is minimum efficient scale relative to the size of the total market* as shown by the demand curve. The larger the minimum efficient scale relative to market size, the smaller is the number of producers in the industry.

This page titled [11.2: Imperfect competitors](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.3: Imperfect competitors- measures of structure and market power

Sectors of the economy do not fit neatly into the limited number of categories described above. The best we can say in most cases is that they resemble more closely one type of market than another. Consider the example of Canada's brewing sector: It has two large brewers in *Molson-Coors* and *Labatt*, a couple of intermediate sized firms such as *Sleeman*, and an uncountable number of small boutique brew pubs. While such a large number of brewers satisfy one requirement for perfect competition, it would not be true to say that the biggest brewers wield no market power; and this is the most critical element in defining market structure.

By the same token, we could not define this market as a duopoly: Even though there are just two major participants, there are countless others who, together, are important.

One way of defining what a particular structure most closely resembles is to examine the percentage of sales in the market that is attributable to a small number of firms. For example: What share is attributable to the largest three or four firms? The larger the share, the more concentrated the market power. Such a statistic is called a concentration ratio. The *N*-firm concentration ratio is the sales share of the largest *N* firms in that sector of the economy.

The ***N*-firm concentration ratio** is the sales share of the largest *N* firms in that sector of the economy.

Table 11.1 Concentration in Canadian food processing 2011

Sector	% of shipments
Sugar	98
Breakfast cereal	96
Canning	60
Meat processing	23

Source: "Four Firm Concentration Ratios (CR4s) for selected food processing sectors," adapted from Statistics Canada publication Measuring industry concentration in Canada's food processing sectors, Agriculture and Rural Working Paper series no. 70, Catalogue 21-601, <http://www.statcan.gc.ca/pub/21-601-m/21-601-m2004070-eng.pdf>.

Table 11.1 contains information on the 4-firm concentration ratio for several sectors of the Canadian economy. It indicates that, at one extreme, sectors such as breakfast cereals and sugars have a high degree of concentration, whereas meat processing has much less. A high degree of concentration suggests market power, and possibly economies of scale.

This page titled [11.3: Imperfect competitors- measures of structure and market power](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.4: Imperfect competition- monopolistic competition

Monopolistic competition presumes a large number of quite small producers or suppliers, each of whom may have a slightly differentiated product. The competition element of this name signifies that there are many participants, while the monopoly component signifies that each supplier faces a downward-sloping demand. In concrete terms, your local coffee shop that serves "fair trade" coffee has a product that differs slightly from that of neighbouring shops that sell the traditional product. They coexist in the same sector, and probably charge different prices: The fair trade supplier likely charges a higher price, but knows nonetheless that too large a difference between her price and the prices of her competitors will see some of her clientele migrate to those lower-priced establishments. That is to say, she faces a downward-sloping demand curve.

The competition part of the name also indicates that there is *free entry and exit*. There are no barriers to entry. As a consequence, we know at the outset that only normal profits will exist in a long-run equilibrium. Economic profits will be competed away by entry, just as losses will erode due to exit.

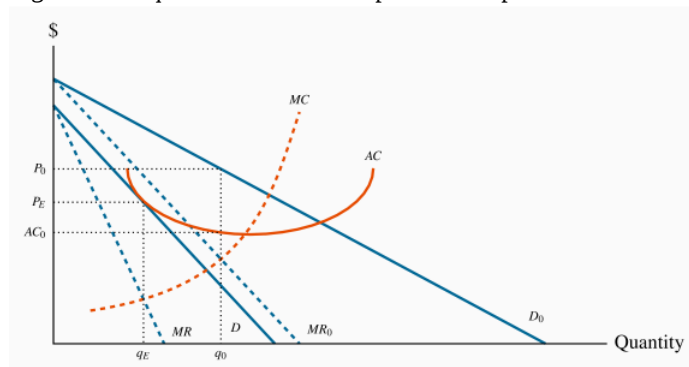
As a general rule then, each firm can influence its market share to some extent by changing its price. Its demand curve is not horizontal because different firms' products are only limited substitutes. A lower price level may draw some new customers away from competitors, but convenience or taste will prevent most patrons from deserting their local businesses. In concrete terms: A pasta special at the local Italian restaurant that reduces the price below the corresponding price at the competing local Thai restaurant will indeed draw clients away from the latter, but the foods are sufficiently different that only some customers will leave the Thai restaurant. The differentiated menus mean that many customers will continue to pay the higher price.

A **differentiated product** is one that differs slightly from other products in the same market.

Given that there are very many firms, the theory also envisages limits to scale economies. Firms are small and, with many competitors, individual firms do not compete strategically with *particular* rivals. Because the various products offered are slightly differentiated, we avoid graphics with a *market* demand, because this would imply that a uniform product is being considered. At the same time the market is a well-defined concept—it might be composed of all those restaurants within a reasonable distance, for example, even though each one is slightly different from the others. The market share of each firm depends on the price that it charges *and* on the number of competing firms. For a given number of suppliers, a shift in industry demand also shifts the demand facing each firm. Likewise, the presence of more firms in the industry reduces the demand facing each one.

Equilibrium is illustrated in Figure 11.2. Here D_0 is the initial demand facing a representative firm, and MR_0 is the corresponding marginal revenue curve. Profit is maximized where $MC=MR$, and the price P_0 is obtained from the demand curve corresponding to the output q_0 . Total profit is the product of output times the difference between price and average cost, which equals $q_0 \times (P_0 - AC_0)$.

Figure 11.2 Equilibrium for a monopolistic competitor



Profits exist at the initial equilibrium (q_0, P_0) . Hence, new firms enter and reduce the share of the total market faced by each firm, thereby shifting back their demand curve. A final equilibrium is reached where economic profits are eliminated: At $AC=P_E$ and $MR=MC$.

With free entry, such profits attract new firms. The increased number of firms reduces the share of the market that any one firm can claim. That is, the firm's demand curve shifts inwards when entry occurs. As long as (economic) profits exist, this process continues. For entry to cease, average cost must equal price. A final equilibrium is illustrated by the combination (P_E, q_E) , where the demand has shifted inward to D .

At this long-run equilibrium, two conditions must hold: First, the optimal pricing rule must be satisfied—that is $MC=MR$; second it must be the case that only normal profits are made at the final equilibrium. Economic profits are competed away as a result of free entry. Graphically this implies that ATC must equal price at the output where $MC=MR$. In turn this implies that the ATC is tangent to the demand curve where $P=ATC$. While this could be proven mathematically, it is easy to intuit why this tangency must exist: If ATC merely intersected the demand curve at the output where $MC=MR$, we could find some other output where the demand price would be above ATC , suggesting that profits could be made at such an output. Clearly that could not represent an equilibrium.

The **monopolistically competitive equilibrium** in the long run requires the firm's demand curve to be tangent to the ATC curve at the output where $MR=MC$.

This page titled [11.4: Imperfect competition- monopolistic competition](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.5: Imperfect competition- economies of scope and platforms

The communications revolution has impacted market structure in modern economies profoundly: it has facilitated economies of scope, meaning that firms may yield more collective profit if merged than if operating independently.

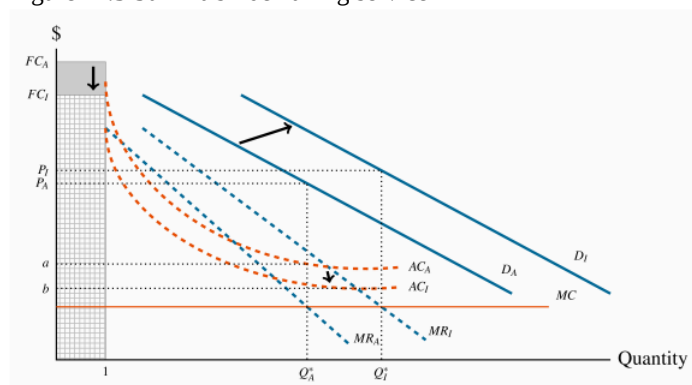
Economies of Scope

Imagine an aspiring entrepreneur who envisages a revolution of the traditional taxi sector of the economy. He decides to develop a smartphone application that will match independent income-seeking vehicle owners (drivers) with individuals seeking transport (passengers) from point A to point B. We know how this adventure evolves. In one case it takes the form of the corporation Uber, in another the corporation Lyft, and others worldwide.

These corporations have grown in leaps and bounds and have taken business from the conventional taxi corporations. As of 2019 they cannot turn a profit, yet the stock market continues to bet upon future success: investors believe that when these corporations evolve into fully integrated multi-product suppliers, both costs will decline and demand will increase for each component of the business. In the case of transportation companies, they aim to become a 'one-stop-shop' for mobility services. Uber is not only a ride-hailing service, it also transports meals through its Uber-eats platform, and is developing the electric scooter and electric bike markets in addition. In some local markets it is linked to public transport services. All of this is being achieved through a single smartphone application. The objective is to simplify movement for persons, by providing multiple options on a variety of transport modes, accessed through a single portal.

This phenomenon is described in Figure 11.3. The subscripts *A* and *I* represent market conditions when the service supplier is operating Alone or in an Integrated corporation. The initial equilibrium is defined by the *A* demand and cost conditions. The profit maximizing output occurs when $MC_A = MR_A$, leading to a price P_A and a quantity Q_A^* . Each unit of the good yields a profit margin of $(P_A - a)$.

Figure 11.3 Summa's ride hailing service



Demand for a particular product increases when the autonomous supplier (A) merges with another firm to become an integrated firm (I), because customers switch to firms that offer several different services from the same platform: the demand curve shifts outward, from D_A to D_I . With integration, the fixed costs fall and average costs fall, even with marginal costs constant. Output and profit increase, and concentration in the marketplace rises.

This firm now merges with another transportation corporation - perhaps a food delivery service, perhaps an electric bike service. Since each firm has a similar type of fixed cost, these costs can be reduced by the merger. In technical terms, the merged firms, or merged operations, share a common hardware-cum-software platform. Each firm will therefore incur lower average costs, even if marginal costs remain unchanged: the AC curve declines to AC_I . In addition to the decline in average costs, each firm sees an increase in its customer base, because transportation service buyers find it preferable to choose their mode of transport through a single portal rather than through several different modes of access. This is represented by an outward shift in the demand curve for vehicle rides to D_I .

The new profit maximizing equilibrium occurs at Q_I^* . Total profit necessarily increases both because average costs have fallen and the number of buyers willing to buy at any price has risen. The analytics in this figure also describe the benefits accruing to the other firm or firms in the merger.

A **platform** describes a technology that is common to more than one product in a multi-product organization.

We conclude from this analysis that, if scope economies are substantial, it may be difficult for stand-alone firms specializing in just one component of the transportation services sector to remain profitable. It may also be impossible to define a conventional equilibrium in this kind of marketplace. This is because some conglomerate firms may have different component producers in their suite of firms. For example, Lyft may not have a food delivery service, but it may have a limousine or bus service. What is critical for an equilibrium is that firms of a particular type, whether they are part of a conglomerate or not, be able to compete with corresponding firms. This means that their cost structure must be similar.

As a further example: Amazon initially was primarily an on-line book seller. But it expanded to include the sale of other products. And once it became a 'market for everything' the demand side of the market exploded in parallel with the product line, because it becomes easy to shop for 'anything' or even different objects on a single site. Only Walmart, in North America, comes close to being able to compete with Amazon.

This page titled [11.5: Imperfect competition- economies of scope and platforms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.6: Strategic behaviour- Oligopoly and games

Under perfect competition or monopolistic competition, there are so many firms in the industry that each one can ignore the immediate effect of its own actions on particular rivals. However, in an oligopolistic industry *each firm must consider how its actions affect the decisions of its relatively few competitors*. Each firm must guess how its rivals will react. Before discussing what constitutes an intelligent guess, we investigate whether they are likely to collude or compete. Collusion is a means of reducing competition with a view to increasing profit.

Collusion is an explicit or implicit agreement to avoid competition with a view to increasing profit.

A particular form of collusion occurs when firms co-operate to form a cartel, as we saw in the last chapter. Collusion is more difficult if there are many firms in the industry, if the product is not standardized, or if demand and cost conditions are changing rapidly. In the absence of collusion, each firm's demand curve depends upon how competitors react: If Air Canada contemplates offering customers a seat sale on a particular route, how will West Jet react? Will it, too, make the same offer to buyers? If Air Canada thinks about West Jet's likely reaction, will it go ahead with the contemplated promotion? A conjecture is a belief that one firm forms about the strategic reaction of another competing firm.

A **conjecture** is a belief that one firm forms about the strategic reaction of another competing firm.

Good poker players will attempt to anticipate their opponents' moves or reactions. Oligopolists are like poker players, in that they try to anticipate their rivals' moves. To study interdependent decision making, we use game theory. A game is a situation in which contestants plan strategically to maximize their payoffs, taking account of rivals' behaviour.

A **game** is a situation in which contestants plan strategically to maximize their payoffs, taking account of rivals' behaviour.

The *players* in the game try to maximize their own *payoffs*. In an oligopoly, the firms are the players and their payoffs are their profits. Each player must choose a strategy, which is a plan describing how a player moves or acts in different situations.

A **strategy** is a game plan describing how a player acts, or moves, in each possible situation.

Equilibrium outcomes

How do we arrive at an equilibrium in these games? Let us begin by defining a commonly used concept of equilibrium. A Nash equilibrium is one in which each player chooses the best strategy, given the strategies chosen by the other players, and there is no incentive to move or change choice.

A **Nash equilibrium** is one in which each player chooses the best strategy, given the strategies chosen by the other player, and there is no incentive for any player to move.

In such an equilibrium, no player wants to change strategy, since the other players' strategies were already figured into determining each player's own best strategy. This concept and theory are attributable to the Princeton mathematician John Nash, who was popularized by the Hollywood movie version of his life, *A Beautiful Mind*.

In most games, each player's best strategy depends on the strategies chosen by their opponents. Occasionally, a player's best strategy is independent of those chosen by rivals. Such a strategy is called a dominant strategy.

A **dominant strategy** is a player's best strategy, independent of the strategies adopted by rivals.

We now illustrate these concepts with the help of two different games. These games differ in their outcomes and strategies. Table 11.2 contains the domestic happiness game¹. Will and Kate are attempting to live in harmony, and their happiness depends upon each of them carrying out domestic chores such as shopping, cleaning and cooking. The first element in each pair defines Will's outcome, the second Kate's outcome. If both contribute to domestic life they each receive a happiness or utility level of 5 units. If one contributes and the other does not the happiness levels are 2 for the contributor and 6 for the non-contributor, or 'free-rider'. If neither contributes happiness levels are 3 each. When each follows the same strategy the payoffs are on the diagonal, when they follow different strategies the payoffs are on the off-diagonal. Since the elements of the table define the payoffs resulting from various choices, this type of matrix is called a payoff matrix.

A **payoff matrix** defines the rewards to each player resulting from particular choices.

So how is the game likely to unfold? In response to Will's choice of a contribute strategy, Kate's utility maximizing choice involves lazing: She gets 6 units by not contributing as opposed to 5 by contributing. Instead, if Will decides to be lazy what is in Kate's best interest? Clearly it is to be lazy also because that strategy yields 3 units of happiness compared to 2 units if she contributes. In sum,

Kate's best strategy is to be lazy, regardless of Will's behaviour. So the strategy of not contributing is a *dominant strategy*, in this particular game.

Will also has a dominant strategy – identical to Kate's. This is not surprising since the payoffs are symmetric in the table. Hence, since each has a dominant strategy of not contributing the Nash equilibrium is in the bottom right cell, where each receives a payoff of 3 units. Interestingly, this equilibrium is not the one that yields maximum combined happiness.

Table 11.2 A game with dominant strategies

		Kate's choice	
		Contribute	Laze
Will's choice	Contribute	5,5	2,6
	Laze	6,2	3,3

The first element in each cell denotes the payoff or utility to Will; the second element the utility to Kate.

The reason that the equilibrium yields less utility for each player in this game is that the game is competitive: Each player tends to their own interest and seeks the best outcome conditional on the choice of the other player. This is evident from the (5,5) combination. From this position Kate would do better to defect to the Laze strategy, because her utility would increase².

To summarize: This game has a unique equilibrium and each player has a dominant strategy. But let us change the payoffs just slightly to the values in Table 11.3. The off-diagonal elements have changed. The contributor now gets no utility as a result of his or her contributions: Even though the household is a better place, he or she may be so annoyed with the other person that no utility flows to the contributor.

Table 11.3 A game without dominant strategies

		Kate's choice	
		Contribute	Laze
Will's choice	Contribute	5,5	0,4
	Laze	4,0	3,3

The first element in each cell denotes the payoff or utility to Will; the second element the utility to Kate.

What are the optimal choices here? Starting again from Will choosing to contribute, what is Kate's best strategy? It is to contribute: She gets 5 units from contributing and 4 from lazing, hence she is better contributing. But what is her best strategy if Will decides to laze? It is to laze, because that yields her 3 units as opposed to 0 by contributing. This set of payoffs therefore *contains no dominant strategy* for either player.

As a result of there being no dominant strategy, there arises the possibility of more than one equilibrium outcome. In fact there are two equilibria in this game now: If the players find themselves both contributing and obtaining a utility level of (5,5) it would not be sensible for either one to defect to a laze option. For example, if Kate decided to laze she would obtain a payoff of 4 utils rather than the 5 she enjoys at the (5,5) equilibrium. By the same reasoning, if they find themselves at the (laze, laze) combination there is no incentive to move to a contribute strategy.

Once again, it is to be emphasized that the twin equilibria emerge in a competitive environment. If this game involved cooperation or collusion the players should be able to reach the (5,5) equilibrium rather than the (3,3) equilibrium. But in the competitive environment we cannot say *ex ante* which equilibrium will be attained.

Repeated games

This game illustrates the tension between collusion and competition. While we have developed the game in the context of the household, it can equally be interpreted in the context of a profit maximizing game between two market competitors. Suppose the numbers define profit levels rather than utility as in Table 11.4. The 'contribute' option can be interpreted as 'cooperate' or 'collude', as we described for a cartel in the previous chapter. They collude by agreeing to restrict output, sell that restricted output at a higher price, and in turn make a greater total profit which they split between themselves. The combined best profit outcome (5,5) arises when each firm restricts its output.

Table 11.4 Collusion possibilities

		Firm K's profit	
		Low output	High output
Firm W's profit	Low output	5,5	2,6
	High output	6,2	3,3

The first element in each cell denotes the profit to Firm W; the second element the profit to Firm K.

But again there arises an incentive to defect: If Firm W agrees to maintain a high price and restrict output, then Firm K has an incentive to renege and increase output, hoping to improve its profit through the willingness of Firm W to restrict output. Since the game is symmetric, each firm has an incentive to renege. Each firm has a dominant strategy – high output, and there is a unique equilibrium (3,3).

Obviously there arises the question of whether these firms can find an operating mechanism that would ensure they each generate a profit of 5 units rather than 3 units, while remaining purely self-interested. This question brings us to the realm of repeated games. For example, suppose that firms make strategic choices each quarter of the year. If firm K had 'cheated' on the collusive strategy it had agreed with firm W in the previous quarter, what would happen in the following quarter? Would firms devise a strategy so that cheating would not be in the interest of either one, or would the competitive game just disintegrate into an unpredictable pattern? These are interesting questions and have provoked a great deal of thought among game theorists. But they are beyond our scope at the present time.

A repeated game is one that is repeated in successive time periods and where the knowledge that the game will be repeated influences the choices and outcomes in earlier periods.

We now examine what might happen in one-shot games of the type we have been examining, but in the context of many possible choices. In particular, instead of assuming that each firm can choose a high or low output, how would the outcome of the game be determined if each firm can choose an output that can lie anywhere between a high and low output? In terms of the demand curve for the market, this means that the firms can choose some output and price that is consistent with demand conditions: There may be an infinite number of choices. This framing of a game enables us to explore new concepts in strategic behavior.

This page titled [11.6: Strategic behaviour- Oligopoly and games](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.7: Strategic behaviour- Duopoly and Cournot games

The duopoly model that we frequently use in economics to analyze competition between a small number of competitors is fashioned after the ideas of French economist Augustin Cournot. Consequently it has come to be known as the Cournot duopoly model. While the maximizing behaviour that is incorporated in this model can apply to a situation with several firms rather than two, we will develop the model with two firms. This differs slightly from the preceding section, where each firm has simply a choice between a high or low output.

The critical element of the Cournot approach is that the firms each determine their optimal strategy – one that maximizes profit – by reacting optimally to their opponent's strategy, which in this case involves their choice of output.

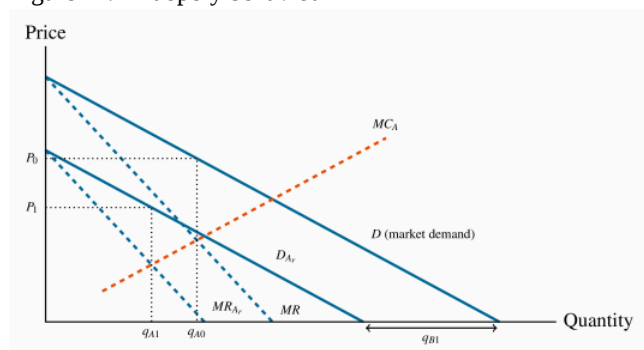
Cournot behaviour involves each firm reacting optimally in their choice of output to their competitors' output decisions.

A central element here is the reaction function of each firm, which defines the optimal output choice conditional upon their opponent's choice.

Reaction functions define the optimal choice of output conditional upon a rival's output choice.

We can develop an optimal strategy with the help of Figure 11.4. D is the market demand, and two firms supply this market. If B supplies a zero output, then A would face the whole demand, and would maximize profit where $MC=MR$. Let this output be defined by q_{A0} . We transfer this output combination to Figure 11.5, where the output of each firm is on one of the axes— A on the vertical axis and B on the horizontal. This particular combination of zero output for B and q_{A0} for A is represented on the vertical axis as the point q_{A0} .

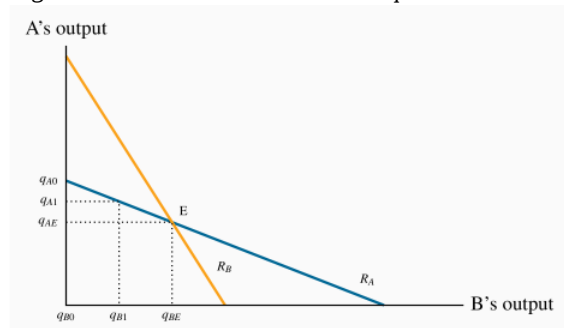
Figure 11.4 Duopoly behaviour



When one firm, B , chooses a specific output, e.g. q_{B1} , then A 's residual demand D_{Ar} is the difference between the market demand and q_{B1} . A 's profit is maximized at q_{A1} – where $MC = MR_{Ar}$. This is an optimal reaction by A to B 's choice. For all possible choices by B , A can form a similar optimal response. The combination of these responses forms A 's reaction function.

Instead, suppose that B produces a quantity q_{B1} in Figure 11.4. This reduces the demand curve facing A correspondingly from D to D_{Ar} , which we call A 's *residual demand*. When subject to such a choice by B , firm A maximizes profit by producing where $MR_{Ar} = MC$, where MR_{Ar} is the marginal revenue corresponding to the residual demand D_{Ar} . The optimum for A is now q_{A1} , and this pair of outputs is represented by the combination (q_{A1}, q_{B1}) in Figure 11.5.

Figure 11.5 Reaction functions and equilibrium



The reaction function for A (R_A) defines the optimal output response for A to any output choice by B . The reaction function for B is defined similarly. The equilibrium occurs at the intersection of R_A and R_B . Any other combination will induce one firm to change

its output, and therefore could not be an equilibrium.

Firm A forms a similar optimal response for every possible output level that B could choose, and these responses define A's reaction function. The reaction function illustrated for A in Figure 11.5 is thus the locus of all optimal response outputs on the part of A. The downward-sloping function makes sense: The more B produces, the smaller is the residual market for A, and therefore the less A will produce.

But A is just one of the players in the game. If B acts in the same optimizing fashion, B too can formulate a series of optimal reactions to A's output choices. The combination of such choices would yield a reaction function for B . This is plotted as R_B in Figure 11.5.

An equilibrium is defined by the intersection of the two reaction functions, in this case by the point E. At this output level *each firm is making an optimal decision, conditional upon the choice of its opponent*. Consequently, neither firm has an incentive to change its output; therefore it can be called the Nash equilibrium.

Any other combination of outputs on either reaction function would lead one of the players to change its output choice, and therefore could not constitute an equilibrium. To see this, suppose that B produces an output greater than q_{BE} ; how will A react? A's reaction function indicates that it should choose a quantity to supply less than q_{AE} . If so, how will B respond in turn to that optimal choice? It responds with a quantity read from its reaction function, and this will be less than the amount chosen at the previous stage. By tracing out such a sequence of reactions it is clear that the output of each firm will move to the equilibrium q_E .

Application Box 11.1 Cournot: Fixed costs and brand

Why do we observe so many industries on the national, and even international, stages with only a handful of firms? For example, Intel produces more than half of the world's computer chips, and AMD produces a significant part of the remainder. Why are there only two major commercial aircraft producers in world aviation – Boeing and Airbus? Why are there only a handful of major North American suppliers in pharmaceuticals, automobile tires, soda pop, internet search engines and wireless telecommunications?

The answer lies primarily in the nature of modern product development. Product development (fixed) costs, coupled with a relatively small marginal cost of production, leads to markets where there is enough space for only a few players. The development cost for a new cell phone, or a new aircraft, or a new computer-operating system may run into billions, while the cost of producing each unit may in fact be constant. The enormous development cost associated with many products explains not only why there may be a small number of firms in the domestic market for the product, but also why the number of firms in some sectors is small worldwide.

The Cournot model yields an outcome that lies between monopoly (or collusion/cartel) and competitive market models. It does not necessarily assume that the firms are identical in terms of their cost structure, although the lower-cost producer will end up with a larger share of the market.

The next question that arises is whether this duopoly market will be sustained as a duopoly, or if entry may take place. In particular, if economic profits accrue to the participants will such profits be competed away by the arrival of new producers, or might there be barriers of either a 'natural' or 'constructed' type that operate against new entrants?

This page titled [11.7: Strategic behaviour- Duopoly and Cournot games](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.8: Strategic behaviour- Entry, exit and potential competition

At this point we inquire about the potential entry and impact of new firms – firms who might enter the industry if conditions were sufficiently enticing, meaning the presence of economic profits. One way of examining entry in this oligopolistic world is to envisage potential entry barriers as being either intended or unintended, though the difference between the two can be blurred. Broadly, an unintended or 'natural' barrier is one related to scale economies and the size of the market. An intended barrier involves a strategic decision on the part of the firm to prevent entry.

Unintended entry barriers

Oligopolists tend to have substantial fixed costs, accompanied by declining average costs up to high output levels. Such a cost structure 'naturally' gives rise to a supply side with a small number of suppliers. For examples, given demand and cost structures, could Vancouver support two professional soccer teams; could Calgary support two professional hockey teams; could Montreal sustain two professional football teams? The answer to each of these questions is likely 'no'. Because given the cost structure of these markets, it would not be possible to induce twice as many spectators without reducing the price per game ticket to such a degree that revenue would be insufficient to cover costs. (We will neglect for the moment that the governing bodies of these sports also have the power to limit entry.) Fixed costs include stadium costs, staff payrolls and player payrolls. In fact most costs in these markets are relatively fixed. Market size relative to fixed and variable costs is not large enough to sustain two teams in most cities. Exceptions in reality are huge urban areas such as New York and Los Angeles.

Accordingly, it is possible that the existing team, or teams, may earn economic profit from their present operation; but such profit does not entice further entry, because the market structure is such that the entry of an additional team could lead to each team making losses.

Intended entry barriers

Patent Law

This is one form of protection for incumbent firms. Research and development is required for the development of many products in the modern era. Pharmaceuticals are an example. If innovations were not protected, firms and individuals would not be incentivized to devote their energies and resources to developing new drugs. Society would be poorer as a result. Patent protection is obviously a legal form of protection. At the same time, patent protection can be excessive. If patents provide immunity from replication or copying for an excessive period of time - for longer than required to recoup R & D costs - then social welfare declines because monopoly profits are being generated as a result of output restriction at too high a price.

Advertising

Advertising is a second form of entry deterrence. In this instance firms attempt to market their product as being distinctive and even enviable. For example, *Coca-Cola* and *PepsiCo* invest hundreds of millions annually to project their products in this light. They sponsor sports, artistic and cultural events. Entry into the cola business is not impossible, but brand image is so strong for these firms that potential competitors would have a very low probability of entering this sector profitably. Likewise, in the 'energy-drinks' market, *Red Bull* spends hundreds of millions of dollars per annum on Formula One racing, kite surfing contests, mountain biking events and other extreme sports. In doing this it is reinforcing its brand image and distinguishing its product from Pepsi or Coca-Cola. This form of advertising is one of product differentiation and enables the manufacturer to maintain a higher price for its products by convincing its buyers that there are no close substitutes.

Predatory pricing

This form of pricing constitutes an illegal form of entry deterrence. It involves an incumbent charging an artificially low price for its product in the event of entry of a new competitor. This is done with a view to making it impossible for the entrant to earn a profit. Given that incumbents have generally greater resources than entrants, they can survive a battle of losses for a more prolonged period, thus ultimately driving out the entrant.

An iconic example of predatory pricing is that of Amazon deciding to take on a startup called Quidsi that operated the website *diapers.com*.³ The latter was proving to be a big hit with consumers in 2009 and Amazon decided that it was eating into Amazon profits on household and baby products. Amazon reacted by cutting its own prices dramatically, to the point where it was ready to

loose a huge amount of money in order to grind Quidsi into the ground. The ultimate outcome was that Quidsi capitulated and sold to Amazon.

Whether this was a legal tactic or not we do not know, but it underlines the importance of war chests.

Maintaining a war chest

Many large corporations maintain a mountain of cash. This might seem like an odd thing to do when it could be paying that cash out to owners in the form of dividends. But there are at least two reasons for not doing this. First, personal taxes on dividends are frequently higher than taxes on capital gains; accordingly if a corporation can transform its cash into capital gain by making judicious investments, that strategy ultimately yields a higher post-tax return to the stock holders. A second reason is that a cash war chest serves as a credible threat to competitors of the type described involving Amazon and Quidsi above.

Network externalities

These externalities arise when the existing number of buyers itself influences the total demand for a product. *Facebook* is now a classic example. An individual contemplating joining a social network has an incentive to join one where she has many existing 'friends'. Not everyone views the *Microsoft* operating system (OS) as the best. Many prefer a simpler system such as *Linux* that also happens to be free. However, the fact that almost every new computer (that is not *Apple*) coming onto the market place uses Microsoft OS, there is an incentive for users to continue to use it because it is so easy to find a technician to repair a breakdown.

Transition costs and loyalty cards

Transition costs can be erected by firms who do not wish to lose their customer base. Cell-phone plans are a good example. Contract-termination costs are one obstacle to moving to a new supplier. Some carriers grant special low rates to users communicating with other users within the same network, or offer special rates for a block of users (perhaps within a family). Tim Hortons and other coffee chains offer loyalty cards that give one free cup of coffee for every eight purchased. These suppliers are not furnishing love to their caffeine consumers, they are providing their consumers with an incentive not to switch to a competing supplier. Air miles rewards operate on a similar principle. So too do loyalty cards for hotel chains.

How do competitors respond to these loyalty programs? Usually by offering their own. Hilton and Marriot each compete by offering a free night after a given points threshold is reached.

Over-investment

An *over-investment* strategy means that an existing supplier generates additional production capacity through investment in new plant or capital. This is costly to the incumbent and is intended as a signal to any potential entrant that this capacity could be brought on-line immediately should a potential competitor contemplate entry. For example, a ski-resort owner may invest in a new chair-lift, even if she does not use it frequently. The existence of the additional capacity may scare potential entrants. A key component of this strategy is that the incumbent firm invests ahead of time – and inflicts a cost on itself. The incumbent does not simply say "I will build another chair-lift if you decide to develop a nearby mountain into a ski hill." That policy does not carry the same degree of credibility as actually incurring the cost of construction ahead of time. However, such a strategy may not always be feasible: It might be just too costly to pre-empt entry by putting spare capacity in place. Spare capacity is not so different from brand development through advertising; both are types of sunk cost. The threats associated with the incumbent's behaviour become a credible threat because the incumbent incurs costs up front.

A **credible threat** is one that is effective in deterring specific behaviours; a competitor must believe that the threat will be implemented if the competitor behaves in a certain way.

Lobbying

In our chapter on monopoly we stressed the role of political/lobbying activity. Large firms invariably employ public relations firms, and maintain their own public relations departments. The role of these units is not simply to portray a positive image of the corporation to the public; it is to maintain and increase whatever market power such firms already possess. It is as much in the interest of an oligopolistic firm as a monopolist to prevent entry and preserve supernormal profits.

In analyzing perfect competition, we saw that free entry is critical to maintaining normal profits. Lobbying is designed to obstruct entry, and it is also designed to facilitate mergers and acquisitions. The economist Thomas Philippon has written about the increasing concentration of economic power in recent decades in the hands of a small number of corporations in many sectors of

the North American economy. He argues that this concentration of power contributes to making the distribution of income more favorable to corporate interests and less favorable to workers. In his recent book ("The Great Reversal: How America Gave up Free Markets"), he shows that, contrary to traditional beliefs, Europe is now much more competitive than the US in most sectors of the economy. More broadband suppliers result in rates in Europe that are about half of US rates. Whereas in the US four airlines control 80% of the market, In Europe they control 40%. If scale economies were the prime determinant of corporate concentration we should not expect such large differences. Likewise, if globalization and technological change were the main determinants of corporate concentration, we should expect experiences in Europe and North America to be similar. But they are not. Hence, it is reasonable to conclude that entry barriers in North America are more effective, or that regulatory forces are stronger in Europe.

This page titled [11.8: Strategic behaviour- Entry, exit and potential competition](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.9: Matching markets- design

Markets are institutions that facilitate the exchange of goods and services. They act as clearing houses. The normal medium of exchange is money in some form. But many markets deal in exchanges that do not involve money and frequently involve matching: Graduating medical students are normally matched with hospitals in order that graduates complete their residency requirement; in many jurisdictions in the US applicants for places in public schools that form a pool within a given school-board must go through an application process that sorts the applicants into the different schools within the board; patients in need of a new kidney must be matched with kidney donors.

These markets are clearinghouses and have characteristics that distinguish them from traditional currency-based markets that we have considered to this point.

- The good or service being traded is generally *heterogeneous*. For example, patients in search of a kidney donor must be medically compatible with the eventual donor if the organ transplant is not to be rejected. Hospitals may seek residents in particular areas of health, and they must find residents who are, likewise, seeking such placements. Students applying to public schools may be facing a choice between schools that focus upon science or upon the arts. Variety is key.
- Frequently the idea of a market that is mediated by money is *repugnant*. For example, the only economy in the world that permits the sale of human organs is Iran. Elsewhere the idea of a monetary payment for a kidney is unacceptable. A market in which potential suppliers of kidneys registered their reservation prices and demanders registered their willingness to pay is incompatible with our social mores. Consequently, potential living donors or actual deceased donors must be directly matched with a patient in need. While some individuals believe that a market in kidneys would do more good than harm, because a monetary payment might incentivize the availability of many more organs and therefore save many more lives, virtually every society considers the downside to such a trading system to outweigh the benefits.
- Modern matching markets are more frequently *electronically mediated*, and the communications revolution has led to an increase in the efficiency of these markets.

The Economics prize in memory of Alfred Nobel was awarded to Alvin Roth and Lloyd Shapley in 2011 in recognition of their contributions to designing markets that function efficiently in the matching of demanders and suppliers of the goods and services. What do we mean by an efficient mechanism? One way is to define it is similar to how we described the market for apartments in Chapter 5: following an equilibrium in the market, is it possible to improve the wellbeing of one participant without reducing the wellbeing of another? We showed in that example that the market performed efficiently: a different set of renters getting the apartments would reduce total surplus in the system.

Consider a system in which medical graduates are matched with hospitals, and the decision process results in the potential for improvement: Christina obtains a residency in the local University Hospital while Ulrich obtains a residency at the Childrens' Hospital. But Christina would have preferred the Childrens' and Ulrich would have preferred the University. The matching algorithm here was not efficient because, at the end of the allocation process, there is scope for gains for each individual. Alvin Roth devised a matching mechanism that surmounts this type of inefficiency. He called it the deferred acceptance algorithm.

Roth also worked on the matching of kidney donors to individuals in need of a kidney. The fundamental challenge in this area is that a patient in need of a kidney may have a family member, say a sibling, who is willing to donate a kidney, but the siblings are not genetically compatible. The patient's immune system may attack the implantation of a 'foreign' organ. One solution to this incompatibility is to find matching pairs of donors that come from a wider choice set. Two families in each of which there is patient and a donor may be able to cross-donate: donor in Family A can donate to patient in Family B, and donor in Family B can donate to patient in Family A, in the sense that the donor organs will not be rejected by recipients' immune systems. Hence if many patient-donor families register in a clearinghouse, a computer algorithm can search for matching pairs. Surgical operations may be performed simultaneously in order to prevent one donor from backing out following his sibling's receipt of a kidney.

A more recent development concerns 'chains'. In this case a good Samaritan ('unaligned donor') offers a kidney while seeking nothing in return. The algorithm then seeks a match for the good Samaritan's kidney among all of the recipient-donor couples registered in the data bank. Having found (at least) one, the algorithm seeks a recipient for the kidney that will come from the first recipient's donor partner. And so on. It turns out that an algorithm which seeks to maximize the potential number of participating pairs is fraught with technical and ethical challenges: should a young patient, who could benefit from the organ for a whole lifetime, get priority over an older patient, who will benefit for fewer years of life, even if the older patient is in greater danger of dying in the absence of a transplant? This is an ethical problem.

Examples where these algorithms have achieved more than a dozen linked transplants are easy to find on an internet search - they are called *chains*, for the obvious reason.

Consider the following efficiency aspect of the exchange. Suppose a patient has two siblings, each of whom is willing to donate (though only one of the two actually will); should such a patient get priority in the computer algorithm over a patient who has just a single sibling willing to donate? The answer may be yes; the dual donor patient should get priority because if his two siblings have different blood types, this greater variety on the supply side increases the chances for matching in the system as a whole and is therefore beneficial. If a higher priority were not given to the dual-donor patient, there would be an incentive for him to name just one potential donor, and that would impact the efficiency of the whole matching algorithm.

It is not always recognized that the discipline of Economics explores social problems of the nature we have described here, despite the fact that the discipline has developed the analytical tools to address them.

This page titled [11.9: Matching markets- design](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.10: Conclusion

Monopoly and perfect competition are interesting paradigms; but few markets resemble them in the real world. In this chapter we addressed some of the complexities that define the economy we inhabit: It is characterized by strategic planning, entry deterrence, differentiated products and so forth.

Entry and exit are critical to competitive markets. Frequently entry is blocked because of scale economies – an example of a natural or unintended entry barrier. In addition, incumbents can formulate numerous strategies to limit entry.

Firms act strategically – particularly when there are just a few participants in the market. Before acting, firms make conjectures about how their competitors will react, and incorporate such reactions into their own planning. Competition between suppliers can frequently be analyzed in terms of a game, and such games usually have an equilibrium outcome. The Cournot duopoly model that we developed is a game between two competitors in which an equilibrium market output is determined from a pair of reaction functions.

Scale economies are critical. Large development costs or setup costs may mean that the market can generally support just a limited number of producers. In turn this implies that potential new (small-scale) firms cannot benefit from the scale economies and will not survive competition from large-scale suppliers.

Product differentiation is critical. If small differences exist between products produced in markets where there is free entry we get a monopolistically competitive structure. In these markets long-run profits are 'normal' and firms operate with some excess capacity. It is not possible to act strategically in this kind of market.

The modern economy also has sectors that have successfully erected barriers. These barriers lead to fewer competitors than could efficiently supply the market. Ultimately the owners of capital are the beneficiaries of these barriers and consumers suffer from higher prices.

This page titled [11.10: Conclusion](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.11: Key Terms

Imperfectly competitive firms face a downward-sloping demand curve, and their output price reflects the quantity sold.

Oligopoly defines an industry with a small number of suppliers.

Monopolistic competition defines a market with many sellers of products that have similar characteristics. Monopolistically competitive firms can exert only a small influence on the whole market.

Duopoly defines a market or sector with just two firms.

Concentration ratio: N -firm concentration ratio is the sales share of the largest N firms in that sector of the economy.

Differentiated product is one that differs slightly from other products in the same market.

The **monopolistically competitive equilibrium** in the long run requires the firm's demand curve to be tangent to the ATC curve at the output where $MR=MC$.

Collusion is an explicit or implicit agreement to avoid competition with a view to increasing profit.

Conjecture: a belief that one firm forms about the strategic reaction of another competing firm.

Game: a situation in which contestants plan strategically to maximize their profits, taking account of rivals' behaviour.

Strategy: a game plan describing how a player acts, or moves, in each possible situation.

Nash equilibrium: one in which each player chooses the best strategy, given the strategies chosen by the other player, and there is no incentive for any player to move.

Dominant strategy: a player's best strategy, whatever the strategies adopted by rivals.

Payoff matrix: defines the rewards to each player resulting from particular choices.

Credible threat: one that, after the fact, is still optimal to implement.

Cournot behaviour involves each firm reacting optimally in their choice of output to their competitors' decisions.

Reaction functions define the optimal choice of output conditional upon a rival's output choice.

This page titled [11.11: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

11.12: Exercises for Chapter 11

EXERCISE 11.1

Imagine that the biggest four firms in each of the sectors listed below produce the amounts defined in each cell. Compute the three-firm and four-firm concentration ratios for each sector, and rank the sectors by degree of industry concentration.

Sector	Firm 1	Firm 2	Firm 3	Firm 4	Total market
Shoes	60	45	20	12	920
Chemicals	120	80	36	24	480
Beer	45	40	3	2	110
Tobacco	206	84	30	5	342

EXERCISE 11.2

You own a company in a monopolistically competitive market. Your marginal cost of production is \$12 per unit. There are no fixed costs. The demand for your own product is given by the equation $P=48-(1/2)Q$.

- Plot the demand curve, the marginal revenue curve, and the marginal cost curve.
- Compute the profit-maximizing output and price combination.
- Compute total revenue and total profit [*Hint*: Remember $AC=MC$ here].
- In this monopolistically competitive industry, can these profits continue indefinitely?

EXERCISE 11.3

Two firms in a particular industry face a market demand curve given by the equation $P=100-(1/3)Q$. The marginal cost is \$40 per unit and the marginal revenue is $MR=100-(2/3)Q$. The quantity intercepts for demand and MR are 300 and 150.

- Draw the demand curve and MR curve to scale on a diagram. Then insert the MC curve.
- If these firms got together to form a cartel, what output would they produce and what price would they charge?
- Assuming they each produce half of the total what is their individual profit?

EXERCISE 11.4

The classic game theory problem is the "prisoners' dilemma." In this game, two criminals are apprehended, but the police have only got circumstantial evidence to prosecute them for a small crime, without having the evidence to prosecute them for the major crime of which they are suspected. The interrogators then pose incentives to the crooks-incentives to talk. The crooks are put in separate jail cells and have the option to confess or deny. Their payoff depends upon what course of action each adopts. The payoff matrix is given below. The first element in each box is the payoff (years in jail) to the player in the left column, and the second element is the payoff to the player in the top row.

		B's strategy	
		Confess	Deny
A's strategy	Confess	6,6	0,10
	Deny	10,0	1,1

- Does a "dominant strategy" present itself for each or both of the crooks?
- What is the Nash equilibrium to this game?
- Is the Nash equilibrium unique?
- Was it important for the police to place the crooks in separate cells?

EXERCISE 11.5

Taylormade and Titlelist are considering a production strategy for their new golf drivers. If they each produce a small output, they can price the product higher and make more profit than if they each produce a large output. Their payoff/profit matrix is given below.

		Taylormade strategy	
		Low output	High output
Titleist strategy	Low output	50,50	20,70
	High output	70,20	40,40

- Does either player have a dominant strategy here?
- What is the Nash equilibrium to the game?
- Do you think that a cartel arrangement would be sustainable?

EXERCISE 11.6

Ronnie's Wraps is the only supplier of sandwich food and makes a healthy profit. It currently charges a high price and makes a profit of six units. However, Flash Salads is considering entering the same market. The payoff matrix below defines the profit outcomes for different possibilities. The first entry in each cell is the payoff/profit to Flash Salads and the second to Ronnie's Wraps.

		Ronnie's Wraps	
		High price	Low price
Flash Salads	Enter the market	2,3	-1,1
	Stay out of market	0,6	0,4

- If Ronnie's Wraps threatens to lower its price in response to the entry of a new competitor, should Flash Salads stay away or enter?
- Explain the importance of threat credibility here.

EXERCISE 11.7

Optional: Consider the market demand curve for appliances: $P = 3,200 - (1/4)Q$. There are no fixed production costs, and the marginal cost of each appliance is $MC = \$400$. As usual, the MR curve has a slope that is twice as great as the slope of the demand curve.

- Illustrate this market geometrically.
- Determine the output that will be produced in a 'perfectly competitive' market structure where no profits accrue in equilibrium.
- If this market is supplied by a monopolist, illustrate the choice of output.

EXERCISE 11.8

Optional: Consider the outputs you have obtained in Exercise 11.7.

- Can you figure out how many firms would produce at the perfectly competitive output? If not, can you think of a reason?
 - If, in contrast, each firm in that market had to cover some fixed costs, in addition to the variable costs defined by the MC value, would that put a limit on the number of firms that could produce in this market?
- This presentation is inspired by a similar problem in Ted Bergstrom and Hal Varian's book "Workouts".
 - This game is sometimes called the "prisoners' dilemma" game because it can be constructed to reflect a scenario in which two prisoners, under isolated questioning, each confess to a crime and each ends up worse off than if neither confessed.
 - <https://slate.com/technology/2013/10/amazon-book-how-jeff-bezos-went-thermonuclear-on-diapers-com.html>

This page titled [11.12: Exercises for Chapter 11](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

SECTION OVERVIEW

Unit 5: The Factors of Production

Output is produced by the factors of production – labour, capital, and land. In Chapter 12 we explore how the markets for these factors, or productive inputs, operate. We focus primarily on labour and capital.

Chapter 13 deals with human capital – a broader concept than labour. The human capital embodied in a person depends upon her skills, education and experience, and this combination determines the rewards, or return, she obtains when employed. This chapter also explores the distribution of income in Canada and the degree of inequality in that distribution. It concludes with a brief overview of the 'top one percent'.

12: Labour and capital

- 12.1: Labour - a derived demand
- 12.2: The supply of labour
- 12.3: Labour market equilibrium and mobility
- 12.4: Capital - concepts
- 12.5: The capital market
- 12.6: Land
- 12.7: Key Terms
- 12.8: Exercises for Chapter 12

13: Human capital and the income distribution

- 13.1: Human capital
- 13.2: Productivity and education
- 13.3: On-the-job training
- 13.4: Education as signalling
- 13.5: Education returns and quality
- 13.6: Discrimination
- 13.7: The income distribution
- 13.8: Wealth and capitalism
- 13.9: Key Terms
- 13.10: Exercises for Chapter 13

This page titled [Unit 5: The Factors of Production](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12: Labour and capital

Chapter 12: Labour and capital

In this chapter we will explore:

12.1	The demand for labour
12.2	Labour supply
12.3	Market equilibrium and labour mobility
12.4	The concepts of capital
12.5	The capital market
12.6	Land

The chapter deals with the markets for the factors of production—labour and capital. The analysis will be presented in terms of the demand, supply and market equilibrium for each. While this is a standard analytical approach in microeconomics, the markets for labour and capital differ from goods and services markets.

In the first instance, goods and services are purchased and consumed by the buyers. In contrast, labour and capital are used as inputs in producing those 'final' goods and services. So the value of labour and capital to a producer depends in part upon the value of the products that the labour and capital are used to produce. Economists say that the value of the factors of production *derives* from the value of the products they ultimately produce.

Secondly, labour and capital offer services. When an employer hires a worker, that worker supplies her time and skills and energy to the employer. When a piece of equipment is rented to a producer, or purchased by a producer, that equipment provides a stream of productive services also. The employer does not purchase the worker, given that we do not live in a society where slavery is legal. In contrast, she may decide to purchase the capital, or else rent it.

The third characteristic of these markets is the time dimension associated with labour and capital. Specifically, once built, a machine will customarily have a lifetime of several years, during which it depreciates in value. Furthermore, it may become obsolete on account of technological change before the end of its anticipated life. Labour too may become obsolete, or at least lose some of its value with the passage of time, if the skills embodied in the labour cease to be required in the economy.

12.1 Labour – a derived demand

The value of labour springs from the value of its use, that is the value placed upon goods and services that it produces – product prices. The wage is the price that equilibrates the supply and demand for a given type of labour, and it reflects the value of that labour in production. Formally, the demand for labour (and capital) is thus a derived demand, in contrast to being a 'final' demand.

Demand for labour: a derived demand, reflecting the value of the output it produces.

We must distinguish between the long run and the short run in our analysis of factor markets. On the *supply side* certain factors of production are fixed in the short run. For example, the supply of radiologists can be increased only over a period of years. While one hospital may be able to attract radiologists from another hospital to meet a shortage, this does not increase the supply in the economy as a whole.

On the *demand side* there is the conventional difference between the short and long run: In the short run some of a firm's factors of production, such as capital, are fixed, and therefore the demand for labour differs from when all factors are variable – the long run.

Demand in the short run

Table 12.1 contains information from the example developed in Chapter 8. It can be used to illustrate how a firm reacts in the short run to a change in an input price, or to a change in the output price. The response of a producer to a change in the wage rate constitutes a demand function for labour – a schedule relating the quantity of the input demanded to different input prices. The output produced by the various numbers of workers yields a marginal product curve, whose values are stated in column 3. The marginal product of labour, MP_L , as developed in Chapter 8, is the additional output resulting from one more worker being employed, while holding constant the other (fixed) factors. But what is the *dollar value* to the firm of an additional worker? It is the

additional *value of output* resulting from the additional employee – the price of the output times the worker's marginal contribution to output, his *MP*. We term this the value of the marginal product.

The **value of the marginal product** is the marginal product multiplied by the price of the good produced.

Table 12.1 Short-run production and labour demand

Workers	Output	MP_L	$VMP_L = MP_L \times P$	Marginal profit = ($VMP_L - \text{wage}$)
(1)	(2)	(3)	(4)	(5)
0	0			
1	15	15	1050	150
2	40	25	1750	750
3	70	30	2100	1100
4	110	40	2800	1800
5	145	35	2450	1450
6	175	30	2100	1100
7	200	25	1750	750
8	220	20	1400	400
9	235	15	1050	50
10	240	5	350	negative

Each unit of labour costs \$1,000; output sells at a fixed price of \$70 per unit.

In this example the MP_L first rises as more labour is employed, and then falls. With each unit of output selling for \$70 the value of the marginal product of labour (VMP_L) is given in column 4. The first worker produces 15 units each week, and since each unit sells for a price of \$70, his production value to the firm is \$1,050 ($= \70×15). A second worker produces 25 units, so his value to the firm is \$1,750, and so forth. If the weekly wage of each worker is \$1,000 then the firm can estimate its marginal profit from hiring each additional worker. This is the difference between the value of the marginal product and the wage paid, and is given in the final column of the table.

It is profitable to hire more workers as long as the cost of an extra worker is less than the VMP_L . The equilibrium amount of labour to employ is therefore 9 units in this example. If the firm were to hire one more worker the contribution of that worker to its profit would be negative (\$350 – \$1,000), and if it hired one worker less it would forego the opportunity to make an additional profit of \$50 on the 9th unit (\$1,050 – \$1,000).

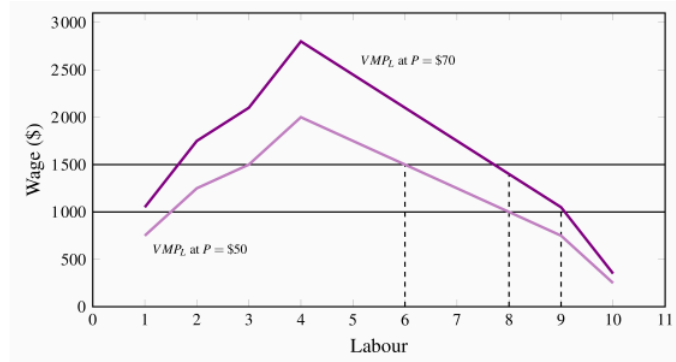
Profit maximizing hiring rule:

- If the VMP_L of next worker > wage, hire more labour.
- If the $VMP_L < \text{wage}$, hire less labour.

To this point we have determined the profit maximizing amount of labour to employ when the output price and the wage are given. However, a demand function for labour reflects the demand for labour at many different wage rates, just as the demand function for any product reflects the quantity demanded at various prices. Accordingly, suppose the wage rate is \$1,500 per week rather than \$1,000. The optimal amount of labour to employ in this case is determined in exactly the same manner: Employ the amount of labour where its contribution is marginally profitable. Clearly the optimal amount to employ is 7 units: The value of the seventh worker to the firm is \$1,750 and the value of the eighth worker is \$1,400. Hence it would not be profitable to employ the eighth, because his marginal contribution to profit would be negative. Following the same procedure we could determine the optimal amount of labour to employ at different wages. This implies that the VMP_L function is the demand for labour function because it determines the most profitable amount of labour to employ at any wage.

The optimal amount of labour to hire is illustrated in Figure 12.1. The wage and VMP_L curves come from Table 12.1. The VMP_L curve has an upward sloping segment, reflecting increasing productivity, and then a regular downward slope as developed in Chapter 8. At employment levels where the VMP_L is greater than the wage additional labour should be employed. But when the VMP_L falls below the wage rate employment should stop. If labour is divisible into very small units, the optimal employment decision is where the MP_L function intersects the wage line.

Figure 12.1 The demand for labour



The optimal hiring decision is defined by the condition that the value of the MP_L is greater than or equal to the wage paid. 9 workers are employed when the wage is \$1,000 and the price of output is \$70; 6 workers are employed when the wage is \$1,500 and the price of output is \$50.

Figure 12.1 also illustrates what happens to hiring when the output price changes. Consider a reduction in its price to \$50 from \$70. The profit impact of such a change is negative because the value of each worker's output has declined. Accordingly, the demand curve must reflect this by shifting inward (down), as in the figure. At various wage rates, less labour is now demanded. The new VMP_L schedule can be derived in Table 12.1 as before: It is the MP_L schedule multiplied by the lower value (\$50) of the final good.

In this example the firm is a perfect competitor in the output market, because the price of the good it produces is fixed. It can produce and sell more of the good without this having an impact on the price of the good in the marketplace. Where the firm is not a perfect competitor it faces a declining MR function. In this case the value of the MP_L is the product of MR and MP_L rather than P and MP_L . To distinguish the different output markets we use the term *marginal revenue product of labour* (MRP_L) when the demand for the output slopes downward. But the optimizing principle remains the same: The firm should calculate the value of each additional unit of labour, and hire up to the point where the additional revenue produced by the worker exceeds or equals the additional cost of that worker.

The **marginal revenue product of labour** is the additional revenue generated by hiring one more unit of labour where the marginal revenue declines.

Demand in the long run

In Chapter 8 we proposed that firms choose their factors of production in accordance with cost-minimizing principles. In producing a specific output, firms choose the least-cost combination of labour and plant size. But how is this choice affected when the price of labour or capital changes? While adjustment to price changes may require a long period of time, we know that if one factor becomes more (less) expensive, the firm will likely change the mix of capital and labour away from (towards) that factor. For example, when the accuracy and prices of production robots began to fall in the nineteen nineties, auto assemblers reduced their labour and used robots instead. When computers and computer software improved and declined in price, clerical workers were replaced by computers that were operated by accountants. But such adjustments and responses do not occur overnight.

In the short run a higher wage increases costs, but the firm is constrained in its choice of inputs by a fixed plant size. In the long run, a wage increase will induce the firm to use relatively more capital than when labour was less expensive in producing a given output. But despite the new choice of inputs, a rise in the cost of any input must increase the total cost of producing any output.

A change in the price of any factor has two impacts on firms: In the first place producers will *substitute* away from the factor whose price increases; second, there will be an impact on output and a change in the price of the final good it produces. Since the cost structure increases when the price of an input rises, the supply curve in the market for the good must reflect this – any given output will now be supplied at a higher price. With a downward sloping demand, this shift in supply must increase the price of the good and reduce the amount sold. This second effect can be called an *output effect*.

Monopsony

Some firms may have to pay a higher wage in order to employ more workers. Think of Hydro Quebec building a dam in Northern Quebec. Not every hydraulic engineer would be equally happy working there as in Montreal. Some engineers may demand only a small wage premium to work in the North, but others will demand a high premium. If so, Hydro Quebec must pay a higher wage to attract more workers – it faces an upward sloping supply of labour curve. Hydro Quebec is the sole buyer in this particular market and is called a monopsonist – a single buyer. Our general optimizing principle governing the employment of labour still holds, even if we have different names for the various functions: Hire any factor of production up to the point where the cost of an additional unit equals the value generated for the firm by that extra worker. The essential difference here is that when a firm faces an upward sloping labour supply it will have to pay more to attract additional workers and *also pay more to its existing workers*. This will impact the firm's willingness to hire additional workers.

A **monopsonist** is the sole buyer of a good or service and faces an upward-sloping supply curve.

Application Box 12.1 Monopsonies

Monopsonies are more than a curiosity; they exist in the real world. An excellent example is the cannabis market in Canada. Virtually every province has set up a trading agency that has the sole right to purchase cannabis from growers; growers and processors are not permitted to sell directly to retailers; they may only sell to the monopsony by law. In turn, these provincial cannabis monopsonies are frequently retail monopolists in that the agency owns all of the retail outlets in the province.

Firm versus industry demand

The demand for labour within an industry, or sector of the economy, is obtained from the sum of the demands by each individual firm. It is analogous to the goods market, but with a subtle difference. In Chapter 3 we obtained a market demand by summing individual demands horizontally. The same could be done here: At lower (or higher) wages, each firm will demand more (or less) labour. However, if all firms employ more labour in order to increase their output, the price of the output will likely decline. This in turn will moderate the demand for labour – it is slightly less valuable now that the price of the output it produces has fallen. This is a subtle point, and we can reasonably think of the demand for labour in a given sector of the economy as the sum of the demands on the part of the employers in that sector.

12.2 The supply of labour

Most prime-age individuals work, but some do not. The decision to join the labour force is called the participation decision. Of those who do participate in the labour force, some individuals work full time, others work part time, and yet others cannot find a job. The unemployment rate is the fraction of the labour force actively seeking employment that is not employed.

The **participation rate** for the economy is the fraction of the population in the working age group that joins the labour force.

The **labour force** is that part of the population either employed or seeking employment.

The **unemployment rate** is the fraction of the labour force actively seeking employment that is not employed.

Data on participation rates in Canada are given in Table 12.2 below for specific years in the modern era. The overall participation for men and women combined has increased since 1977 from 60.8% to 65.8% This aggregated rate camouflages different patterns for men and women. The rates for women have been rising while the rates for men have fallen. Women today are more highly educated, and their role in society and the economy is viewed very differently than in the earlier period. Female participation has increased both because of changing social norms, a rise in household productivity, the development of service industries designed to support home life, and the development of the institution of daycare for young children.

In contrast, male participation rates declined over the period, largely offsetting the increase in female participation. Fewer individuals in total are retiring before the age of 55 in the most recent decades. This reflects both the greater number of females in the market place, and perhaps also a recognition that many households have not saved enough to fund a retirement period that has become longer as a result of increased longevity.

Table 12.2 Labour force participation rate, Canada 1977-2015

Year	Total	Men	Women	All > 55	Unemployment
1977	60.8	80.2	42.1	30.5	5.9
1990	66.6	77.1	56.8	25.9	7.4

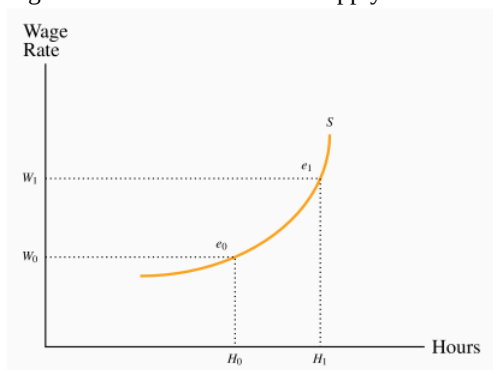
1994	65.4	74.9	56.8	24.5	9.2
2001	66.1	73.4	59.2	26.4	6.2
2008	67.5	73.6	61.6	34.2	5.1
2015	66.2	72.0	60.6	37.3	6.0
2019	65.8	71.0	60.8	38.0	4.4

Source: Statistics Canada, CANSIM 14-10-0287-02
September of each year, for individuals aged ≥ 25 , unless stated.

At the micro level, the participation rate of individuals depends upon several factors. First, the wage rate that an individual can earn in the market is crucial. If that wage is low, then the individual may be more efficient in producing home services directly, rather than going into the labour market, earning a modest income and having to pay for home services. Second, there are fixed costs associated with working. A decision to work means that the individual must have work clothing, must undertake the costs of travel to work, and pay for daycare if there are children in the family. Third, the participation decision depends upon non-labour income. If the individual in question has a partner who earns a substantial amount, or if she has investment income, she will have less incentive to participate. Fourth, it depends inversely upon the tax rate.

The supply curve relates the supply decision to the price of labour – the wage rate. Economists who have studied the labour market tell us that the individual supply curve is upward sloping: As the wage increases, the individual wishes to supply more labour. From the point e_0 on the supply function in Figure 12.2, let the wage increase from W_0 to W_1 .

Figure 12.2 Individual labour supply



A wage increase from W_0 to W_1 induces the individual to *substitute* away from leisure, which is now more expensive, and work *more*. But the higher wage also means the individual can work fewer hours for a given standard of living; therefore the income effect induces *fewer hours*. On balance the substitution effect tends to dominate and the supply curve therefore slopes upward.

The individual offers more labour, H_1 , at the higher wage. What is the economic intuition behind the higher amount of labour supplied? Like much of choice theory there are two impacts associated with a higher price. First, the higher wage makes leisure more expensive relative to working. That midweek game of golf has become more expensive in terms of what the individual could earn. So the individual should *substitute* away from the more expensive 'good', leisure, towards labour. But at the same time, in order to generate a given income target the individual can work fewer hours at the higher wage. This is a type of *income effect*, indicating that income is greater at a higher wage regardless of the amount worked, and this induces the individual to work less. The fact that we draw the labour supply curve with a positive slope means that the substitution effect is the more important of the two. That is what statistical research has revealed.

Elasticity of the supply of labour

The value of the supply elasticity depends upon how the market in question is defined. In particular, it depends upon how large or small a given sector of the economy is, and whether we are considering the short run or the long run.

Suppose an industry is small relative to the whole economy and employs workers with common skills. These industries tend to pay the 'going wage'. For example, very many students are willing to work at the going rate for telemarketing firms, which compose a small sector of the economy. This means that the supply curve of such labour, *as far as that sector is concerned*, is in effect horizontal – infinitely elastic.

But some industries may not be small relative to the total labour supply. And in order to get more labour to work in such large sectors it may be necessary to provide the inducement of a higher wage: Additional workers may have to be attracted from another sector by means of higher wages. To illustrate: Consider the behaviour of two related sectors in housing – new construction and home restoration. In order to employ more plumbers and carpenters, new home builders may have to offer higher wages to induce them to move from the renovation sector. In this case the new housing industry's labour supply curve slopes upwards.

In the time dimension, a longer period is always associated with more flexibility. In this context, the supply of labour to any sector is more elastic, because it may take time for workers to move from one sector to another. Or, in cases where skills must be built up: When a sectoral expansion bids up the wages of information technology (IT) workers, more school leavers are likely to develop IT skills. Time will be required before additional graduates are produced, but in the long run, such additional supply will moderate the short-run wage increases.

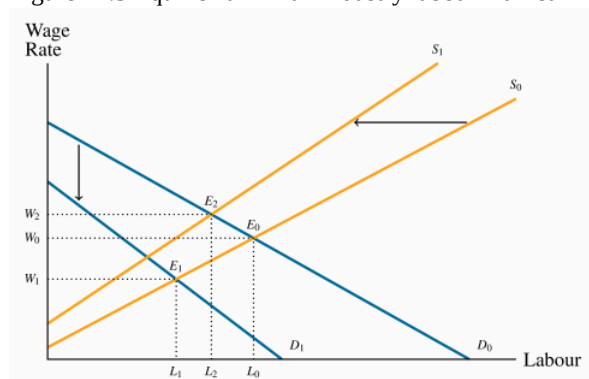
Wages can be defined as being before-tax or after-tax. The after-tax, or take-home, wage is more important than the gross wage in determining the quantity of labour to be supplied. If taxes on additional hours of work are very high, workers are more likely to supply less hours than if tax rates are lower.

12.3 Labour market equilibrium and mobility

The fact that labour is a derived demand differentiates the labour market's equilibrium from the goods-market equilibrium. Let us investigate this with the help of Figure 12.3; it contains supply and demand functions for one particular industry – the cement industry, let us assume.

In Figure 12.1 we illustrated the impact on the demand for labour of a decline in the price of the output produced – a decline in the output price reduced the value of the marginal product of labour. In the current example, suppose that a slowdown in construction results in a decline in the price of cement. The impact of this price fall is to reduce the output value of each worker in the cement producing industry, because their output now yields a lower price. This decline in the VMP_L is represented in Figure 12.3 as a shift from D_0 to D_1 , which results in the new equilibrium E_1 .

Figure 12.3 Equilibrium in an industry labour market



A fall in the *price of the good* produced in a particular industry reduces the value of the MP_L . Demand for labour thus falls from D_0 to D_1 and a new equilibrium E_1 results. Alternatively, from E_0 , an increase in wages in another sector of the economy induces some labour to move to that sector. This is represented by the shift of S_0 to S_1 and the new equilibrium E_2 .

As a second example: Suppose that wages in some other sectors of the economy increase. The impact of this on the cement sector is that the supply of labour to the cement sector is reduced. In Chapter 3 we showed that a change in other prices may *shift* the demand or supply curve of interest. In Figure 12.3 supply shifts from S_0 to S_1 and the equilibrium goes from E_0 to E_2 .

How large are these impacts likely to be? That will depend upon how mobile labour is between sectors: Spillover effects will be smaller if labour is less mobile. This brings us naturally to the concepts of transfer earnings and rent.

Transfer earnings and rent

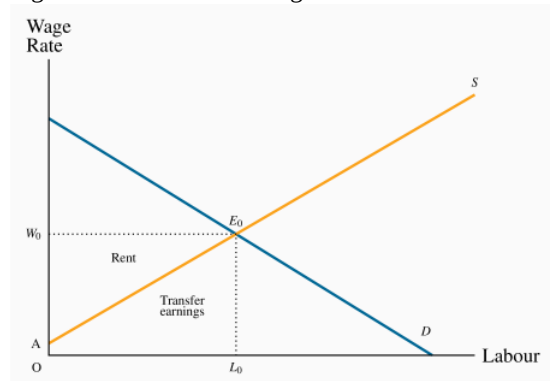
Consider the case of a performing violinist whose wage is \$80,000. If, as a best alternative, she can earn \$60,000 as a music teacher then her rent is \$20,000 and her transfer earnings \$60,000: Her rent is the excess she currently earns above the best alternative. Another violinist in the same orchestra, earning the same amount, who could earn \$55,000 as a teacher has rent of \$25,000. The alternative is also called the reservation wage. The violinists should not work in the orchestra unless they earn at least what they can earn in the next best alternative.

Transfer earnings are the amount that an individual can earn in the next highest paying alternative job.

Rent is the excess remuneration an individual currently receives above the next best alternative. This alternative is the **reservation wage**.

These concepts are illustrated in Figure 12.4. In this illustration, different individuals are willing to work for different amounts, but all are paid the same wage w_0 . The market labour supply curve by definition defines the wage for which each individual is willing to work. Thus the rent earned by labour in this market is the sum of the excess of the wage over each individual's transfer earnings – the area W_0E_0A . This area is also what we called producer or supplier surplus in Chapter 5.

Figure 12.4 Transfer earnings and rent



Rent is the excess of earnings over reservation wages. Each individual earns W_0 and is willing to work for the amount defined by the labour supply curve. Hence rent is W_0E_0A and transfer earnings OAE_0L_0 . Rent is thus the term for supplier surplus in this market.

Free labour markets?

Real-world labour markets are characterized by trade unions, minimum wage laws, benefit regulations, severance packages, parental leave, sick-day allowances and so forth. So can we really claim that markets work in the way we have described them – essentially as involving individual agents demanding and supplying labour? While labour markets are not completely 'free' in the conventional sense, the important issue is whether these interventions, that are largely designed to protect workers, have a large or small impact on the market. One reason why unemployment rates are generally higher in European economies than in Canada and the US is that labour markets are less subject to controls, and workers have a less supportive social safety net in North America.

Application Box 12.2 Are high salaries killing professional sports?

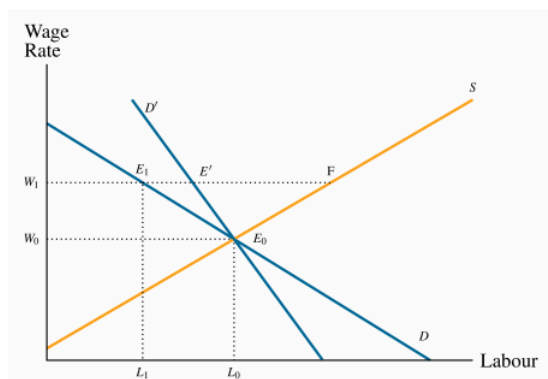
It is often said that the agents of professional players are killing their sport by demanding unreasonable salaries. On occasion, the major leagues are threatened with strikes, even though players are paid millions each year. In fact, wages are high because the derived demand is high. Fans are willing to pay high ticket prices, and television rights generate huge revenues. Combined, these revenues not only make ownership profitable, but increase the demand for the top players.

The lay person may be horrified at thirty-million dollar annual salaries. But in reality, many players receiving such salaries may be earning less than their marginal product! If Tom Brady did not play for the New England Patriots the team would have a lower winning record, attract fewer fans and make less profit. If Brady is paid \$25m per season, he is being paid less than his marginal product if the team were to lose \$40m in revenue as a result of his absence.

Given this, why do some teams incur financial losses? In fact very few teams make losses: Cries of poverty on the part of owners are more frequently part of the bargaining process, and revenue sharing means that very few teams do not make a profit.

The impact of 'frictions', such as unionization and minimum wages, in the labour market can be understood with the help of Figure 12.5. The initial 'free market' equilibrium is at E_0 , assuming that the workers are not unionized. In contrast, if the workers in this industry form a union, and negotiate a higher wage, for example w_1 rather than w_0 , then fewer workers will be employed. But how big will this reduction be? Clearly it depends on the elasticities of demand and supply. With the demand curve D , the excess supply at the wage w_1 is the difference E_1F . However, if the demand curve is less elastic, as illustrated by the curve D' , the excess supply is $E'F$. The excess supply also depends upon the supply elasticity. It is straightforward to see that a less elastic (more vertical) supply curve through E_0 would result in less excess supply.

Figure 12.5 Market interventions



E_0 is the equilibrium in the absence of a union. If the presence of a union forces the wage to W_1 fewer workers are employed. The magnitude of the decline from L_0 to L_1 depends on the elasticity of demand for labour. The excess supply at the wage W_1 is $(F-E_1)$. With a less elastic demand curve (D') the excess supply is reduced to $(F-E')$.

Beyond elasticity, the magnitude of the excess supply will also depend upon the degree to which the minimum wage, or the union-negotiated wage, lies above the equilibrium. That is, a larger value of the difference $(W_1 - W_0)$ results in more excess supply than a smaller difference.

While the above discussion pertains to unionization, it could equally well be interpreted in a minimum-wage context. If this figure describes the market for low-skill labour, and the government intervenes by setting a legal minimum at W_1 , then this will induce some degree of excess supply, depending upon the actual value of W_1 and the elasticities of supply and demand.

Despite the fact that a higher wage may induce some excess supply, it may increase total earnings. In Chapter 4 we saw that the dollar value of expenditure on a good increases when the price rises if the demand is inelastic. In the current example the 'good' is labour. Hence, a union-negotiated wage increase, or a higher minimum wage will each increase total remuneration if the demand for labour is inelastic. A case which has stirred great interest is described in Application Box 12.3.

Application Box 12.3 David Card on minimum wage

David Card is a famous Canadian-born labour economist who has worked at Princeton University and University of California, Berkeley. He is a winner of the prestigious Clark medal, an award made annually to an outstanding economist under the age of forty. Among his many contributions to the discipline, is a study of the impact of minimum wage laws on the employment of fast-food workers. With Alan Krueger as his co-researcher, Card examined the impact of the 1992 increase in the minimum wage in New Jersey and contrasted the impact on employment changes with neighbouring Pennsylvania, which did not experience an increase. They found virtually no difference in employment patterns between the two states. This research generated so much interest that it led to a special conference. Most economists now believe that modest changes in the level of the minimum wage have a small impact on employment levels.

Since about 2015, numerous labor-friendly movements favoring higher wages for low-paid workers have proposed a \$15 minimum in both Canada and the US. Some political parties have supported this movement, as have specific cities and municipalities and governments. While any increase in the minimum wage must by definition help those working, care must be exercised in implementing particularly large increases. This is because large increases in particular areas or spheres may induce production units to move outside of the area covered, and thereby shift jobs to lower-wage areas.

12.4 Capital – concepts

The share of national income accruing to capital is more substantial than commonly recognized. National income in Canada is divided 60-40, favoring labour. This leaves a very large component going to the owners of capital. The stock of physical capital includes assembly-line machinery, rail lines, dwellings, consumer durables, school buildings and so forth. It is the stock of produced goods used as inputs to the production of other goods and services.

Physical capital is the stock of produced goods that are inputs in the production of other goods and services.

Physical capital is distinct from land in that the former is produced, whereas land is not. These in turn differ from *financial wealth*, which is not an input to production. We add to the capital stock by undertaking investment. But, because capital depreciates, investment in new capital goods is required merely to stand still. Depreciation accounts for the difference between gross and net investment.

Gross investment is the production of new capital goods and the improvement of existing capital goods.

Net investment is gross investment minus depreciation of the existing capital stock.

Depreciation is the annual change in the value of a physical asset.

Since capital is a stock of productive assets we must distinguish between the value of services that flow from capital and the value of capital assets themselves.

A **stock** is the quantity of an asset at a point in time.

A **flow** is the stream of services an asset provides during a period of time.

When a car is rented it provides the driver with a service; the car is the asset, or stock of capital, and the driving, or ability to move from place to place, is the service that flows from the use of the asset. When a photocopier is leased it provides a stream of services to the user. The copier is the asset; it represents a stock of physical capital. The printed products result from the service the copier provides per unit of time.

The price of an asset is what a purchaser pays for the asset. The owner then obtains the future stream of capital services it provides. Buying a car for \$30,000 entitles the owner to a stream of future transport services. The term rental rate defines the cost of the services from capital.

Capital services are the production inputs generated by capital assets.

The **rental rate** is the cost of using capital services.

The **price of an asset** is the financial sum for which the asset can be purchased.

But what determines the *price* of a productive asset? The price must reflect the value of future services that the capital provides. But we cannot simply add up these future values, because a dollar today is more valuable than a dollar several years from now. The key to valuing an asset lies in understanding how to compute the *present value* of a future income stream.

Present values and discounting

When capital is purchased it generates a stream of dollar values (returns) in the future. A critical question is: How is the price that should be paid for capital today related to the benefits that capital will bring in the future? Consider the simplest of examples: A business is contemplating buying a computer. This business has a two-year horizon. It believes that the purchase of the computer will yield a return of \$500 in the first year (today), \$500 in the second year (one period into the future), and have a scrap value of \$200. What is the maximum price the entrepreneur should pay for the computer? The answer is obtained by *discounting* the future returns to the present. Since a dollar today is worth more than a dollar tomorrow, we cannot simply add the dollar values from different time periods.

The value today of \$500 received a year from now is less than \$500, because if you had this amount today you could invest it at the going rate of interest and end up with more than \$500 tomorrow. For example, if the rate of interest is 10% ($= 0.1$), then \$500 today is worth \$550 next period. By the same reasoning, \$500 tomorrow is worth less than \$500 today. Formally, the value next period of any amount is that amount plus the interest earned; in this case the value next period of \$500 today is $\$500 \times (1 + r) = \$500 \times 1.1 = \$550$, where r is the interest rate. It follows that if we multiply a given sum by $(1+r)$ to obtain its value next period, then we must divide a sum received next period to obtain its value today. Hence the value today of \$500 next period is simply $\$500 / (1 + r) = \$500 / 1.1 = \$454.54$. To see that this must be true, note that if you have \$454.54 today you can invest it and obtain \$500 next period if the interest rate is 10%. In general:

Value next period	$= \text{value this period} \times (1 + \text{interest rate})$
Value this period	$= \frac{\text{Value next period}}{(1 + \text{interest rate})}$

This rule carries over to any number of future periods. The value of a sum of money today two periods into the future is obtained by multiplying the today value by $(1 + \text{interest rate})$ twice. Or the value of a sum of money today that will be received two periods from now is that sum divided by $(1 + \text{interest rate})$ twice. And so on, for any number of time periods. So if the amount is received twenty years into the future, its value today would be obtained by dividing that sum by $(1 + \text{interest rate})$ twenty times; if received ' n ' periods into the future it must be divided by $(1 + \text{interest rate})$ ' n ' times.

Two features of this discounting are to be noted: First, if the interest rate is high, the value today of future sums is smaller than if the interest rate is low. Second, sums received far in the future are worth much less than sums received in the near future.

Let us return to our initial example, assuming the interest rate is 0.1 (or 10%). The value of the year 1 return is \$500. The value of the year 2 return today is \$454.54, and the scrap value in today's terms is \$181.81. The value of all returns discounted to today is thus \$1,136.35.

Table 12.3 Present value of an asset ($i = 10\%$)

Year	Annual return	Scrap value	Discounted values
Year 1	500		500
Year 2	500	200	454.54 + 181.81
Asset value today			1,136.35

The **present value of a stream of future earnings** is the sum of each year's earnings divided by one plus the interest rate ' n ' times, where ' n ' is the number of years in the future when the amount will be received.

We are now in a position to determine how much the buyer should be willing to pay for the computer. Clearly if the value of the computer today, measured in terms of future returns to the entrepreneur's business, is \$1,136.35, then the potential buyer should be willing to pay any sum less than that amount. Paying more makes no economic sense.

Discounting is a technique used in countless applications. It underlies the prices we are willing to pay for corporate stocks: Analysts make estimates of future earnings of corporations; they then *discount those earnings back to the present*, and suggest that we not pay more for a unit of stock than indicated by the present value of future earnings.

12.5 The capital market

Demand

The analysis of the demand for the services of capital parallels closely that of labour demand: The rental rate for capital replaces the wage rate and capital services replace the hours of labour. It is important to keep in mind the distinction we drew above between capital services on the one hand and the amount of capital on the other. Capital services are produced by capital assets, just as work is produced by humans. Terms that are analogous to the marginal product of labour emerge naturally: The marginal product of capital (MP_K) is the output produced by one additional unit of capital services, with other inputs held constant. The value of this marginal product (VMP_K) is its value in the market place. It is the MP_K multiplied by the price of output.

The MP_K must eventually decline with a fixed amount of other factors of production. So, if the price of output is fixed for the firm, it follows that the VMP_K must also decline. We could pursue an analysis of the short-run demand for capital services, assuming labour was fixed, that would completely mirror the short-run demand for labour that we have already developed. But this would not add any new insights, so we move on to the supply side.

The **marginal product of capital** is the output produced by one additional unit of capital services, with all other inputs being held constant.

The **value of the marginal product of capital** is the marginal product of capital multiplied by the price of the output it produces.

Supply

We can grasp the key features of the market for capital by recognizing that the *flow* of capital services is determined by the capital *stock*: More capital means more services. The analysis of supply is complex because we must distinguish between the long run and the short run, and also between the supply to an industry and the supply in the whole economy.

In the *short run* the total supply of capital assets, and therefore services, is fixed to the *economy*, since new production capacity cannot come on stream overnight: The short-run supply of services is therefore vertical. In contrast, *a particular industry* in the short run faces a positively sloped supply: By offering a higher rental rate for trucks, one industry can bid them away from others.

The *long run* is a period of sufficient length to permit an addition to the capital stock. A supplier of capital, or capital services, must estimate the likely return he will get on the equipment he is contemplating having built. To illustrate: He is analyzing the purchase or construction of an earthmover that will cost \$100,000. Assuming that the annual maintenance and depreciation costs are \$10,000, and that the interest rate is 5% (implying that annual interest cost is \$5,000), it follows that the annual cost of owning such

a machine is \$15,000. If the entrepreneur is to undertake the investment she must therefore earn at least this amount annually (by renting it to others, or using it herself), and this is what is termed the required rental. We can think of it as the opportunity cost of ownership.

The **required rental** covers the sum of *maintenance, depreciation and interest costs*.

Prices and returns

In the long run, capital services in any sector of the economy must earn the required rental. If they earn more, entrepreneurs will be induced to build or purchase additional capital goods; if they earn less, owners of capital will allow machines to depreciate, or move the machines to other sectors of the economy.

As an example, the price of oil on world markets fell by half during 2015; from about \$100US per barrel to \$50US. At this price, many oil wells were no longer profitable, and oil drilling equipment was decommissioned. Technically, the value of the marginal product of capital declined, because the price of the good it was producing declined. In the near and medium term, no new investment in capital goods will take place in the oil drilling sector of the economy. If the price of oil should increase in the future, some of the decommissioned capital will be brought back into service. But some of this capital will deteriorate or depreciate and simply 'die', and be sold for scrap metal – particularly the older vintage capital. Only when the stock of oil drilling equipment is reduced by depreciation and decay to the required level will any new investment in this form of capital take place.

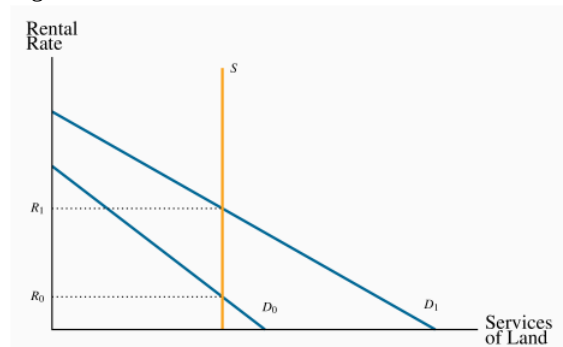
Note that the capital in this example is sector-specific. Drilling equipment cannot be easily redirected for use in other sectors. In contrast, earth movers can move from one sector of the economy to another with greater ease. An earth mover can be used to dig foundations for housing or commercial buildings; it can be used for strip mining; to build roads and bridges; to build tennis courts, golf courses and public parks. Such equipment may thus be moved to other sectors of the economy if in one particular sector the capital no longer can earn the required rental.

The prices of capital goods in the long run will be determined by the supply and demand for the services they provide. If the value of the services, as determined by supply and demand is high, then the price of assets will reflect this.

12.6 Land

Land is an input used in production, though is not a capital good in the way we defined capital goods earlier – production inputs that are themselves produced in the economy. Land is relatively fixed in supply *to the economy*, even in the long run. While this may not be literally true – the Netherlands reclaimed from the sea a great quantity of low-lying farmland, and fertilizers can turn marginal land into fertile land – it is a good approximation to reality. Figure 12.6 shows the derived demand D_0 for land services. With a fixed supply S , the equilibrium rental is R_0 .

Figure 12.6 The market for land services



The supply of land is relatively fixed, and therefore the return to land is primarily demand determined. Shifts in demand give rise to differences in returns.

In contrast to this economy-wide perspective, consider now a retailer who rents space in a commercial mall. The area around the mall experiences a surge in development and more people are shopping and doing business there. The retailer finds that she sells more, but also finds that her rent increases on account of the additional demand for space by commercial enterprises in the area. Her landlord is able to charge a higher rent because so many potential clients wish to rent space in the area. Consequently, despite the additional commerce in the area, the retailer's profit increase will be moderated by the higher rents she must pay: The demand for retail space is a *derived demand*. The situation can be explained with reference to Figure 12.6 again. On account of growth in this area, the demand for retail space shifts from D_0 to D_1 . Space in the area is restricted, and thus the vertical supply curve describes

the supply side well. So with little or no possibility of higher prices bringing forth additional supply, the additional demand makes for a steep price (rent) increase.

Land has many uses and the returns to land must reflect this. Land in downtown Vancouver is priced higher than land in rural Saskatchewan. Land cannot be moved from the latter to the former location however, and therefore the rent differences represent an equilibrium. In contrast, land in downtown Winnipeg that is used for a parking lot may not be able to compete with the use of that land for office development. Therefore, for it to remain as a parking lot, the rental must reflect its high opportunity cost. This explains why parking fees in big US cities such as Boston or New York may run to \$40 per day. If the parking owners could not obtain this fee, they could profitably sell the land to a developer. Ultimately it is the *value in its most productive use* that determines the price of land.

Key Terms

Demand for labour: a derived demand, reflecting the demand for the output of final goods and services.

Value of the marginal product is the marginal product multiplied by the price of the good produced.

Marginal revenue product of labour is the additional revenue generated by hiring one more unit of labour where the marginal revenue declines.

Monopsonist is the sole buyer of a good or service and faces an upward-sloping supply curve.

Participation rate: the fraction of the population in the working age group that joins the labour force.

The **labour force** is that part of the population either employed or seeking employment.

Unemployment rate: the fraction of the labour force actively seeking employment that is not employed.

Transfer earnings are the amount that an individual can earn in the next highest paying alternative job.

Rent is the excess remuneration an individual currently receives above the next best alternative. This alternative is the reservation wage.

Physical capital is the stock of produced goods that are inputs to the production of other goods and services.

Gross investment is the production of new capital goods and the improvement of existing capital goods.

Net investment is gross investment minus depreciation of the existing capital stock.

Depreciation is the annual change in the value of a physical asset.

Stock is the quantity of an asset at a point in time.

Flow is the stream of services an asset provides during a period of time.

Capital services are the production inputs generated by capital assets.

Rental rate: the cost of using capital services.

Asset price: the financial sum for which the asset can be purchased.

Present value of a stream of future earnings: the sum of each year's earnings divided by one plus the interest rate raised to the appropriate power.

Marginal product of capital is the output produced by one additional unit of capital services, with all other inputs being held constant.

Value of the marginal product of capital is the marginal product of capital multiplied by the price of the output it produces.

Required rental covers the sum of maintenance, depreciation and interest costs.

Exercises for Chapter 12

EXERCISE 12.1

Aerodynamics is a company specializing in the production of bicycle shirts. It has a fixed capital stock, and sells its shirts for \$20 each. It pays a weekly wage of \$400 per worker. Aerodynamics must maximize its profits by determining the optimal number of employees to hire. The marginal product of each worker can be inferred from the table below. Determine the optimal number of employees. [Hint: You must determine the VMP_L schedule, having first computed the MP_L .]

Employment	0	1	2	3	4	5	6
Total output	0	20	50	75	95	110	120
MP_L							
VMP_L							

EXERCISE 12.2

Suppose that, in Exercise 12.1 above, wages are not fixed. Instead the firm must pay \$50 more to employ each individual worker: The first worker is willing to work for \$250, the second for \$300, the third for \$350, etc. But once employed, each worker actually earns the same wage. Determine the optimal number of workers to be employed. [Hint: You must recognize that each worker earns the same wage; so when one additional worker is hired, the wage must increase to all workers employed.]

EXERCISE 12.3

Consider the following supply and demand equations for berry pickers. Demand: $W=22-0.4L$; supply: $W=10+0.2L$.

1. For values of $L = 1, 5, 10, 15, \dots, 30$, calculate the corresponding wage in each of the supply and demand functions.
2. Using the data from part (a), plot and identify the equilibrium wage and quantity of labour.
3. Illustrate in the diagram the areas defining transfer earnings and rent.
4. Compute the transfer earnings and rent components of the total wage bill.

EXERCISE 12.4

The rows of the following table describe the income stream for three different capital investments. The income flows accrue in years 1 and 2. Only year 2 returns need to be discounted. The rate of interest is the first entry in each row, and the project cost is the final entry.

Interest rate	Year 1	Year 2	Cost
8%	8,000	9,000	16,000
6%	0	1,000	900
10%	4,000	5,000	11,000

1. For each investment calculate the present value of the stream of services.
2. Decide whether or not the investment should be undertaken.

EXERCISE 12.5

Nihilist Nicotine is a small tobacco farm in south-western Ontario. It has three plots of land, each with a different productivity, in that the annual yield differs across plots. The output from each plot is given in the table below. Each plot is the same size and requires 3 workers and one machine to harvest the leaves. The cost of these inputs is \$10,000. If the price of each kilogram of leaves is \$4, how many plots should be planted?

Land plot	Leaf yield in kilograms
One	3,000
Two	2,500
Three	2,000

EXERCISE 12.6

The timing of wine sales is a frequent problem encountered by vintners. This is because many red wines improve with age. Let us suppose you own a particular vintage and you envisage that each bottle should increase in value by 10% the first year, 9% the second year, 8% the third year, etc.

1. Suppose the interest rate is 5%, for how many years would you hold the wine if there is no storage cost?
2. If in addition to interest rate costs, there is a cost of storing the wine that equals 2% of the wine's value each year, for how many years would you hold the wine before selling?

EXERCISE 12.7

Optional: The industry demand for plumbers is given by the equation $W=50-0.08L$, and there is a fixed supply of 300 qualified plumbers.

1. Draw a diagram illustrating the supply, demand and equilibrium, knowing that the quantity intercept for the demand equation is 625.
2. Solve the supply and demand equations for the equilibrium wage, W .
3. If the plumbers now form a union, and supply their labour at a wage of \$30 per hour, illustrate the new equilibrium on your diagram and calculate the new level of employment.

This page titled [12: Labour and capital](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis](#) and [Ian Irvine](#) (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.1: Labour - a derived demand

The value of labour springs from the value of its use, that is the value placed upon goods and services that it produces – product prices. The wage is the price that equilibrates the supply and demand for a given type of labour, and it reflects the value of that labour in production. Formally, the demand for labour (and capital) is thus a derived demand, in contrast to being a 'final' demand.

Demand for labour: a derived demand, reflecting the value of the output it produces.

We must distinguish between the long run and the short run in our analysis of factor markets. On the *supply side* certain factors of production are fixed in the short run. For example, the supply of radiologists can be increased only over a period of years. While one hospital may be able to attract radiologists from another hospital to meet a shortage, this does not increase the supply in the economy as a whole.

On the *demand side* there is the conventional difference between the short and long run: In the short run some of a firm's factors of production, such as capital, are fixed, and therefore the demand for labour differs from when all factors are variable – the long run.

Demand in the short run

Table 12.1 contains information from the example developed in Chapter 8. It can be used to illustrate how a firm reacts in the short run to a change in an input price, or to a change in the output price. The response of a producer to a change in the wage rate constitutes a demand function for labour – a schedule relating the quantity of the input demanded to different input prices. The output produced by the various numbers of workers yields a marginal product curve, whose values are stated in column 3. The marginal product of labour, MP_L , as developed in Chapter 8, is the additional output resulting from one more worker being employed, while holding constant the other (fixed) factors. But what is the *dollar value* to the firm of an additional worker? It is the additional *value of output* resulting from the additional employee – the price of the output times the worker's marginal contribution to output, his MP . We term this the value of the marginal product.

The **value of the marginal product** is the marginal product multiplied by the price of the good produced.

Table 12.1 Short-run production and labour demand

Workers	Output	MP_L	$VMP_L = MP_L \times P$	Marginal profit = (VMP_L –wage)
(1)	(2)	(3)	(4)	(5)
0	0			
1	15	15	1050	150
2	40	25	1750	750
3	70	30	2100	1100
4	110	40	2800	1800
5	145	35	2450	1450
6	175	30	2100	1100
7	200	25	1750	750
8	220	20	1400	400
9	235	15	1050	50
10	240	5	350	negative

Each unit of labour costs \$1,000; output sells at a fixed price of \$70 per unit.

In this example the MP_L first rises as more labour is employed, and then falls. With each unit of output selling for \$70 the value of the marginal product of labour (VMP_L) is given in column 4. The first worker produces 15 units each week, and since each unit sells for a price of \$70, his production value to the firm is \$1,050 ($= \70×15). A second worker produces 25 units, so his value to the firm is \$1,750, and so forth. If the weekly wage of each worker is \$1,000 then the firm can estimate its marginal profit from hiring each

additional worker. This is the difference between the value of the marginal product and the wage paid, and is given in the final column of the table.

It is profitable to hire more workers as long as the cost of an extra worker is less than the VMP_L . The equilibrium amount of labour to employ is therefore 9 units in this example. If the firm were to hire one more worker the contribution of that worker to its profit would be negative (\$350 – \$1,000), and if it hired one worker less it would forego the opportunity to make an additional profit of \$50 on the 9th unit (\$1,050 – \$1,000).

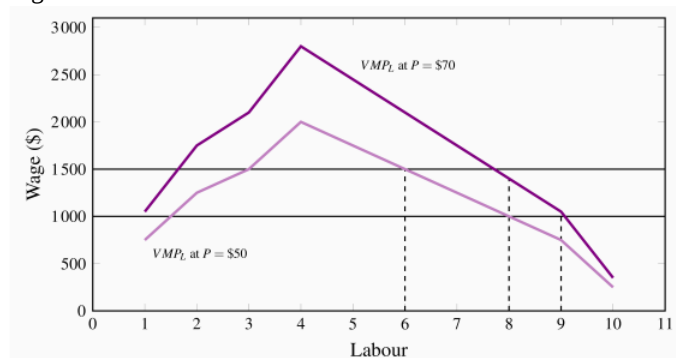
Profit maximizing hiring rule:

- If the VMP_L of next worker > wage, hire more labour.
- If the VMP_L < wage, hire less labour.

To this point we have determined the profit maximizing amount of labour to employ when the output price and the wage are given. However, a demand function for labour reflects the demand for labour at many different wage rates, just as the demand function for any product reflects the quantity demanded at various prices. Accordingly, suppose the wage rate is \$1,500 per week rather than \$1,000. The optimal amount of labour to employ in this case is determined in exactly the same manner: Employ the amount of labour where its contribution is marginally profitable. Clearly the optimal amount to employ is 7 units: The value of the seventh worker to the firm is \$1,750 and the value of the eighth worker is \$1,400. Hence it would not be profitable to employ the eighth, because his marginal contribution to profit would be negative. Following the same procedure we could determine the optimal amount of labour to employ at different wages. This implies that the VMP_L function is the demand for labour function because it determines the most profitable amount of labour to employ at any wage.

The optimal amount of labour to hire is illustrated in Figure 12.1. The wage and VMP_L curves come from Table 12.1. The VMP_L curve has an upward sloping segment, reflecting increasing productivity, and then a regular downward slope as developed in Chapter 8. At employment levels where the VMP_L is greater than the wage additional labour should be employed. But when the VMP_L falls below the wage rate employment should stop. If labour is divisible into very small units, the optimal employment decision is where the MP_L function intersects the wage line.

Figure 12.1 The demand for labour



The optimal hiring decision is defined by the condition that the value of the MP_L is greater than or equal to the wage paid. 9 workers are employed when the wage is \$1,000 and the price of output is \$70; 6 workers are employed when the wage is \$1,500 and the price of output is \$50.

Figure 12.1 also illustrates what happens to hiring when the output price changes. Consider a reduction in its price to \$50 from \$70. The profit impact of such a change is negative because the value of each worker's output has declined. Accordingly, the demand curve must reflect this by shifting inward (down), as in the figure. At various wage rates, less labour is now demanded. The new VMP_L schedule can be derived in Table 12.1 as before: It is the MP_L schedule multiplied by the lower value (\$50) of the final good.

In this example the firm is a perfect competitor in the output market, because the price of the good it produces is fixed. It can produce and sell more of the good without this having an impact on the price of the good in the marketplace. Where the firm is not a perfect competitor it faces a declining MR function. In this case the value of the MP_L is the product of MR and MP_P rather than P and MP_P . To distinguish the different output markets we use the term marginal revenue product of labour (MRP_L) when the demand for the output slopes downward. But the optimizing principle remains the same: The firm should calculate the value of each additional unit of labour, and hire up to the point where the additional revenue produced by the worker exceeds or equals the additional cost of that worker.

The **marginal revenue product of labour** is the additional revenue generated by hiring one more unit of labour where the marginal revenue declines.

Demand in the long run

In Chapter 8 we proposed that firms choose their factors of production in accordance with cost-minimizing principles. In producing a specific output, firms choose the least-cost combination of labour and plant size. But how is this choice affected when the price of labour or capital changes? While adjustment to price changes may require a long period of time, we know that if one factor becomes more (less) expensive, the firm will likely change the mix of capital and labour away from (towards) that factor. For example, when the accuracy and prices of production robots began to fall in the nineteen nineties, auto assemblers reduced their labour and used robots instead. When computers and computer software improved and declined in price, clerical workers were replaced by computers that were operated by accountants. But such adjustments and responses do not occur overnight.

In the short run a higher wage increases costs, but the firm is constrained in its choice of inputs by a fixed plant size. In the long run, a wage increase will induce the firm to use relatively more capital than when labour was less expensive in producing a given output. But despite the new choice of inputs, a rise in the cost of any input must increase the total cost of producing any output.

A change in the price of any factor has two impacts on firms: In the first place producers will *substitute* away from the factor whose price increases; second, there will be an impact on output and a change in the price of the final good it produces. Since the cost structure increases when the price of an input rises, the supply curve in the market for the good must reflect this – any given output will now be supplied at a higher price. With a downward sloping demand, this shift in supply must increase the price of the good and reduce the amount sold. This second effect can be called an *output effect*.

Monopsony

Some firms may have to pay a higher wage in order to employ more workers. Think of Hydro Quebec building a dam in Northern Quebec. Not every hydraulic engineer would be equally happy working there as in Montreal. Some engineers may demand only a small wage premium to work in the North, but others will demand a high premium. If so, Hydro Quebec must pay a higher wage to attract more workers – it faces an upward sloping supply of labour curve. Hydro Quebec is the sole buyer in this particular market and is called a monopsonist – a single buyer. Our general optimizing principle governing the employment of labour still holds, even if we have different names for the various functions: Hire any factor of production up to the point where the cost of an additional unit equals the value generated for the firm by that extra worker. The essential difference here is that when a firm faces an upward sloping labour supply it will have to pay more to attract additional workers and *also pay more to its existing workers*. This will impact the firm's willingness to hire additional workers.

A **monopsonist** is the sole buyer of a good or service and faces an upward-sloping supply curve.

Application Box 12.1 Monopsonies

Monopsonies are more than a curiosity; they exist in the real world. An excellent example is the cannabis market in Canada. Virtually every province has set up a trading agency that has the sole right to purchase cannabis from growers; growers and processors are not permitted to sell directly to retailers; they may only sell to the monopsony by law. In turn, these provincial cannabis monopsonies are frequently retail monopolists in that the agency owns all of the retail outlets in the province.

Firm versus industry demand

The demand for labour within an industry, or sector of the economy, is obtained from the sum of the demands by each individual firm. It is analogous to the goods market, but with a subtle difference. In Chapter 3 we obtained a market demand by summing individual demands horizontally. The same could be done here: At lower (or higher) wages, each firm will demand more (or less) labour. However, if all firms employ more labour in order to increase their output, the price of the output will likely decline. This in turn will moderate the demand for labour – it is slightly less valuable now that the price of the output it produces has fallen. This is a subtle point, and we can reasonably think of the demand for labour in a given sector of the economy as the sum of the demands on the part of the employers in that sector.

This page titled [12.1: Labour - a derived demand](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by Douglas Curtis and Ian Irvine (Lyryx) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.2: The supply of labour

Most prime-age individuals work, but some do not. The decision to join the labour force is called the participation decision. Of those who do participate in the labour force, some individuals work full time, others work part time, and yet others cannot find a job. The unemployment rate is the fraction of the labour force actively seeking employment that is not employed.

The **participation rate** for the economy is the fraction of the population in the working age group that joins the labour force.

The **labour force** is that part of the population either employed or seeking employment.

The **unemployment rate** is the fraction of the labour force actively seeking employment that is not employed.

Data on participation rates in Canada are given in Table 12.2 below for specific years in the modern era. The overall participation for men and women combined has increased since 1977 from 60.8% to 65.8% This aggregated rate camouflages different patterns for men and women. The rates for women have been rising while the rates for men have fallen. Women today are more highly educated, and their role in society and the economy is viewed very differently than in the earlier period. Female participation has increased both because of changing social norms, a rise in household productivity, the development of service industries designed to support home life, and the development of the institution of daycare for young children.

In contrast, male participation rates declined over the period, largely offsetting the increase in female participation. Fewer individuals in total are retiring before the age of 55 in the most recent decades. This reflects both the greater number of females in the market place, and perhaps also a recognition that many households have not saved enough to fund a retirement period that has become longer as a result of increased longevity.

Table 12.2 Labour force participation rate, Canada 1977-2015

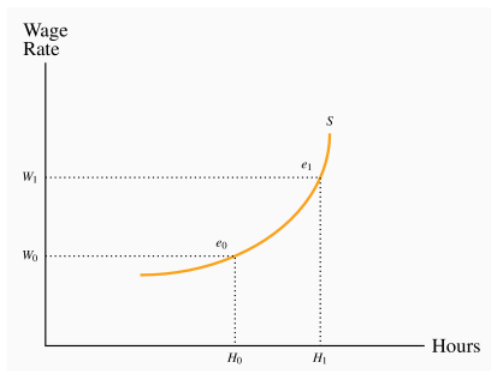
Year	Total	Men	Women	All > 55	Unemployment
1977	60.8	80.2	42.1	30.5	5.9
1990	66.6	77.1	56.8	25.9	7.4
1994	65.4	74.9	56.8	24.5	9.2
2001	66.1	73.4	59.2	26.4	6.2
2008	67.5	73.6	61.6	34.2	5.1
2015	66.2	72.0	60.6	37.3	6.0
2019	65.8	71.0	60.8	38.0	4.4

Source: Statistics Canada, CANSIM 14-10-0287-02
September of each year, for individuals aged ≥ 25 , unless stated.

At the micro level, the participation rate of individuals depends upon several factors. First, the wage rate that an individual can earn in the market is crucial. If that wage is low, then the individual may be more efficient in producing home services directly, rather than going into the labour market, earning a modest income and having to pay for home services. Second, there are fixed costs associated with working. A decision to work means that the individual must have work clothing, must undertake the costs of travel to work, and pay for daycare if there are children in the family. Third, the participation decision depends upon non-labour income. If the individual in question has a partner who earns a substantial amount, or if she has investment income, she will have less incentive to participate. Fourth, it depends inversely upon the tax rate.

The supply curve relates the supply decision to the price of labour – the wage rate. Economists who have studied the labour market tell us that the individual supply curve is upward sloping: As the wage increases, the individual wishes to supply more labour. From the point e_0 on the supply function in Figure 12.2, let the wage increase from w_0 to w_1 .

Figure 12.2 Individual labour supply



A wage increase from W_0 to W_1 induces the individual to *substitute* away from leisure, which is now more expensive, and work *more*. But the higher wage also means the individual can work fewer hours for a given standard of living; therefore the income effect induces *fewer hours*. On balance the substitution effect tends to dominate and the supply curve therefore slopes upward.

The individual offers more labour, H_1 , at the higher wage. What is the economic intuition behind the higher amount of labour supplied? Like much of choice theory there are two impacts associated with a higher price. First, the higher wage makes leisure more expensive relative to working. That midweek game of golf has become more expensive in terms of what the individual could earn. So the individual should *substitute* away from the more expensive 'good', leisure, towards labour. But at the same time, in order to generate a given income target the individual can work fewer hours at the higher wage. This is a type of *income effect*, indicating that income is greater at a higher wage regardless of the amount worked, and this induces the individual to work less. The fact that we draw the labour supply curve with a positive slope means that the substitution effect is the more important of the two. That is what statistical research has revealed.

Elasticity of the supply of labour

The value of the supply elasticity depends upon how the market in question is defined. In particular, it depends upon how large or small a given sector of the economy is, and whether we are considering the short run or the long run.

Suppose an industry is small relative to the whole economy and employs workers with common skills. These industries tend to pay the 'going wage'. For example, very many students are willing to work at the going rate for telemarketing firms, which compose a small sector of the economy. This means that the supply curve of such labour, *as far as that sector is concerned*, is in effect horizontal – infinitely elastic.

But some industries may not be small relative to the total labour supply. And in order to get more labour to work in such large sectors it may be necessary to provide the inducement of a higher wage: Additional workers may have to be attracted from another sector by means of higher wages. To illustrate: Consider the behaviour of two related sectors in housing – new construction and home restoration. In order to employ more plumbers and carpenters, new home builders may have to offer higher wages to induce them to move from the renovation sector. In this case the new housing industry's labour supply curve slopes upwards.

In the time dimension, a longer period is always associated with more flexibility. In this context, the supply of labour to any sector is more elastic, because it may take time for workers to move from one sector to another. Or, in cases where skills must be built up: When a sectoral expansion bids up the wages of information technology (IT) workers, more school leavers are likely to develop IT skills. Time will be required before additional graduates are produced, but in the long run, such additional supply will moderate the short-run wage increases.

Wages can be defined as being before-tax or after-tax. The after-tax, or take-home, wage is more important than the gross wage in determining the quantity of labour to be supplied. If taxes on additional hours of work are very high, workers are more likely to supply less hours than if tax rates are lower.

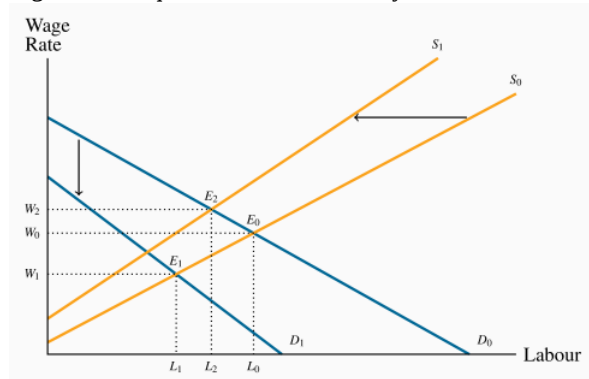
This page titled [12.2: The supply of labour](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.3: Labour market equilibrium and mobility

The fact that labour is a derived demand differentiates the labour market's equilibrium from the goods-market equilibrium. Let us investigate this with the help of Figure 12.3; it contains supply and demand functions for one particular industry – the cement industry, let us assume.

In Figure 12.1 we illustrated the impact on the demand for labour of a decline in the price of the output produced – a decline in the output price reduced the value of the marginal product of labour. In the current example, suppose that a slowdown in construction results in a decline in the price of cement. The impact of this price fall is to reduce the output value of each worker in the cement producing industry, because their output now yields a lower price. This decline in the VMP_L is represented in Figure 12.3 as a shift from D_0 to D_1 , which results in the new equilibrium E_1 .

Figure 12.3 Equilibrium in an industry labour market



A fall in the *price of the good* produced in a particular industry reduces the value of the MP_L . Demand for labour thus falls from D_0 to D_1 and a new equilibrium E_1 results. Alternatively, from E_0 , an increase in wages in another sector of the economy induces some labour to move to that sector. This is represented by the shift of S_0 to S_1 and the new equilibrium E_2 .

As a second example: Suppose that wages in some other sectors of the economy increase. The impact of this on the cement sector is that the supply of labour to the cement sector is reduced. In Chapter 3 we showed that a change in other prices may *shift* the demand or supply curve of interest. In Figure 12.3 supply shifts from S_0 to S_1 and the equilibrium goes from E_0 to E_2 .

How large are these impacts likely to be? That will depend upon how mobile labour is between sectors: Spillover effects will be smaller if labour is less mobile. This brings us naturally to the concepts of transfer earnings and rent.

Transfer earnings and rent

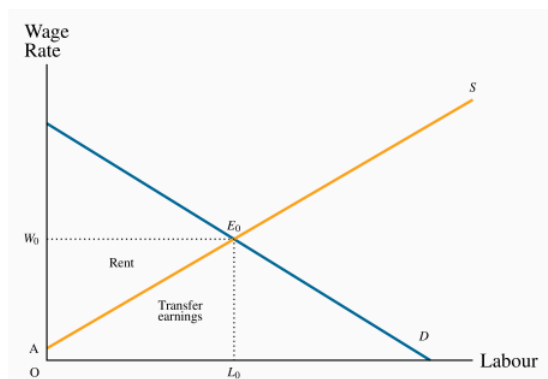
Consider the case of a performing violinist whose wage is \$80,000. If, as a best alternative, she can earn \$60,000 as a music teacher then her rent is \$20,000 and her transfer earnings \$60,000: Her rent is the excess she currently earns above the best alternative. Another violinist in the same orchestra, earning the same amount, who could earn \$55,000 as a teacher has rent of \$25,000. The alternative is also called the reservation wage. The violinists should not work in the orchestra unless they earn at least what they can earn in the next best alternative.

Transfer earnings are the amount that an individual can earn in the next highest paying alternative job.

Rent is the excess remuneration an individual currently receives above the next best alternative. This alternative is the **reservation wage**.

These concepts are illustrated in Figure 12.4. In this illustration, different individuals are willing to work for different amounts, but all are paid the same wage w_0 . The market labour supply curve by definition defines the wage for which each individual is willing to work. Thus the rent earned by labour in this market is the sum of the excess of the wage over each individual's transfer earnings – the area w_0E_0A . This area is also what we called producer or supplier surplus in Chapter 5.

Figure 12.4 Transfer earnings and rent



Rent is the excess of earnings over reservation wages. Each individual earns W_0 and is willing to work for the amount defined by the labour supply curve. Hence rent is W_0E_0A and transfer earnings OAE_0L_0 . Rent is thus the term for supplier surplus in this market.

Free labour markets?

Real-world labour markets are characterized by trade unions, minimum wage laws, benefit regulations, severance packages, parental leave, sick-day allowances and so forth. So can we really claim that markets work in the way we have described them – essentially as involving individual agents demanding and supplying labour? While labour markets are not completely 'free' in the conventional sense, the important issue is whether these interventions, that are largely designed to protect workers, have a large or small impact on the market. One reason why unemployment rates are generally higher in European economies than in Canada and the US is that labour markets are less subject to controls, and workers have a less supportive social safety net in North America.

Application Box 12.2 Are high salaries killing professional sports?

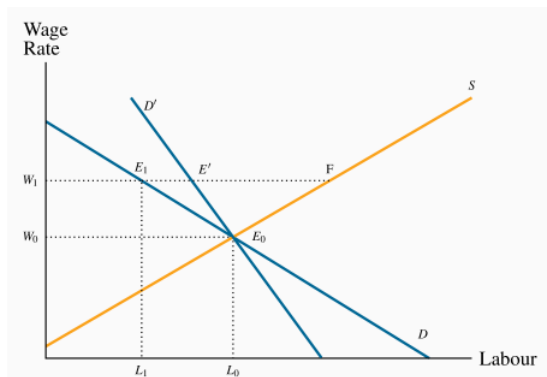
It is often said that the agents of professional players are killing their sport by demanding unreasonable salaries. On occasion, the major leagues are threatened with strikes, even though players are paid millions each year. In fact, wages are high because the derived demand is high. Fans are willing to pay high ticket prices, and television rights generate huge revenues. Combined, these revenues not only make ownership profitable, but increase the demand for the top players.

The lay person may be horrified at thirty-million dollar annual salaries. But in reality, many players receiving such salaries may be earning less than their marginal product! If Tom Brady did not play for the New England Patriots the team would have a lower winning record, attract fewer fans and make less profit. If Brady is paid \$25m per season, he is being paid less than his marginal product if the team were to lose \$40m in revenue as a result of his absence.

Given this, why do some teams incur financial losses? In fact very few teams make losses: Cries of poverty on the part of owners are more frequently part of the bargaining process, and revenue sharing means that very few teams do not make a profit.

The impact of 'frictions', such as unionization and minimum wages, in the labour market can be understood with the help of Figure 12.5. The initial 'free market' equilibrium is at E_0 , assuming that the workers are not unionized. In contrast, if the workers in this industry form a union, and negotiate a higher wage, for example W_1 rather than W_0 , then fewer workers will be employed. But how big will this reduction be? Clearly it depends on the elasticities of demand and supply. With the demand curve D , the excess supply at the wage W_1 is the difference E_1F . However, if the demand curve is less elastic, as illustrated by the curve D' , the excess supply is $E'F$. The excess supply also depends upon the supply elasticity. It is straightforward to see that a less elastic (more vertical) supply curve through E_0 would result in less excess supply.

Figure 12.5 Market interventions



E_0 is the equilibrium in the absence of a union. If the presence of a union forces the wage to W_1 fewer workers are employed. The magnitude of the decline from L_0 to L_1 depends on the elasticity of demand for labour. The excess supply at the wage W_1 is $(F-E_1)$. With a less elastic demand curve (D') the excess supply is reduced to $(F-E')$.

Beyond elasticity, the magnitude of the excess supply will also depend upon the degree to which the minimum wage, or the union-negotiated wage, lies above the equilibrium. That is, a larger value of the difference $(W_1 - W_0)$ results in more excess supply than a smaller difference.

While the above discussion pertains to unionization, it could equally well be interpreted in a minimum-wage context. If this figure describes the market for low-skill labour, and the government intervenes by setting a legal minimum at W_1 , then this will induce some degree of excess supply, depending upon the actual value of W_1 and the elasticities of supply and demand.

Despite the fact that a higher wage may induce some excess supply, it may increase total earnings. In Chapter 4 we saw that the dollar value of expenditure on a good increases when the price rises if the demand is inelastic. In the current example the 'good' is labour. Hence, a union-negotiated wage increase, or a higher minimum wage will each increase total remuneration if the demand for labour is inelastic. A case which has stirred great interest is described in Application Box 12.3.

Application Box 12.3 David Card on minimum wage

David Card is a famous Canadian-born labour economist who has worked at Princeton University and University of California, Berkeley. He is a winner of the prestigious Clark medal, an award made annually to an outstanding economist under the age of forty. Among his many contributions to the discipline, is a study of the impact of minimum wage laws on the employment of fast-food workers. With Alan Krueger as his co-researcher, Card examined the impact of the 1992 increase in the minimum wage in New Jersey and contrasted the impact on employment changes with neighbouring Pennsylvania, which did not experience an increase. They found virtually no difference in employment patterns between the two states. This research generated so much interest that it led to a special conference. Most economists now believe that modest changes in the level of the minimum wage have a small impact on employment levels.

Since about 2015, numerous labor-friendly movements favoring higher wages for low-paid workers have proposed a \$15 minimum in both Canada and the US. Some political parties have supported this movement, as have specific cities and municipalities and governments. While any increase in the minimum wage must by definition help those working, care must be exercised in implementing particularly large increases. This is because large increases in particular areas or spheres may induce production units to move outside of the area covered, and thereby shift jobs to lower-wage areas.

This page titled [12.3: Labour market equilibrium and mobility](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.4: Capital - concepts

The share of national income accruing to capital is more substantial than commonly recognized. National income in Canada is divided 60-40, favoring labour. This leaves a very large component going to the owners of capital. The stock of physical capital includes assembly-line machinery, rail lines, dwellings, consumer durables, school buildings and so forth. It is the stock of produced goods used as inputs to the production of other goods and services.

Physical capital is the stock of produced goods that are inputs in the production of other goods and services.

Physical capital is distinct from land in that the former is produced, whereas land is not. These in turn differ from *financial wealth*, which is not an input to production. We add to the capital stock by undertaking investment. But, because capital depreciates, investment in new capital goods is required merely to stand still. Depreciation accounts for the difference between gross and net investment.

Gross investment is the production of new capital goods and the improvement of existing capital goods.

Net investment is gross investment minus depreciation of the existing capital stock.

Depreciation is the annual change in the value of a physical asset.

Since capital is a stock of productive assets we must distinguish between the value of services that flow from capital and the value of capital assets themselves.

A **stock** is the quantity of an asset at a point in time.

A **flow** is the stream of services an asset provides during a period of time.

When a car is rented it provides the driver with a service; the car is the asset, or stock of capital, and the driving, or ability to move from place to place, is the service that flows from the use of the asset. When a photocopier is leased it provides a stream of services to the user. The copier is the asset; it represents a stock of physical capital. The printed products result from the service the copier provides per unit of time.

The price of an asset is what a purchaser pays for the asset. The owner then obtains the future stream of capital services it provides. Buying a car for \$30,000 entitles the owner to a stream of future transport services. The term rental rate defines the cost of the services from capital.

Capital services are the production inputs generated by capital assets.

The **rental rate** is the cost of using capital services.

The **price of an asset** is the financial sum for which the asset can be purchased.

But what determines the *price* of a productive asset? The price must reflect the value of future services that the capital provides. But we cannot simply add up these future values, because a dollar today is more valuable than a dollar several years from now. The key to valuing an asset lies in understanding how to compute the *present value* of a future income stream.

Present values and discounting

When capital is purchased it generates a stream of dollar values (returns) in the future. A critical question is: How is the price that should be paid for capital today related to the benefits that capital will bring in the future? Consider the simplest of examples: A business is contemplating buying a computer. This business has a two-year horizon. It believes that the purchase of the computer will yield a return of \$500 in the first year (today), \$500 in the second year (one period into the future), and have a scrap value of \$200. What is the maximum price the entrepreneur should pay for the computer? The answer is obtained by *discounting* the future returns to the present. Since a dollar today is worth more than a dollar tomorrow, we cannot simply add the dollar values from different time periods.

The value today of \$500 received a year from now is less than \$500, because if you had this amount today you could invest it at the going rate of interest and end up with more than \$500 tomorrow. For example, if the rate of interest is 10% ($= 0.1$), then \$500 today is worth \$550 next period. By the same reasoning, \$500 tomorrow is worth less than \$500 today. Formally, the value next period of any amount is that amount plus the interest earned; in this case the value next period of \$500 today is $500 \times (1 + r) = 500 \times 1.1 = 550$, where r is the interest rate. It follows that if we multiply a given sum by $(1+r)$ to obtain its value next period, then we must divide a sum received next period to obtain its value today. Hence the value today of \$500 next period is simply

$\$500/(1+r) = \$500/1.1 = \$454.54$. To see that this must be true, note that if you have \$454.54 today you can invest it and obtain \$500 next period if the interest rate is 10%. In general:

Value next period	$= \text{value this period} \times (1 + \text{interest rate})$
Value this period	$= \frac{\text{Value next period}}{(1 + \text{interest rate})}$

This rule carries over to any number of future periods. The value of a sum of money today two periods into the future is obtained by multiplying the today value by $(1 + \text{interest rate})$ twice. Or the value of a sum of money today that will be received two periods from now is that sum divided by $(1 + \text{interest rate})$ twice. And so on, for any number of time periods. So if the amount is received twenty years into the future, its value today would be obtained by dividing that sum by $(1 + \text{interest rate})$ twenty times; if received 'n' periods into the future it must be divided by $(1 + \text{interest rate})$ 'n' times.

Two features of this discounting are to be noted: First, if the interest rate is high, the value today of future sums is smaller than if the interest rate is low. Second, sums received far in the future are worth much less than sums received in the near future.

Let us return to our initial example, assuming the interest rate is 0.1 (or 10%). The value of the year 1 return is \$500. The value of the year 2 return today is \$454.54, and the scrap value in today's terms is \$181.81. The value of all returns discounted to today is thus \$1,136.35.

Table 12.3 Present value of an asset ($i = 10\%$)

Year	Annual return	Scrap value	Discounted values
Year 1	500		500
Year 2	500	200	454.54 + 181.81
Asset value today			1,136.35

The **present value of a stream of future earnings** is the sum of each year's earnings divided by one plus the interest rate 'n' times, where 'n' is the number of years in the future when the amount will be received.

We are now in a position to determine how much the buyer should be willing to pay for the computer. Clearly if the value of the computer today, measured in terms of future returns to the entrepreneur's business, is \$1,136.35, then the potential buyer should be willing to pay any sum less than that amount. Paying more makes no economic sense.

Discounting is a technique used in countless applications. It underlies the prices we are willing to pay for corporate stocks: Analysts make estimates of future earnings of corporations; they then *discount those earnings back to the present*, and suggest that we not pay more for a unit of stock than indicated by the present value of future earnings.

This page titled [12.4: Capital - concepts](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.5: The capital market

Demand

The analysis of the demand for the services of capital parallels closely that of labour demand: The rental rate for capital replaces the wage rate and capital services replace the hours of labour. It is important to keep in mind the distinction we drew above between capital services on the one hand and the amount of capital on the other. Capital services are produced by capital assets, just as work is produced by humans. Terms that are analogous to the marginal product of labour emerge naturally: The marginal product of capital (MP_K) is the output produced by one additional unit of capital services, with other inputs held constant. The value of this marginal product (VMP_K) is its value in the market place. It is the MP_K multiplied by the price of output.

The MP_K must eventually decline with a fixed amount of other factors of production. So, if the price of output is fixed for the firm, it follows that the VMP_K must also decline. We could pursue an analysis of the short-run demand for capital services, assuming labour was fixed, that would completely mirror the short-run demand for labour that we have already developed. But this would not add any new insights, so we move on to the supply side.

The **marginal product of capital** is the output produced by one additional unit of capital services, with all other inputs being held constant.

The **value of the marginal product of capital** is the marginal product of capital multiplied by the price of the output it produces.

Supply

We can grasp the key features of the market for capital by recognizing that the *flow* of capital services is determined by the capital *stock*: More capital means more services. The analysis of supply is complex because we must distinguish between the long run and the short run, and also between the supply to an industry and the supply in the whole economy.

In the *short run* the total supply of capital assets, and therefore services, is fixed to the *economy*, since new production capacity cannot come on stream overnight: The short-run supply of services is therefore vertical. In contrast, *a particular industry* in the short run faces a positively sloped supply: By offering a higher rental rate for trucks, one industry can bid them away from others.

The *long run* is a period of sufficient length to permit an addition to the capital stock. A supplier of capital, or capital services, must estimate the likely return he will get on the equipment he is contemplating having built. To illustrate: He is analyzing the purchase or construction of an earthmover that will cost \$100,000. Assuming that the annual maintenance and depreciation costs are \$10,000, and that the interest rate is 5% (implying that annual interest cost is \$5,000), it follows that the annual cost of owning such a machine is \$15,000. If the entrepreneur is to undertake the investment she must therefore earn at least this amount annually (by renting it to others, or using it herself), and this is what is termed the required rental. We can think of it as the opportunity cost of ownership.

The **required rental** covers the sum of *maintenance*, *depreciation* and *interest costs*.

Prices and returns

In the long run, capital services in any sector of the economy must earn the required rental. If they earn more, entrepreneurs will be induced to build or purchase additional capital goods; if they earn less, owners of capital will allow machines to depreciate, or move the machines to other sectors of the economy.

As an example, the price of oil on world markets fell by half during 2015; from about \$100US per barrel to \$50US. At this price, many oil wells were no longer profitable, and oil drilling equipment was decommissioned. Technically, the value of the marginal product of capital declined, because the price of the good it was producing declined. In the near and medium term, no new investment in capital goods will take place in the oil drilling sector of the economy. If the price of oil should increase in the future, some of the decommissioned capital will be brought back into service. But some of this capital will deteriorate or depreciate and simply 'die', and be sold for scrap metal – particularly the older vintage capital. Only when the stock of oil drilling equipment is reduced by depreciation and decay to the required level will any new investment in this form of capital take place.

Note that the capital in this example is sector-specific. Drilling equipment cannot be easily redirected for use in other sectors. In contrast, earth movers can move from one sector of the economy to another with greater ease. An earth mover can be used to dig foundations for housing or commercial buildings; it can be used for strip mining; to build roads and bridges; to build tennis courts,

golf courses and public parks. Such equipment may thus be moved to other sectors of the economy if in one particular sector the capital no longer can earn the required rental.

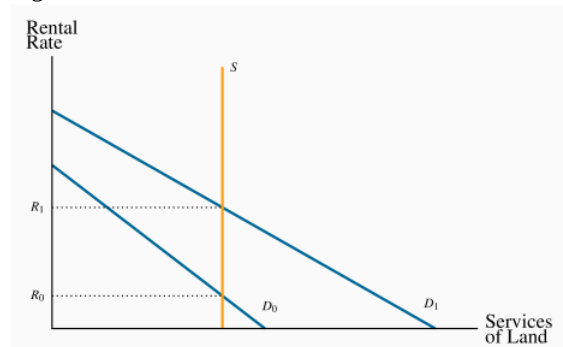
The prices of capital goods in the long run will be determined by the supply and demand for the services they provide. If the value of the services, as determined by supply and demand is high, then the price of assets will reflect this.

This page titled [12.5: The capital market](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.6: Land

Land is an input used in production, though is not a capital good in the way we defined capital goods earlier – production inputs that are themselves produced in the economy. Land is relatively fixed in supply *to the economy*, even in the long run. While this may not be literally true – the Netherlands reclaimed from the sea a great quantity of low-lying farmland, and fertilizers can turn marginal land into fertile land – it is a good approximation to reality. Figure 12.6 shows the derived demand D_0 for land services. With a fixed supply S , the equilibrium rental is R_0 .

Figure 12.6 The market for land services



The supply of land is relatively fixed, and therefore the return to land is primarily demand determined. Shifts in demand give rise to differences in returns.

In contrast to this economy-wide perspective, consider now a retailer who rents space in a commercial mall. The area around the mall experiences a surge in development and more people are shopping and doing business there. The retailer finds that she sells more, but also finds that her rent increases on account of the additional demand for space by commercial enterprises in the area. Her landlord is able to charge a higher rent because so many potential clients wish to rent space in the area. Consequently, despite the additional commerce in the area, the retailer's profit increase will be moderated by the higher rents she must pay: The demand for retail space is a *derived demand*. The situation can be explained with reference to Figure 12.6 again. On account of growth in this area, the demand for retail space shifts from D_0 to D_1 . Space in the area is restricted, and thus the vertical supply curve describes the supply side well. So with little or no possibility of higher prices bringing forth additional supply, the additional demand makes for a steep price (rent) increase.

Land has many uses and the returns to land must reflect this. Land in downtown Vancouver is priced higher than land in rural Saskatchewan. Land cannot be moved from the latter to the former location however, and therefore the rent differences represent an equilibrium. In contrast, land in downtown Winnipeg that is used for a parking lot may not be able to compete with the use of that land for office development. Therefore, for it to remain as a parking lot, the rental must reflect its high opportunity cost. This explains why parking fees in big US cities such as Boston or New York may run to \$40 per day. If the parking owners could not obtain this fee, they could profitably sell the land to a developer. Ultimately it is the *value in its most productive use* that determines the price of land.

This page titled [12.6: Land](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.7: Key Terms

Demand for labour: a derived demand, reflecting the demand for the output of final goods and services.

Value of the marginal product is the marginal product multiplied by the price of the good produced.

Marginal revenue product of labour is the additional revenue generated by hiring one more unit of labour where the marginal revenue declines.

Monopsonist is the sole buyer of a good or service and faces an upward-sloping supply curve.

Participation rate: the fraction of the population in the working age group that joins the labour force.

The **labour force** is that part of the population either employed or seeking employment.

Unemployment rate: the fraction of the labour force actively seeking employment that is not employed.

Transfer earnings are the amount that an individual can earn in the next highest paying alternative job.

Rent is the excess remuneration an individual currently receives above the next best alternative. This alternative is the reservation wage.

Physical capital is the stock of produced goods that are inputs to the production of other goods and services.

Gross investment is the production of new capital goods and the improvement of existing capital goods.

Net investment is gross investment minus depreciation of the existing capital stock.

Depreciation is the annual change in the value of a physical asset.

Stock is the quantity of an asset at a point in time.

Flow is the stream of services an asset provides during a period of time.

Capital services are the production inputs generated by capital assets.

Rental rate: the cost of using capital services.

Asset price: the financial sum for which the asset can be purchased.

Present value of a stream of future earnings: the sum of each year's earnings divided by one plus the interest rate raised to the appropriate power.

Marginal product of capital is the output produced by one additional unit of capital services, with all other inputs being held constant.

Value of the marginal product of capital is the marginal product of capital multiplied by the price of the output it produces.

Required rental covers the sum of maintenance, depreciation and interest costs.

This page titled [12.7: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

12.8: Exercises for Chapter 12

EXERCISE 12.1

Aerodynamics is a company specializing in the production of bicycle shirts. It has a fixed capital stock, and sells its shirts for \$20 each. It pays a weekly wage of \$400 per worker. Aerodynamics must maximize its profits by determining the optimal number of employees to hire. The marginal product of each worker can be inferred from the table below. Determine the optimal number of employees. [Hint: You must determine the VMP_L schedule, having first computed the MP_L .]

Employment	0	1	2	3	4	5	6
Total output	0	20	50	75	95	110	120
MP_L							
VMP_L							

EXERCISE 12.2

Suppose that, in Exercise 12.1 above, wages are not fixed. Instead the firm must pay \$50 more to employ each individual worker: The first worker is willing to work for \$250, the second for \$300, the third for \$350, etc. But once employed, each worker actually earns the same wage. Determine the optimal number of workers to be employed. [Hint: You must recognize that each worker earns the same wage; so when one additional worker is hired, the wage must increase to all workers employed.]

EXERCISE 12.3

Consider the following supply and demand equations for berry pickers. Demand: $W=22-0.4L$; supply: $W=10+0.2L$.

- For values of $L = 1, 5, 10, 15, \dots, 30$, calculate the corresponding wage in each of the supply and demand functions.
- Using the data from part (a), plot and identify the equilibrium wage and quantity of labour.
- Illustrate in the diagram the areas defining transfer earnings and rent.
- Compute the transfer earnings and rent components of the total wage bill.

EXERCISE 12.4

The rows of the following table describe the income stream for three different capital investments. The income flows accrue in years 1 and 2. Only year 2 returns need to be discounted. The rate of interest is the first entry in each row, and the project cost is the final entry.

Interest rate	Year 1	Year 2	Cost
8%	8,000	9,000	16,000
6%	0	1,000	900
10%	4,000	5,000	11,000

- For each investment calculate the present value of the stream of services.
- Decide whether or not the investment should be undertaken.

EXERCISE 12.5

Nihilist Nicotine is a small tobacco farm in south-western Ontario. It has three plots of land, each with a different productivity, in that the annual yield differs across plots. The output from each plot is given in the table below. Each plot is the same size and requires 3 workers and one machine to harvest the leaves. The cost of these inputs is \$10,000. If the price of each kilogram of leaves is \$4, how many plots should be planted?

Land plot	Leaf yield in kilograms
One	3,000
Two	2,500

Three

2,000

EXERCISE 12.6

The timing of wine sales is a frequent problem encountered by vintners. This is because many red wines improve with age. Let us suppose you own a particular vintage and you envisage that each bottle should increase in value by 10% the first year, 9% the second year, 8% the third year, etc.

- Suppose the interest rate is 5%, for how many years would you hold the wine if there is no storage cost?
- If in addition to interest rate costs, there is a cost of storing the wine that equals 2% of the wine's value each year, for how many years would you hold the wine before selling?

EXERCISE 12.7

Optional: The industry demand for plumbers is given by the equation $W=50-0.08L$, and there is a fixed supply of 300 qualified plumbers.

- Draw a diagram illustrating the supply, demand and equilibrium, knowing that the quantity intercept for the demand equation is 625.
- Solve the supply and demand equations for the equilibrium wage, W .
- If the plumbers now form a union, and supply their labour at a wage of \$30 per hour, illustrate the new equilibrium on your diagram and calculate the new level of employment.

This page titled [12.8: Exercises for Chapter 12](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13: Human capital and the income distribution

Chapter 13: Human capital and the income distribution

In this chapter we will explore:

13.1	The concept of human capital
13.2	Productivity and education
13.3	On-the-job training
13.4	Education as signaling
13.5	Returns to education and education quality
13.6	Discrimination
13.7	The income and earnings distribution in Canada
13.8	Wealth and capitalism

Individuals with different characteristics earn different amounts because their productivity levels differ. While it is convenient to work with a single marginal productivity of labour function to illustrate the functioning of the labour market, as we did in Chapter 12, for the most part wages and earnings vary by education level and experience, and sometimes by ethnicity and gender. In this chapter we develop an understanding of the sources of these differentials, and how they are reflected in the distribution of income.

13.1 Human capital

Human capital, HK, is the stock of knowledge and ability accumulated by a worker that determines future productivity and earnings. It depends on many different attributes – education, experience, intelligence, interpersonal skills etc. While human capital influences own earnings, it also impacts the productivity of the economy at large, and is therefore a vital force in determining long-run growth. Canada has been investing heavily in human capital in recent decades, and this suggests that future productivity and earnings will benefit accordingly.

Human capital is the stock of knowledge and ability accumulated by a worker that determines future productivity and earning.

Several features of Canada's recent human capital accumulation are noteworthy. First, Canada's enrollment rate in post-secondary education now exceeds the US rate, and that of virtually every economy in the world. Second is the fact that the number of women in third-level institutions exceeds the number of men. Almost 60% of university students are women. Third, international testing of high-school students sees Canadian students performing well, which indicates that the quality of the Canadian educational system appears to be high. These are positive aspects of a system that frequently comes under criticism. Nonetheless, the distribution of income that emerges from market forces in Canada has become more unequal.

Let us now try to understand individuals' *economic* motivation for embarking on the accumulation of human capital, and in the process see why different groups earn different amounts. We start by analyzing the role of education and then turn to on-the-job training. At the outset, we recognize that many individuals acquire knowledge for its own sake or in order to improve the quality of their lives. Education can improve one's appreciation of art, literature and the sciences.

13.2 Productivity and education

Human capital is the result of past investment that raises future incomes. A critical choice for individuals is to decide upon exactly how much additional human capital to accumulate. The cost of investing in another year of school is the *direct cost*, such as school fees, plus the *indirect*, or *opportunity*, cost, which can be measured by the foregone earnings during that extra year. The benefit of the additional investment is that the future flow of earnings is augmented. Consequently, wage differentials should reflect different degrees of education-dependent productivity.

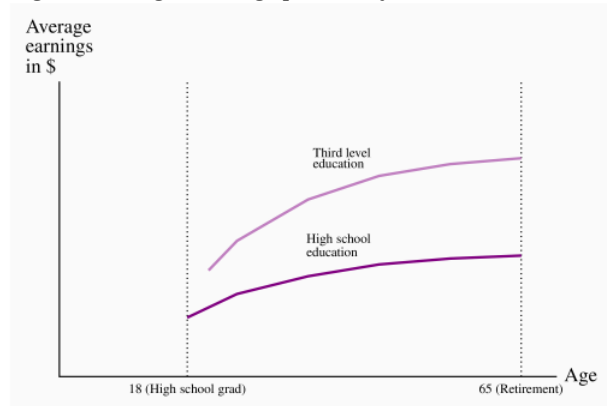
Age-earnings profiles

Figure 13.1 illustrates two typical age-earnings profiles for individuals with different levels of education. These profiles define the typical pattern of earnings over time, and are usually derived by examining averages across individuals in surveys. Two aspects are

clear: People with more education not only earn more, but the spread tends to grow with time. Less educated, healthy young individuals who work hard may earn a living wage but, unlike their more educated counterparts, they cannot look forward to a wage that rises substantially over time. More highly-educated individuals go into jobs and occupations that take a longer time to master: Lawyers, doctors and most professionals not only undertake more schooling than truck drivers, they also spend many years learning on the job, building up a clientele and accumulating expertise.

Age-earnings profiles define the pattern of earnings over time for individuals with different characteristics.

Figure 13.1 Age-Earnings profiles by education level



Individuals with a higher level of education earn more than individuals with a 'standard' level of education. In addition, the differential grows over time.

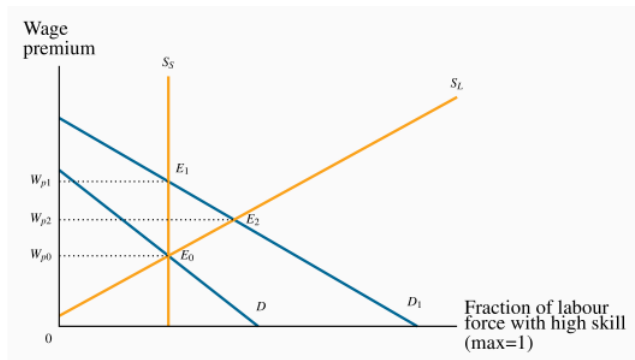
The education premium

Individuals with different education levels earn different wages. The education premium is the difference in earnings between the more and less highly educated. Quantitatively, Professors Kelly Foley and David Green have recently proposed that the completion of a college or trade certification adds about 15% to one's income, relative to an individual who has completed high school. A Bachelor's degree brings a premium of 20-25%, and a graduate degree several percentage points more¹. The failure to complete high school penalizes individuals to the extent of about 10%. These are average numbers, and they vary depending upon the province of residence, time period and gender. Nonetheless the findings underline that more human capital is associated with higher earnings. The earnings premium depends upon both the supply and demand of high HK individuals. *Ceteris paribus*, if high-skill workers are heavily in demand by employers, then the premium should be greater than if lower-skill workers are more in demand.

Education premium: the difference in earnings between the more and less highly educated.

The distribution of earnings has become more unequal in Canada and the US in recent decades, and one reason that has been proposed for this development is that the modern economy demands more high-skill workers; in particular that technological change has a bigger impact on productivity when combined with high-skill workers than with low-skill workers. Consider Figure 13.2 which contains supply and demand functions with a twist. We imagine that there are two types of labour: One with a high level of human capital, the other with a lower level. The vertical axis measures the wage *premium* of the high-education group (which can be measured in dollars or percentage terms), and the horizontal axis measures the *fraction of the total labour force that is of the high-skill type*. D is the *relative demand* for the high skill workers, in this example for the economy as a whole. There is some degree of substitution between high and low-skill workers in the modern economy. We do not propose that several low-skill workers can perform the work of one neuro-surgeon; but several individual households (low-skill) could complete their income tax submissions in the same time as one skilled tax specialist. In this example there is a degree of substitutability. In a production environment, a high-skill manager, equipped with technology and capital, can perform the tasks of several line workers.

Figure 13.2 The education/skill premium



A shift in demand increases the wage premium in the short run (from E_0 to E_1) by more than in the long run (to E_2). In the short run, the percentage of the labour force (S_S) that is highly skilled is fixed. In the long run it (S_L) is variable and responds to the wage premium.

The demand curve D defines the premium that employers (demanders) are willing to pay to the higher skill group. The negative slope indicates that if demanders were to employ a high proportion of skilled workers, the premium they would be willing to pay would be less than if they demanded a smaller share of high-skilled workers, and a larger share of lower-skilled workers. The wage premium for high HK individuals at any given time is determined by the intersection of supply and demand.

In the short run the make-up of the labour force is fixed, and this is reflected in the vertical supply curve S_S . The equilibrium is at E_0 , and W_{p0} is the premium, or excess, paid to the higher-skill worker over the lower-skill worker. In the long run it is possible for the economy to change the composition of its labour supply: If the wage premium increases, more individuals will find it profitable to train as high-skill workers. That is to say, the fraction of the total that is high-skill increases. It follows that the long-run supply curve slopes upwards.

So what happens when there is an increase in the demand for high-skill workers relative to low-skill workers? The demand curve shifts upward to D_1 , and the new equilibrium is at E_1 . The supply mix is fixed in the short run, so there is an increase in the wage premium. But over time, some individuals who might have been just indifferent between educating themselves more and going into the workplace with lower skill levels now find it worthwhile to pursue further education. Their higher anticipated returns to the additional human capital they invest in now exceed the additional costs of more schooling, whereas before the premium increase these additional costs and benefits were in balance. In Figure 13.2 the new short-run equilibrium at E_1 has a corresponding wage premium of W_{p1} . In the long run, after additional supply has reached the market, the increased premium is moderated to W_{p2} at the equilibrium E_2 .

This figure displays what many economists believe has happened in North America in recent decades: The demand for high HK individuals has increased, and the additional supply has not been as great. Consequently the wage premium for the high-skill workers has increased. As we describe later in this chapter, that is not the only perspective on what has happened.

Are students credit-constrained or culture-constrained?

The foregoing analysis assumes that students and potential students make rational decisions on the costs and benefits of further education and act accordingly. It also assumes implicitly that individuals can borrow the funds necessary to build their human capital: If the additional returns to further education are worthwhile, individuals should borrow the money to make the investment, just as entrepreneurs do with physical capital.

However, there is a key difference in the credit markets. If an entrepreneur fails in her business venture the lender will have a claim on the physical capital. But a bank cannot repossess a human being who drops out of school without having accumulated the intended human capital. Accordingly, the traditional lending institutions are frequently reluctant to lend the amount that students might like to borrow—students are credit constrained. The sons and daughters of affluent families therefore find it easier to attend university, because they are more likely to have a supply of funds domestically. Governments customarily step into the breach and supply loans and bursaries to students who have limited resources. While funding frequently presents an obstacle to attending a third-level institution, a stronger determinant of attendance is the education of the parents, as detailed in Application Box 13.1.

Application Box 13.1 Parental education and university attendance in Canada

The biggest single determinant of university attendance in the modern era is parental education. A recent study* of who goes to university examined the level of parental education of young people 'in transition' – at the end of their high school – for the years

1991 and 2000.

For the year 2000 they found that, if a parent had not completed high school, there was just a 12% chance that their son would attend university and an 18% chance that a daughter would attend. In contrast, for parents who themselves had completed a university degree, the probability that a son would also attend university was 53% and for a daughter 62%. Hence, the probability of a child attending university was roughly four times higher if the parent came from the top educational category rather than the bottom category! Furthermore the authors found that this probability gap opened wider between 1991 and 2000.

In the United States, Professor Sear Reardon of Stanford University has followed the performance of children from low-income households and compared their achievement with children from high-income households. He has found that the achievement gap between these groups of children has increased substantially over the last three decades. The reason for this growing separation is not because children from low-income households are performing worse in school, it is because high-income parents invest much more of their time and resources in educating their children, both formally in the school environment, and also in extra-school activities.

*Finnie, R., C. Laporte and E. Lascelles. "Family Background and Access to Post-Secondary Education: What Happened in the Nineties?" Statistics Canada Research Paper, Catalogue number 11F0019MIE-226, 2004

Reardon, Sean, "The Great Divide", New York Times, April 8, 2015.

13.3 On-the-job training

As is clear from Figure 13.1, earnings are raised both by education and experience. Learning on the job is central to the age-earnings profiles of the better educated, and is less important for those with lower levels of education. On-the-job training improves human capital through work experience. If on-the-job training increases worker productivity, who should pay for this learning – firms or workers? To understand who should pay, we distinguish between two kinds of skills: Firm-specific skills that raise a worker's productivity in a particular firm, and general skills that enhance productivity in many jobs or firms.

Firm-specific HK could involve knowing how particular components of a somewhat unique production structure functions, whereas general human capital might involve an understanding of engineering or architectural principles that can be applied universally. As for who should pay for the accumulation of skills: An employer should be willing to undertake most of the cost of firm-specific skills, because they are of less value to the worker should she go elsewhere. Firms offering general or transferable training try to pass the cost on to the workers, perhaps by offering a wage-earnings profile that starts very low, but that rises over time. Low-wage apprenticeships are examples. Hence, whether an employee is a medical doctor in residence, a plumber in an apprenticeship or a young lawyer in a law partnership, she 'pays' for the accumulation of her portable HK by facing a low wage when young. Workers are willing to accept such an earnings profile because their projected future earnings will compensate for lower initial earning.

On-the-job training improves human capital through work experience.

Firm-specific skills raise a worker's productivity in a particular firm.

General skills enhance productivity in many jobs or firms.

13.4 Education as signalling

An alternative view of education springs from the theory of signalling. This is a provocative theory that proposes education may be worthwhile, even if it generates little additional skills. The theory recognizes that individuals possess different abilities. However, firms cannot easily recognize the more productive workers without actually hiring them and finding out *ex-post* – sometimes a costly process. Signalling theory says that, in pursuing more education, people who know they are more capable thereby send a signal to potential employers that they are the more capable workers. Education therefore screens out the low-productivity workers from the high-productivity (more educated) workers. Firms pay more to the latter, because firms know that the high-ability workers are those with the additional education.

Signalling is the decision to undertake an action in order to reveal information.

Screening is the process of obtaining information by observing differences in behaviour.

To be effective, the process must separate the two types. Why don't lower-ability workers go to university and pretend they are of the high-ability type? Primarily because that strategy could backfire: Such individuals are less likely to succeed at school and there

are costs associated with school in the form of school fees, books and foregone earnings. While they may have lower innate skills, they are likely smart enough to recognize a bad bet.

On balance economists believe that further education does indeed add to productivity, although there may be an element of screening present: An engineering degree (we should hope) increases an individual's understanding of mechanical forces so that she can design a bridge that will not collapse, in addition to telling a potential employer that the student is smart!

Finally, it should be evident that if education raises productivity, it is also good for society and the economy at large.

13.5 Education returns and quality

How can we be sure that further education really does generate the returns, in the form of higher future incomes, to justify the investment? For many years econometricians proposed that an extra year of schooling might offer a return in the region of 10% – quite a favourable return in comparison with what is frequently earned on physical capital. Doubters then asked if the econometric estimation might be subject to bias – what if the additional earnings of those with more education are simply attributable to the fact that it is the innately more capable individuals who both earn more and who have more schooling? And since we cannot observe who has more innate ability, how can we be sure that it is the education itself, rather than just differences in ability, that generate the extra income?

This is a classical problem in inference: Does correlation imply causation? The short answer to this question is that education economists are convinced that the time invested in additional schooling does indeed produce additional rewards, even if it is equally true that individuals who are innately smarter do choose to invest in that way. Furthermore, it appears that the returns to graduate education are higher than the returns to undergraduate education.

What can be said of the quality of different educational systems? Are educational institutions in different countries equally good at producing knowledgeable students? Or, viewed another way: Has a grade nine student in Canada the same skill set as a grade nine student in France or Hong Kong? An answer to this question is presented in Table 13.1, which contains results from the Program for International Student Assessment (PISA) – an international survey of 15-year old student abilities in mathematics, science and literacy. This particular table presents the results for a sample of the countries that were surveyed. The results indicate that Canadian students perform well in all three dimensions of the test.

Table 13.1 Mean scores in PISA tests

Country	Math	Science	Reading
Australia	494	510	503
Austria	497	495	485
Belgium	507	502	499
Canada	516	528	527
Denmark	511	502	500
Finland	511	531	526
France	493	495	499
Germany	506	509	509
Greece	454	455	467
Hong Kong	548	523	527
Ireland	504	503	521
Italy	490	481	485
Japan	532	538	516
Korea	524	516	517
Mexico	408	416	423
New Zealand	495	513	509

Norway	502	498	513
Spain	486	493	496
Sweden	494	493	500
Switzerland	521	506	492
Turkey	420	425	428
United States	470	496	497
United Kingdom	492	509	498

Source: <https://www.oecd.org/pisa/pisa-2015-results-in-focus.pdf>

An interesting paradox arises at this point: If productivity growth in Canada has lagged behind some other economies in recent decades, as many economists believe, how can this be explained if Canada produces many well-educated high-skill workers? The answer may be that there is a considerable time lag before high participation rates in third-level education and high quality make themselves felt on the national stage in the form of elevated productivity. The evidence suggests a good productivity future, in so far as it depends upon human capital. At the same time, if investment in human capital is not matched by investment in physical capital then the human capital may not be able to perform to its ability.

13.6 Discrimination

Wage differences are a natural response to differences in human capital. But we frequently observe wage differences that might be discriminatory. For example, women on average earn less than men with similar qualifications; older workers may be paid less than those in their prime years; immigrants may be paid less than native-born Canadians, and ethnic minorities may be paid less than traditional white workers. The term discrimination describes an earnings differential that is attributable to a trait other than human capital.

If two individuals have the same HK, in the broadest sense of having the same capability to perform a particular task, then a wage premium paid to one represents discrimination. Correctly measured then, the discrimination premium between individuals from these various groups is the differential in earnings after correcting for HK differences. Thousands of studies have been undertaken on discrimination, and most conclude that discrimination abounds. Women, particularly those who have children, are paid less than men, and frequently face a 'glass ceiling' – a limit on their promotion possibilities within organizations.

Discrimination implies an earnings differential that is attributable to a trait other than human capital.

In contrast, women no longer face discrimination in university and college admissions, and form a much higher percentage of the student population than men in many of the higher paying professions such as medicine and law. Immigrants to Canada also suffer from a wage deficit. This is especially true for the most recent cohorts of working migrants who now come predominantly, not from Europe, as was once the case, but from China, South Asia, Africa and the Caribbean. For similarly-measured HK as Canadian-born individuals, these migrants frequently have an initial wage deficit of 30%, and require a period of more than twenty years to catch-up. How much of this differential might be due to the quality of the education or human capital received abroad is difficult to determine.

13.7 The income distribution

How does all of our preceding discussion play out when it comes to the income distribution? That is, when we examine the incomes of all individuals or households in the economy, how equally or unequally are they distributed?

The study of inequality is a critical part of economic analysis. It recognizes that income differences that are in some sense 'too large' are not good for society. Inordinately large differences can reflect poverty and foster social exclusion and crime. Economic growth that is concentrated in the hands of the few can increase social tensions, and these can have economic as well as social or psychological costs. Crime is one reflection of the divide between 'haves' and 'have-nots'. It is economically costly; but so too is child poverty. Impoverished children rarely achieve their social or economic potential and this is a loss both to the individual and society at large.

In this section we will first describe a subset of the basic statistical tools that economists use to measure inequality. Second, we will examine how income inequality has evolved in recent decades. We shall see that, while the picture is complex, market income

inequality has indeed increased in Canada. Third, we shall investigate some of the proposed reasons for the observed increase in inequality. Finally we will examine if the government offsets the inequality that arises from the marketplace through its taxation and redistribution policies.

It is to be emphasized that income inequality is just one proximate measure of the distribution of wellbeing. The extent of poverty is another such measure. Income is not synonymous with happiness but, that being said, income inequality can be computed reliably, and it provides a good measure of households' control over economic resources.

Theory and measurement

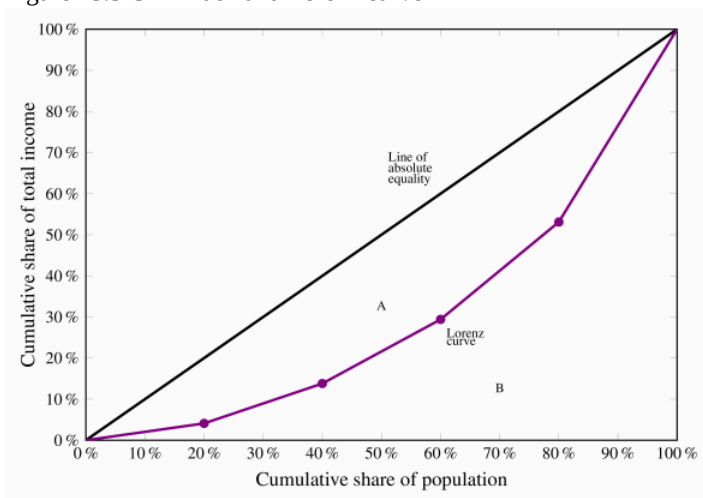
Let us rank the market incomes of all households in the economy from poor to rich, and categorize this ordering into different quantiles or groups. With five such quantiles the shares are called *quintiles*. The richest group forms the highest quintile, while the poorest group forms the lowest quintile. Such a representation is given in Table 13.2. The first numerical column displays the income in each quintile as a percentage of total income. If we wanted a finer breakdown, we could opt for decile (ten), or even vintile (twenty) shares, rather than quintile shares. These data can be graphed in a variety of ways. Since the data are in share, or percentage, form, we can compare, in a meaningful manner, distributions from economies that have different average income levels.

Table 13.2 Quintile shares of total family income in Canada, 2011

	Quintile share of total income	Cumulative share
First quintile	4.1	4.1
Second quintile	9.6	13.7
Third quintile	15.3	29.0
Fourth quintile	23.8	52.8
Fifth quintile	47.2	100.0
Total	100	

Source: Statistics Canada, CANSIM Matrix 2020405. These combinations are represented by the circles in the figure.

Figure 13.3 Gini index and Lorenz curve



The more equal are the income shares, the closer is the Lorenz curve to the diagonal line of equality. The Gini index is the ratio of the area A to the area (A+B). The Lorenz curve plots the cumulative percentage of total income against the cumulative percentage of the population.

An informative way of presenting these data graphically is to plot the cumulative share of income against the cumulative share of the population. This is given in the final column, and also presented graphically in Figure 13.3. The bottom quintile has 4.1% of total income. The bottom two quintiles together have 13.7% (4.1% + 9.6%), and so forth. By joining the coordinate pairs represented by the circles, a Lorenz curve is obtained. Relative to the diagonal line it is a measure of how unequally incomes are distributed: If

everyone had the same income, each 20% of the population would have 20% of total income and by joining the points for such a distribution we would get a straight diagonal line joining the corners of the box. In consequence, if the Lorenz curve is further from the line of equality the distribution is less equal than if the Lorenz curve is close to the line of equality.

Lorenz curve describes the cumulative percentage of the income distribution going to different quantiles of the population.

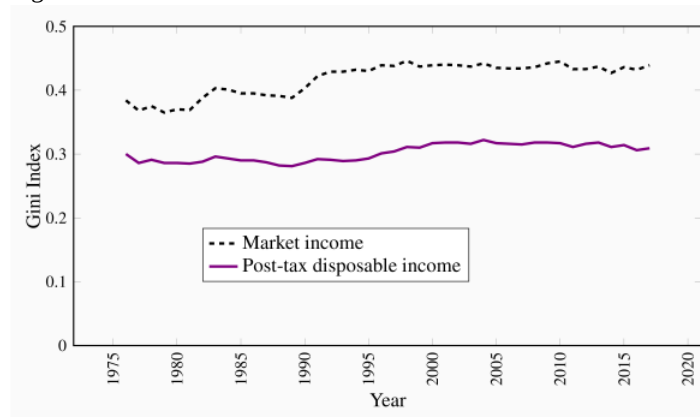
This suggests that the area A relative to the area (A + B) forms a measure of inequality in the income distribution. This fraction obviously lies between zero and one, and it is called the Gini index. A larger value of the Gini index indicates that inequality is greater. We will not delve into the mathematical formula underlying the Gini, but for this set of numbers its value is 0.4.

Gini index: a measure of how far the Lorenz curve lies from the line of equality. Its maximum value is one; its minimum value is zero.

The Gini index is what is termed *summary index* of inequality – it encompasses a lot of information in one number. There exist very many other such summary statistics.

It is important to recognize that very different Gini index values emerge for a given economy by using different income definitions of the variable going into the calculations. For example, the quintile shares of the *earnings of individuals* rather than the *incomes of households* could be very different. Similarly, the shares of income *post tax* and *post transfers* will differ from their shares on a *pre-tax, pre-transfer* basis.

Figure 13.4 Gini index Canada 1976–2015



Source: Statistics Canada, CANSIM Table 206-0033

Figure 13.4 contains Gini index values for two different definitions of income from 1976 to 2011. The upper line represents the Gini index values for households where the income measure is market income; the lower line defines the Gini values when income is defined as post-tax and post-transfer incomes. The latter income measure deducts taxes paid and adds income such as Employment Insurance or Social Assistance benefits. Two messages emerge from this graphic: The first is that the distribution of market incomes displays more inequality than the distribution of incomes after the government has intervened. In the latter case incomes are more equally distributed than in the former. The second message to emerge is that inequality increased over time – the Gini values are larger in the years after approximately 2000 than in the earlier years, although the increase in market income inequality is greater than the increase in income inequality based on a 'post-government' measure of income.

This is a very brief description of recent events. It is also possible to analyze inequality among women and men, for example, as well as among individuals and households. But the essential message remains clear: Definitions are important; in particular the distinction between incomes generated in the market place and incomes after the government has intervened through its tax and transfer policies.

Application Box 13.2 The very rich

McMaster University Professor Michael Veall and his colleague Emmanuel Saez, from University of California, Berkeley, have examined the evolution of the top end of the Canadian earnings distribution in the twentieth century. Using individual earnings from a database built upon tax returns, they show how the share of the very top of the distribution declined in the nineteen thirties and forties, remained fairly stable in the decades following World War II, and then increased from the eighties to the present time. The increase in share is particularly strong for the top 1% and even stronger for the top one tenth of the top 1%. These changes are driven primarily by changes in earnings, not on stock options awarded to high-level corporate employees. The authors conclude

that the change in this region of the distribution is attributable to changes in social norms. Whereas, in the nineteen eighties, it was expected that a top executive would earn perhaps a half million dollars, the 'norm' has become several million dollars in the present day. Such high remuneration became a focal point of public discussion after so many banks in the United States in 2008 and 2009 required government loans and support in order to avoid collapse. It also motivated the many 'occupy' movements of 2011 and 2012, and the US presidential race in 2019.

Saez, E. and M. Veall. "The evolution of high incomes in Canada, 1920-2000." Department of Economics research paper, McMaster University, March 2003.

In the international context, Canada is neither a strongly egalitarian economy nor one characterized by great income inequality. OECD data indicate that the economies with the lowest Gini index values are the Czech and Slovak republics and Iceland, with values in the neighborhood of 0.25 based on a post-government measure of income. Canada has a Gini index of .31, the US a value of .38 and at the upper end are economies such as Mexico and Chile with values of .47 (<https://data.oecd.org/inequality/income-inequality.htm>).

Economic forces

The increase in inequality of earnings in the market place in Canada has been reflected in many other developed economies – to a greater degree in the US and to a lesser extent in some European economies. Economists have devoted much energy to studying why, and as a result there are several accepted reasons.

Younger workers and those with lower skill levels have fared poorly in the last three decades. **Globalization and out-sourcing** have put pressure on low-end wages. In effect the workers in the lower tail of the distribution are increasingly competing with workers from low-wage less-developed economies. While this is a plausible causation, the critics of the perspective point out that wages at the bottom have fallen not only for those workers who compete with overseas workers in manufacturing, but also in the domestic services sector right across the economy. Obviously the workers at *McDonalds* have not the same competition from low-wage economies as workers who assemble toys.

A competing perspective is that it is **technological change** that has enabled some workers to do better than others. In explaining why high wage workers in many economies have seen their wages increase, whereas low-wage workers have seen a relative decline, the technological change hypothesis proposes that the form of recent technological change is critical: Change has been such as to require other complementary skills and education in order to benefit from it. For example, the introduction of computer-aided design technology is a benefit to workers who are already skilled and earning a high wage: *Existing high skills and technological change are complementary*. Such technological change is therefore different from the type underlying the production line. Automation in the early twentieth century in Henry Ford's plants improved the wages of lower skilled workers. But in the modern economy it is the highly skilled rather than the low skilled that benefit most from innovation.

A third perspective is that key **institutional changes** manifested themselves in the eighties and nineties, and these had independent impacts on the distribution. In particular, declines in the extent of unionization and changes in the minimum wage had significant impacts on earnings in the middle and bottom of the distribution: If unionization declines or the minimum wage fails to keep up with inflation, these workers will suffer. An alternative 'institutional' player is the government: In Canada the federal government became slightly less supportive, or 'generous', with its array of programs that form Canada's social safety net in the nineteen nineties. This tightening goes some way to explaining the modest inequality increase in the post-government income distribution in Figure 13.3 at this time. Nonetheless, most Canadian provincial governments increased the legal minimum wage in the first decade of the new millennium by substantially more than the rate of inflation. This meant that the economy's low-income workers did not fall further behind.

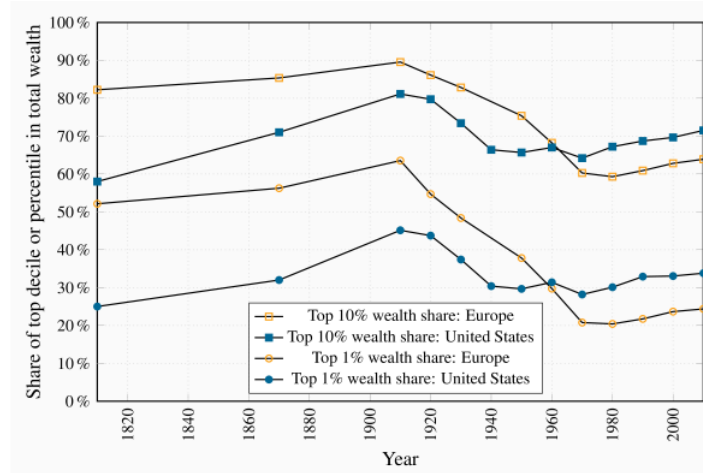
We conclude this overview of distributional issues by pointing out that we have not analyzed the distribution of wealth. Wealth too represents purchasing power, and it is wealth rather than income flows that primarily distinguishes Warren Buffet, Mark Zuckerberg and Bill Gates from the rest of us mortals. A detailed treatment of wealth inequality is beyond the scope of this book. We describe briefly, in the final section, the recent contribution of Thomas Piketty to the inequality debate.

13.8 Wealth and capitalism

In an insightful and popular study of capital accumulation, from both a historical and contemporary perspective, Thomas Piketty draws our attention to the enormous inequality in the distribution of wealth and explores what the future may hold in his book *Capital in the Twenty-First Century*.

The distribution of wealth is universally more unequal than the distribution of incomes or earnings. Gini coefficients in the neighbourhood of 0.8 are commonplace in developed economies. In terms of shares of the wealth pie, such a magnitude may imply that the top 1% of wealth holders own more than one third of all of an economy's wealth, that the top decile may own two-thirds of all wealth, and that the remaining one third is held by the 'bottom 90%'. And within this bottom 90%, virtually all of the remaining wealth is held by the 40% of the population below the top decile, leaving only a few percent of all wealth to the bottom 50% of the population.

Figure 13.5 Wealth inequality in Europe and the US, 1810-2010



Source: <http://piketty.pse.ens.fr/en/capital21c2>

While such an unequal holding pattern may appear shockingly unjust, Piketty informs us that current wealth inequality is not as great in most economies as it was about 1900. Figure 13.5 (Piketty 10.6) above is borrowed from his web site. Wealth was more unequally distributed in Old World Europe than New World America a century ago, but this relativity has since been reversed. A great transformation in the wealth holding pattern of societies in the twentieth century took the form of the emergence of a 'patrimonial middle class', by which he means the emergence of substantial wealth holdings on the part of that 40% of the population below the top decile. This development is noteworthy, but Piketty warns that the top percentiles of the wealth distribution may be on their way to controlling a share of all wealth close to their share in the early years of the twentieth century. He illustrates that two elements are critical to this prediction; first is the rate of growth in the economy relative to the return on capital, the second is inheritances.

To illustrate: Imagine an economy with very low growth, and where the owners of capital obtain an annual return of say 5%. If the owners merely maintain their capital intact and consume the remainder of this 5%, then the pattern of wealth holding will continue to be stable. However, if the holders of wealth can reinvest from their return an amount more than is necessary to replace depreciation then their wealth will grow. And if labour income in this economy is almost static on account of low overall growth, then wealth holders will secure a larger share of the economic pie. In contrast, if economic growth is significant, then labour income may grow in line with income from capital and inequality may remain stable. This summarizes Piketty's famous $(r-g)$ law – inequality depends upon the difference between the return on wealth and the growth rate of the economy. This potential for an ever-expanding degree of inequality is magnified when the stock of capital in the economy is large.

Consider now the role of inheritances. That is to say, do individuals leave large or small inheritances when they die, and how concentrated are such inheritances? If individual wealth accumulation patterns are generated by a desire to save for retirement and old age – during which time individuals decumulate by spending their assets – such motivation should result in small bequests being left to following generations. In contrast, if individuals who are in a position to do so save and accumulate, not just for their old age, but because they have dynastic preferences, or if they take pleasure simply from the ownership of wealth, or even if they are very cautious about running down their wealth in their old age, then we should see substantial inheritances passed on to the sons and daughters of these individuals, thereby perpetuating, and perhaps exacerbating, the inequality of wealth holding in the economy.

Piketty shows that in fact individuals who save substantial amounts tend to leave large bequests; that is they do not save purely for life-cycle motives. In modern economies the annual amount of bequests and gifts from parents to children falls in the range of 10%

to 15% of annual GDP. This may grow in future decades, and since wealth is highly concentrated, these bequests in turn are concentrated among a small number of the following generation – inequality is transmitted from one generation to the next.

As a final observation, if we consider the distribution of income and wealth together, particularly at the very top end, we can see readily that a growing concentration of *income* among the top 1% should ultimately translate itself into greater *wealth* inequality. This is because top earners can save more easily than lower earners. To compound matters, if individuals who inherit wealth also tend to inherit more human capital from their parents than others, the concentration of income and wealth may become yet stronger.

The study of distributional issues in economics has probably received too little attention in the modern era. Yet it is vitally important both in terms of the well-being of the individuals who constitute an economy and in terms of adherence to social norms. Given that utility declines with additions to income and wealth, transfers from those at the top to those at the bottom have the potential to increase total utility in the economy. Furthermore, an economy in which justice is seen to prevail—in the form of avoiding excessive inequality—is more likely to achieve a higher degree of social coherence than one where inequality is large.

Key Terms

Human capital is the stock of expertise accumulated by a worker that determines future productivity and earnings.

Age-earnings profiles define the pattern of earnings over time for individuals with different characteristics.

Education premium: the difference in earnings between the more and less highly educated.

On-the-job training improves human capital through work experience.

Firm-specific skills raise a worker's productivity in a particular firm.

General skills enhance productivity in many jobs or firms.

Signalling is the decision to undertake an action in order to reveal information.

Screening is the process of obtaining information by observing differences in behaviour.

Discrimination implies an earnings differential that is attributable to a trait other than human capital.

Lorenz curve describes the cumulative percentage of the income distribution going to different quantiles of the population.

Gini index: a measure of how far the Lorenz curve lies from the line of equality. Its maximum value is one; its minimum value is zero.

Exercises for Chapter 13

EXERCISE 13.1

Georgina is contemplating entering the job market after graduating from high school. Her future lifespan is divided into two phases: An initial one during which she may go to university, and a second when she will work. Since dollars today are worth more than dollars in the future she discounts the future by 20%, that is the value today of that future income is the income divided by 1.2. By going to university and then working she will earn (i) -\$60,000; (ii) \$600,000. The negative value implies that she will incur costs in educating herself in the first period. In contrast, if she decides to work for both periods she will earn \$30,000 in the first period and \$480,000 in the second.

1. If her objective is to maximize her lifetime earnings, should she go to university or enter the job market immediately?
2. If instead of discounting the future at the rate of 20%, she discounts it at the rate of 50%, what should she do?

EXERCISE 13.2

Imagine that you have the following data on the income distribution for two economies.

	Quintile share of total income	
First quintile	4.1	3.0
Second quintile	9.6	9.0
Third quintile	15.3	17.0
Fourth quintile	23.8	29.0
Fifth quintile	47.2	42.0

Total	100	100
--------------	-----	-----

1. On graph paper, or in a spreadsheet program, plot the Lorenz curves corresponding to the two sets of quintile shares. You must first compute the cumulative shares as we did for Figure 13.3.
2. Can you say, from a visual analysis, which distribution is more equal?

EXERCISE 13.3

The distribution of income in the economy is given in the table below. The first numerical column represents the dollars earned by each quintile. Since the numbers add to 100 you can equally think of the dollar values as shares of the total pie. In this economy the government changes the distribution by levying taxes and distributing benefits.

Quintile	Gross income \$m	Taxes \$m	Benefits \$m
First	4	0	9
Second	11	1	6
Third	19	3	5
Fourth	26	7	3
Fifth	40	15	3
Total	100	26	26

1. Plot the Lorenz curve for gross income to scale.
2. Now subtract the taxes paid and add the benefits received by each quintile. Check that the total income is still \$100. Calculate the cumulative income shares and plot the resulting Lorenz curve. Can you see that taxes and benefits reduce inequality?

EXERCISE 13.4

Consider two individuals, each facing a 45 year horizon at the age of 20. Ivan decides to work immediately and his earnings path takes the following form: Earnings = $20,000 + 1,000t - 10t^2$, where the t is time, and it takes on values from 1 to 25, reflecting the working lifespan.

1. In a spreadsheet enter values 1... 25 in the first column and then compute the value of earnings in each of the 25 years in the second column using the earnings equation.
2. John decides to study some more and only earns a part-time salary in his first few years. He hopes that the additional earnings in future years will compensate for that. His function is given by $10,000 + 2,000t - 12t^2$. In the same spreadsheet compute his annual earnings for 25 years.
3. Plot the two earnings functions you have computed using the 'charts' feature of Excel. Does your graph indicate that John passes Ivan between year 10 and year 11?

EXERCISE 13.5

In the short run one half of the labour force has high skills and one half low skills (in terms of Figure 13.2 this means that the short-run supply curve is vertical at 0.5). The relative demand for the high-skill workers is given by $W = 40 \times (1 - f)$, where W is the wage premium and f is the fraction that is skilled. The premium is measured in percent and f has a maximum value of 1. The W function thus has vertical and horizontal intercepts of $\{40, 1\}$.

1. Illustrate the supply and demand curves graphically, and illustrate the skill premium going to the high-skill workers in the short run by determining the value of W when $f=0.5$.
2. If demand increases to $W = 60 \times (1 - f)$ what is the new premium? Illustrate your answer graphically.

EXERCISE 13.6

Consider the foregoing problem in a long-run context, when the fraction of the labour force that is high-skilled is more elastic with respect to the premium. Let this long-run relative supply function be $W = 40 \times f$.

1. Graph this long-run supply function and verify that it goes through the same initial equilibrium as in Exercise 13.5.

2. Illustrate the long run and short run on the same diagram.
 3. What is the numerical value of the premium in the long run after the increase in demand? Illustrate graphically.
1. Foley, K. and D. Green, 2015, "Why more education will not solve rising inequality (and may make it worse)" , *Institute for Research in Public Policy*, Montreal, Canada.

This page titled [13: Human capital and the income distribution](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.1: Human capital

Human capital, HK, is the stock of knowledge and ability accumulated by a worker that determines future productivity and earnings. It depends on many different attributes – education, experience, intelligence, interpersonal skills etc. While human capital influences own earnings, it also impacts the productivity of the economy at large, and is therefore a vital force in determining long-run growth. Canada has been investing heavily in human capital in recent decades, and this suggests that future productivity and earnings will benefit accordingly.

Human capital is the stock of knowledge and ability accumulated by a worker that determines future productivity and earning.

Several features of Canada's recent human capital accumulation are noteworthy. First, Canada's enrollment rate in post-secondary education now exceeds the US rate, and that of virtually every economy in the world. Second is the fact that the number of women in third-level institutions exceeds the number of men. Almost 60% of university students are women. Third, international testing of high-school students sees Canadian students performing well, which indicates that the quality of the Canadian educational system appears to be high. These are positive aspects of a system that frequently comes under criticism. Nonetheless, the distribution of income that emerges from market forces in Canada has become more unequal.

Let us now try to understand individuals' *economic* motivation for embarking on the accumulation of human capital, and in the process see why different groups earn different amounts. We start by analyzing the role of education and then turn to on-the-job training. At the outset, we recognize that many individuals acquire knowledge for its own sake or in order to improve the quality of their lives. Education can improve one's appreciation of art, literature and the sciences.

This page titled [13.1: Human capital](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.2: Productivity and education

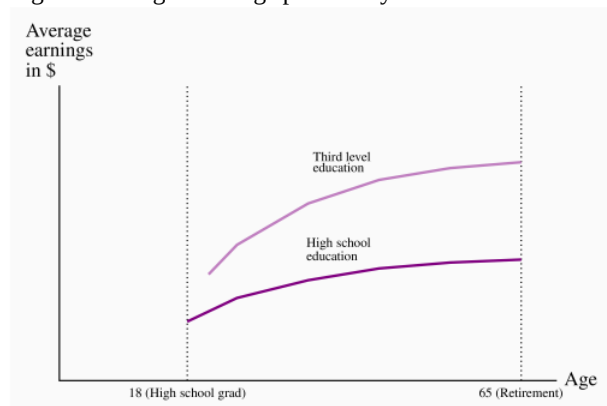
Human capital is the result of past investment that raises future incomes. A critical choice for individuals is to decide upon exactly how much additional human capital to accumulate. The cost of investing in another year of school is the *direct cost*, such as school fees, plus the *indirect*, or *opportunity cost*, which can be measured by the foregone earnings during that extra year. The benefit of the additional investment is that the future flow of earnings is augmented. Consequently, wage differentials should reflect different degrees of education-dependent productivity.

Age-earnings profiles

Figure 13.1 illustrates two typical age-earnings profiles for individuals with different levels of education. These profiles define the typical pattern of earnings over time, and are usually derived by examining averages across individuals in surveys. Two aspects are clear: People with more education not only earn more, but the spread tends to grow with time. Less educated, healthy young individuals who work hard may earn a living wage but, unlike their more educated counterparts, they cannot look forward to a wage that rises substantially over time. More highly-educated individuals go into jobs and occupations that take a longer time to master: Lawyers, doctors and most professionals not only undertake more schooling than truck drivers, they also spend many years learning on the job, building up a clientele and accumulating expertise.

Age-earnings profiles define the pattern of earnings over time for individuals with different characteristics.

Figure 13.1 Age-Earnings profiles by education level



Individuals with a higher level of education earn more than individuals with a 'standard' level of education. In addition, the differential grows over time.

The education premium

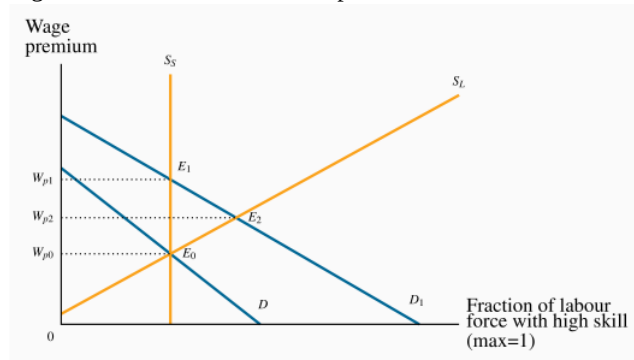
Individuals with different education levels earn different wages. The education premium is the difference in earnings between the more and less highly educated. Quantitatively, Professors Kelly Foley and David Green have recently proposed that the completion of a college or trade certification adds about 15% to one's income, relative to an individual who has completed high school. A Bachelor's degree brings a premium of 20-25%, and a graduate degree several percentage points more¹. The failure to complete high school penalizes individuals to the extent of about 10%. These are average numbers, and they vary depending upon the province of residence, time period and gender. Nonetheless the findings underline that more human capital is associated with higher earnings. The earnings premium depends upon both the supply and demand of high HK individuals. *Ceteris paribus*, if high-skill workers are heavily in demand by employers, then the premium should be greater than if lower-skill workers are more in demand.

Education premium: the difference in earnings between the more and less highly educated.

The distribution of earnings has become more unequal in Canada and the US in recent decades, and one reason that has been proposed for this development is that the modern economy demands more high-skill workers; in particular that technological change has a bigger impact on productivity when combined with high-skill workers than with low-skill workers. Consider Figure 13.2 which contains supply and demand functions with a twist. We imagine that there are two types of labour: One with a high level of human capital, the other with a lower level. The vertical axis measures the wage *premium* of the high-education group (which can be measured in dollars or percentage terms), and the horizontal axis measures the *fraction of the total labour force that is of the high-skill type*. D is the *relative demand* for the high skill workers, in this example for the economy as a whole. There is

some degree of substitution between high and low-skill workers in the modern economy. We do not propose that several low-skill workers can perform the work of one neuro-surgeon; but several individual households (low-skill) could complete their income tax submissions in the same time as one skilled tax specialist. In this example there is a degree of substitutability. In a production environment, a high-skill manager, equipped with technology and capital, can perform the tasks of several line workers.

Figure 13.2 The education/skill premium



A shift in demand increases the wage premium in the short run (from E_0 to E_1) by more than in the long run (to E_2). In the short run, the percentage of the labour force (S_S) that is highly skilled is fixed. In the long run it (S_L) is variable and responds to the wage premium.

The demand curve D defines the premium that employers (demanders) are willing to pay to the higher skill group. The negative slope indicates that if demanders were to employ a high proportion of skilled workers, the premium they would be willing to pay would be less than if they demanded a smaller share of high-skilled workers, and a larger share of lower-skilled workers. The wage premium for high HK individuals at any given time is determined by the intersection of supply and demand.

In the short run the make-up of the labour force is fixed, and this is reflected in the vertical supply curve S_S . The equilibrium is at E_0 , and W_{p0} is the premium, or excess, paid to the higher-skill worker over the lower-skill worker. In the long run it is possible for the economy to change the composition of its labour supply: If the wage premium increases, more individuals will find it profitable to train as high-skill workers. That is to say, the fraction of the total that is high-skill increases. It follows that the long-run supply curve slopes upwards.

So what happens when there is an increase in the demand for high-skill workers relative to low-skill workers? The demand curve shifts upward to D_1 , and the new equilibrium is at E_1 . The supply mix is fixed in the short run, so there is an increase in the wage premium. But over time, some individuals who might have been just indifferent between educating themselves more and going into the workplace with lower skill levels now find it worthwhile to pursue further education. Their higher anticipated returns to the additional human capital they invest in now exceed the additional costs of more schooling, whereas before the premium increase these additional costs and benefits were in balance. In Figure 13.2 the new short-run equilibrium at E_1 has a corresponding wage premium of W_{p1} . In the long run, after additional supply has reached the market, the increased premium is moderated to W_{p2} at the equilibrium E_2 .

This figure displays what many economists believe has happened in North America in recent decades: The demand for high HK individuals has increased, and the additional supply has not been as great. Consequently the wage premium for the high-skill workers has increased. As we describe later in this chapter, that is not the only perspective on what has happened.

Are students credit-constrained or culture-constrained?

The foregoing analysis assumes that students and potential students make rational decisions on the costs and benefits of further education and act accordingly. It also assumes implicitly that individuals can borrow the funds necessary to build their human capital: If the additional returns to further education are worthwhile, individuals should borrow the money to make the investment, just as entrepreneurs do with physical capital.

However, there is a key difference in the credit markets. If an entrepreneur fails in her business venture the lender will have a claim on the physical capital. But a bank cannot repossess a human being who drops out of school without having accumulated the intended human capital. Accordingly, the traditional lending institutions are frequently reluctant to lend the amount that students might like to borrow—students are credit constrained. The sons and daughters of affluent families therefore find it easier to attend university, because they are more likely to have a supply of funds domestically. Governments customarily step into the breach and

supply loans and bursaries to students who have limited resources. While funding frequently presents an obstacle to attending a third-level institution, a stronger determinant of attendance is the education of the parents, as detailed in Application Box 13.1.

Application Box 13.1 Parental education and university attendance in Canada

The biggest single determinant of university attendance in the modern era is parental education. A recent study* of who goes to university examined the level of parental education of young people 'in transition' – at the end of their high school – for the years 1991 and 2000.

For the year 2000 they found that, if a parent had not completed high school, there was just a 12% chance that their son would attend university and an 18% chance that a daughter would attend. In contrast, for parents who themselves had completed a university degree, the probability that a son would also attend university was 53% and for a daughter 62%. Hence, the probability of a child attending university was roughly four times higher if the parent came from the top educational category rather than the bottom category! Furthermore the authors found that this probability gap opened wider between 1991 and 2000.

In the United States, Professor Sear Reardon of Stanford University has followed the performance of children from low-income households and compared their achievement with children from high-income households. He has found that the achievement gap between these groups of children has increased substantially over the last three decades. The reason for this growing separation is not because children from low-income households are performing worse in school, it is because high-income parents invest much more of their time and resources in educating their children, both formally in the school environment, and also in extra-school activities.

*Finnie, R., C. Laporte and E. Lascelles. "Family Background and Access to Post-Secondary Education: What Happened in the Nineties?" Statistics Canada Research Paper, Catalogue number 11F0019MIE-226, 2004

Reardon, Sean, "The Great Divide", New York Times, April 8, 2015.

This page titled [13.2: Productivity and education](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.3: On-the-job training

As is clear from Figure 13.1, earnings are raised both by education and experience. Learning on the job is central to the age-earnings profiles of the better educated, and is less important for those with lower levels of education. On-the-job training improves human capital through work experience. If on-the-job training increases worker productivity, who should pay for this learning – firms or workers? To understand who should pay, we distinguish between two kinds of skills: Firm-specific skills that raise a worker's productivity in a particular firm, and general skills that enhance productivity in many jobs or firms.

Firm-specific HK could involve knowing how particular components of a somewhat unique production structure functions, whereas general human capital might involve an understanding of engineering or architectural principles that can be applied universally. As for who should pay for the accumulation of skills: An employer should be willing to undertake most of the cost of firm-specific skills, because they are of less value to the worker should she go elsewhere. Firms offering general or transferable training try to pass the cost on to the workers, perhaps by offering a wage-earnings profile that starts very low, but that rises over time. Low-wage apprenticeships are examples. Hence, whether an employee is a medical doctor in residence, a plumber in an apprenticeship or a young lawyer in a law partnership, she 'pays' for the accumulation of her portable HK by facing a low wage when young. Workers are willing to accept such an earnings profile because their projected future earnings will compensate for lower initial earning.

On-the-job training improves human capital through work experience.

Firm-specific skills raise a worker's productivity in a particular firm.

General skills enhance productivity in many jobs or firms.

This page titled [13.3: On-the-job training](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.4: Education as signalling

An alternative view of education springs from the theory of signalling. This is a provocative theory that proposes education may be worthwhile, even if it generates little additional skills. The theory recognizes that individuals possess different abilities. However, firms cannot easily recognize the more productive workers without actually hiring them and finding out *ex-post* – sometimes a costly process. Signalling theory says that, in pursuing more education, people who know they are more capable thereby send a signal to potential employers that they are the more capable workers. Education therefore screens out the low-productivity workers from the high-productivity (more educated) workers. Firms pay more to the latter, because firms know that the high-ability workers are those with the additional education.

Signalling is the decision to undertake an action in order to reveal information.

Screening is the process of obtaining information by observing differences in behaviour.

To be effective, the process must separate the two types. Why don't lower-ability workers go to university and pretend they are of the high-ability type? Primarily because that strategy could backfire: Such individuals are less likely to succeed at school and there are costs associated with school in the form of school fees, books and foregone earnings. While they may have lower innate skills, they are likely smart enough to recognize a bad bet.

On balance economists believe that further education does indeed add to productivity, although there may be an element of screening present: An engineering degree (we should hope) increases an individual's understanding of mechanical forces so that she can design a bridge that will not collapse, in addition to telling a potential employer that the student is smart!

Finally, it should be evident that if education raises productivity, it is also good for society and the economy at large.

This page titled [13.4: Education as signalling](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.5: Education returns and quality

How can we be sure that further education really does generate the returns, in the form of higher future incomes, to justify the investment? For many years econometricians proposed that an extra year of schooling might offer a return in the region of 10% – quite a favourable return in comparison with what is frequently earned on physical capital. Doubters then asked if the econometric estimation might be subject to bias – what if the additional earnings of those with more education are simply attributable to the fact that it is the innately more capable individuals who both earn more and who have more schooling? And since we cannot observe who has more innate ability, how can we be sure that it is the education itself, rather than just differences in ability, that generate the extra income?

This is a classical problem in inference: Does correlation imply causation? The short answer to this question is that education economists are convinced that the time invested in additional schooling does indeed produce additional rewards, even if it is equally true that individuals who are innately smarter do choose to invest in that way. Furthermore, it appears that the returns to graduate education are higher than the returns to undergraduate education.

What can be said of the quality of different educational systems? Are educational institutions in different countries equally good at producing knowledgeable students? Or, viewed another way: Has a grade nine student in Canada the same skill set as a grade nine student in France or Hong Kong? An answer to this question is presented in Table 13.1, which contains results from the Program for International Student Assessment (PISA) – an international survey of 15-year old student abilities in mathematics, science and literacy. This particular table presents the results for a sample of the countries that were surveyed. The results indicate that Canadian students perform well in all three dimensions of the test.

Table 13.1 Mean scores in PISA tests

Country	Math	Science	Reading
Australia	494	510	503
Austria	497	495	485
Belgium	507	502	499
Canada	516	528	527
Denmark	511	502	500
Finland	511	531	526
France	493	495	499
Germany	506	509	509
Greece	454	455	467
Hong Kong	548	523	527
Ireland	504	503	521
Italy	490	481	485
Japan	532	538	516
Korea	524	516	517
Mexico	408	416	423
New Zealand	495	513	509
Norway	502	498	513
Spain	486	493	496
Sweden	494	493	500
Switzerland	521	506	492

Turkey	420	425	428
United States	470	496	497
United Kingdom	492	509	498

Source: <https://www.oecd.org/pisa/pisa-2015-results-in-focus.pdf>

An interesting paradox arises at this point: If productivity growth in Canada has lagged behind some other economies in recent decades, as many economists believe, how can this be explained if Canada produces many well-educated high-skill workers? The answer may be that there is a considerable time lag before high participation rates in third-level education and high quality make themselves felt on the national stage in the form of elevated productivity. The evidence suggests a good productivity future, in so far as it depends upon human capital. At the same time, if investment in human capital is not matched by investment in physical capital then the human capital may not be able to perform to its ability.

This page titled [13.5: Education returns and quality](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.6: Discrimination

Wage differences are a natural response to differences in human capital. But we frequently observe wage differences that might be discriminatory. For example, women on average earn less than men with similar qualifications; older workers may be paid less than those in their prime years; immigrants may be paid less than native-born Canadians, and ethnic minorities may be paid less than traditional white workers. The term discrimination describes an earnings differential that is attributable to a trait other than human capital.

If two individuals have the same HK, in the broadest sense of having the same capability to perform a particular task, then a wage premium paid to one represents discrimination. Correctly measured then, the discrimination premium between individuals from these various groups is the differential in earnings after correcting for HK differences. Thousands of studies have been undertaken on discrimination, and most conclude that discrimination abounds. Women, particularly those who have children, are paid less than men, and frequently face a 'glass ceiling' – a limit on their promotion possibilities within organizations.

Discrimination implies an earnings differential that is attributable to a trait other than human capital.

In contrast, women no longer face discrimination in university and college admissions, and form a much higher percentage of the student population than men in many of the higher paying professions such as medicine and law. Immigrants to Canada also suffer from a wage deficit. This is especially true for the most recent cohorts of working migrants who now come predominantly, not from Europe, as was once the case, but from China, South Asia, Africa and the Caribbean. For similarly-measured HK as Canadian-born individuals, these migrants frequently have an initial wage deficit of 30%, and require a period of more than twenty years to catch-up. How much of this differential might be due to the quality of the education or human capital received abroad is difficult to determine.

This page titled [13.6: Discrimination](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis](#) and [Ian Irvine](#) ([Lyryx](#)) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.7: The income distribution

How does all of our preceding discussion play out when it comes to the income distribution? That is, when we examine the incomes of all individuals or households in the economy, how equally or unequally are they distributed?

The study of inequality is a critical part of economic analysis. It recognizes that income differences that are in some sense 'too large' are not good for society. Inordinately large differences can reflect poverty and foster social exclusion and crime. Economic growth that is concentrated in the hands of the few can increase social tensions, and these can have economic as well as social or psychological costs. Crime is one reflection of the divide between 'haves' and 'have-nots'. It is economically costly; but so too is child poverty. Impoverished children rarely achieve their social or economic potential and this is a loss both to the individual and society at large.

In this section we will first describe a subset of the basic statistical tools that economists use to measure inequality. Second, we will examine how income inequality has evolved in recent decades. We shall see that, while the picture is complex, market income inequality has indeed increased in Canada. Third, we shall investigate some of the proposed reasons for the observed increase in inequality. Finally we will examine if the government offsets the inequality that arises from the marketplace through its taxation and redistribution policies.

It is to be emphasized that income inequality is just one proximate measure of the distribution of wellbeing. The extent of poverty is another such measure. Income is not synonymous with happiness but, that being said, income inequality can be computed reliably, and it provides a good measure of households' control over economic resources.

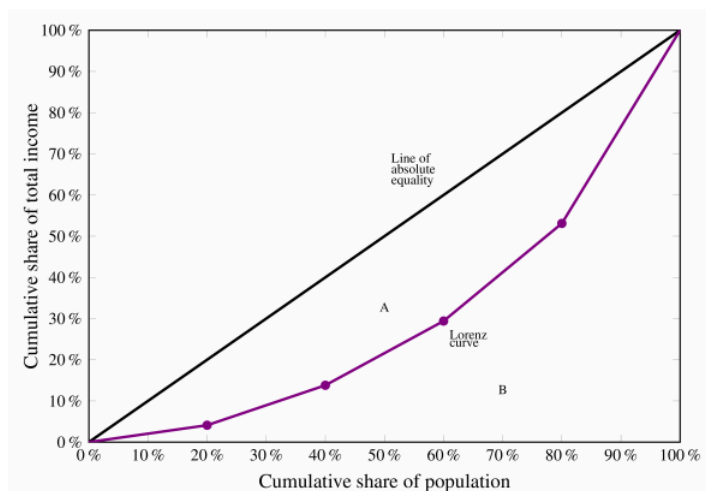
Theory and measurement

Let us rank the market incomes of all households in the economy from poor to rich, and categorize this ordering into different quantiles or groups. With five such quantiles the shares are called *quintiles*. The richest group forms the highest quintile, while the poorest group forms the lowest quintile. Such a representation is given in Table 13.2. The first numerical column displays the income in each quintile as a percentage of total income. If we wanted a finer breakdown, we could opt for decile (ten), or even vintile (twenty) shares, rather than quintile shares. These data can be graphed in a variety of ways. Since the data are in share, or percentage, form, we can compare, in a meaningful manner, distributions from economies that have different average income levels.

Table 13.2 Quintile shares of total family income in Canada, 2011

	Quintile share of total income	Cumulative share
First quintile	4.1	4.1
Second quintile	9.6	13.7
Third quintile	15.3	29.0
Fourth quintile	23.8	52.8
Fifth quintile	47.2	100.0
Total	100	

Source: Statistics Canada, CANSIM Matrix 2020405. These combinations are represented by the circles in the figure.
Figure 13.3 Gini index and Lorenz curve



The more equal are the income shares, the closer is the Lorenz curve to the diagonal line of equality. The Gini index is the ratio of the area A to the area (A+B). The Lorenz curve plots the cumulative percentage of total income against the cumulative percentage of the population.

An informative way of presenting these data graphically is to plot the cumulative share of income against the cumulative share of the population. This is given in the final column, and also presented graphically in Figure 13.3. The bottom quintile has 4.1% of total income. The bottom two quintiles together have 13.7% (4.1% + 9.6%), and so forth. By joining the coordinate pairs represented by the circles, a Lorenz curve is obtained. Relative to the diagonal line it is a measure of how unequally incomes are distributed: If everyone had the same income, each 20% of the population would have 20% of total income and by joining the points for such a distribution we would get a straight diagonal line joining the corners of the box. In consequence, if the Lorenz curve is further from the line of equality the distribution is less equal than if the Lorenz curve is close to the line of equality.

Lorenz curve describes the cumulative percentage of the income distribution going to different quantiles of the population.

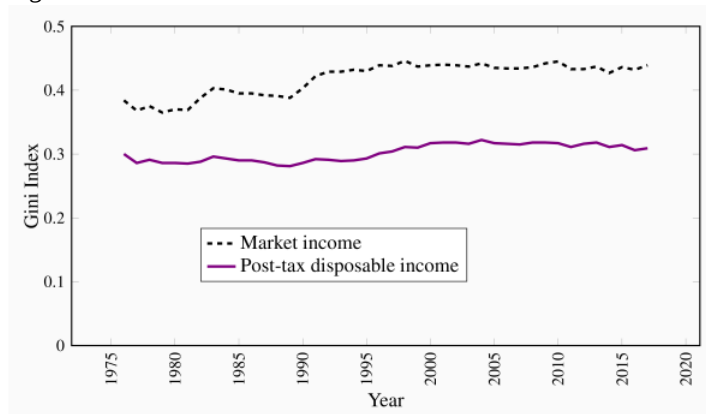
This suggests that the area A relative to the area (A + B) forms a measure of inequality in the income distribution. This fraction obviously lies between zero and one, and it is called the Gini index. A larger value of the Gini index indicates that inequality is greater. We will not delve into the mathematical formula underlying the Gini, but for this set of numbers its value is 0.4.

Gini index: a measure of how far the Lorenz curve lies from the line of equality. Its maximum value is one; its minimum value is zero.

The Gini index is what is termed *summary index* of inequality – it encompasses a lot of information in one number. There exist very many other such summary statistics.

It is important to recognize that very different Gini index values emerge for a given economy by using different income definitions of the variable going into the calculations. For example, the quintile shares of the *earnings of individuals* rather than the *incomes of households* could be very different. Similarly, the shares of income *post tax* and *post transfers* will differ from their shares on a *pre-tax, pre-transfer* basis.

Figure 13.4 Gini index Canada 1976–2015



Source: Statistics Canada, CANSIM Table 206-0033

Figure 13.4 contains Gini index values for two different definitions of income from 1976 to 2011. The upper line represents the Gini index values for households where the income measure is market income; the lower line defines the Gini values when income is defined as post-tax and post-transfer incomes. The latter income measure deducts taxes paid and adds income such as Employment Insurance or Social Assistance benefits. Two messages emerge from this graphic: The first is that the distribution of market incomes displays more inequality than the distribution of incomes after the government has intervened. In the latter case incomes are more equally distributed than in the former. The second message to emerge is that inequality increased over time – the Gini values are larger in the years after approximately 2000 than in the earlier years, although the increase in market income inequality is greater than the increase in income inequality based on a 'post-government' measure of income.

This is a very brief description of recent events. It is also possible to analyze inequality among women and men, for example, as well as among individuals and households. But the essential message remains clear: Definitions are important; in particular the distinction between incomes generated in the market place and incomes after the government has intervened through its tax and transfer policies.

Application Box 13.2 The very rich

McMaster University Professor Michael Veall and his colleague Emmanuel Saez, from University of California, Berkeley, have examined the evolution of the top end of the Canadian earnings distribution in the twentieth century. Using individual earnings from a database built upon tax returns, they show how the share of the very top of the distribution declined in the nineteen thirties and forties, remained fairly stable in the decades following World War II, and then increased from the eighties to the present time. The increase in share is particularly strong for the top 1% and even stronger for the top one tenth of the top 1%. These changes are driven primarily by changes in earnings, not on stock options awarded to high-level corporate employees. The authors conclude that the change in this region of the distribution is attributable to changes in social norms. Whereas, in the nineteen eighties, it was expected that a top executive would earn perhaps a half million dollars, the 'norm' has become several million dollars in the present day. Such high remuneration became a focal point of public discussion after so many banks in the United States in 2008 and 2009 required government loans and support in order to avoid collapse. It also motivated the many 'occupy' movements of 2011 and 2012, and the US presidential race in 2019.

Saez, E. and M. Veall. "The evolution of high incomes in Canada, 1920-2000." Department of Economics research paper, McMaster University, March 2003.

In the international context, Canada is neither a strongly egalitarian economy nor one characterized by great income inequality. OECD data indicate that the economies with the lowest Gini index values are the Czech and Slovak republics and Iceland, with values in the neighborhood of 0.25 based on a post-government measure of income. Canada has a Gini index of .31, the US a value of .38 and at the upper end are economies such as Mexico and Chile with values of .47 (<https://data.oecd.org/inequality/income-inequality.htm>).

Economic forces

The increase in inequality of earnings in the market place in Canada has been reflected in many other developed economies – to a greater degree in the US and to a lesser extent in some European economies. Economists have devoted much energy to studying why, and as a result there are several accepted reasons.

Younger workers and those with lower skill levels have fared poorly in the last three decades. **Globalization and out-sourcing** have put pressure on low-end wages. In effect the workers in the lower tail of the distribution are increasingly competing with workers from low-wage less-developed economies. While this is a plausible causation, the critics of the perspective point out that wages at the bottom have fallen not only for those workers who compete with overseas workers in manufacturing, but also in the domestic services sector right across the economy. Obviously the workers at *McDonalds* have not the same competition from low-wage economies as workers who assemble toys.

A competing perspective is that it is **technological change** that has enabled some workers to do better than others. In explaining why high wage workers in many economies have seen their wages increase, whereas low-wage workers have seen a relative decline, the technological change hypothesis proposes that the form of recent technological change is critical: Change has been such as to require other complementary skills and education in order to benefit from it. For example, the introduction of computer-aided design technology is a benefit to workers who are already skilled and earning a high wage: *Existing high skills and technological change are complementary*. Such technological change is therefore different from the type underlying the production

line. Automation in the early twentieth century in Henry Ford's plants improved the wages of lower skilled workers. But in the modern economy it is the highly skilled rather than the low skilled that benefit most from innovation.

A third perspective is that key **institutional changes** manifested themselves in the eighties and nineties, and these had independent impacts on the distribution. In particular, declines in the extent of unionization and changes in the minimum wage had significant impacts on earnings in the middle and bottom of the distribution: If unionization declines or the minimum wage fails to keep up with inflation, these workers will suffer. An alternative 'institutional' player is the government: In Canada the federal government became slightly less supportive, or 'generous', with its array of programs that form Canada's social safety net in the nineteen nineties. This tightening goes some way to explaining the modest inequality increase in the post-government income distribution in Figure 13.3 at this time. Nonetheless, most Canadian provincial governments increased the legal minimum wage in the first decade of the new millennium by substantially more than the rate of inflation. This meant that the economy's low-income workers did not fall further behind.

We conclude this overview of distributional issues by pointing out that we have not analyzed the distribution of wealth. Wealth too represents purchasing power, and it is wealth rather than income flows that primarily distinguishes Warren Buffet, Mark Zuckerberg and Bill Gates from the rest of us mortals. A detailed treatment of wealth inequality is beyond the scope of this book. We describe briefly, in the final section, the recent contribution of Thomas Piketty to the inequality debate.

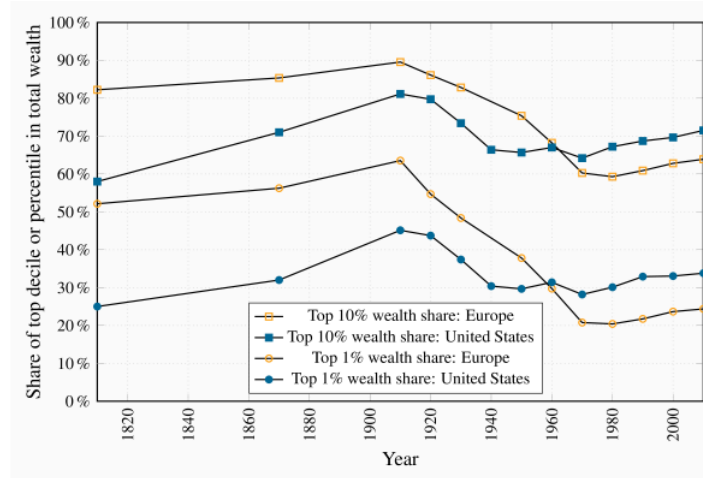
This page titled [13.7: The income distribution](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.8: Wealth and capitalism

In an insightful and popular study of capital accumulation, from both a historical and contemporary perspective, Thomas Piketty draws our attention to the enormous inequality in the distribution of wealth and explores what the future may hold in his book *Capital in the Twenty-First Century*.

The distribution of wealth is universally more unequal than the distribution of incomes or earnings. Gini coefficients in the neighbourhood of 0.8 are commonplace in developed economies. In terms of shares of the wealth pie, such a magnitude may imply that the top 1% of wealth holders own more than one third of all of an economy's wealth, that the top decile may own two-thirds of all wealth, and that the remaining one third is held by the 'bottom 90%'. And within this bottom 90%, virtually all of the remaining wealth is held by the 40% of the population below the top decile, leaving only a few percent of all wealth to the bottom 50% of the population.

Figure 13.5 Wealth inequality in Europe and the US, 1810-2010



Source: <http://piketty.pse.ens.fr/en/capital21c2>

While such an unequal holding pattern may appear shockingly unjust, Piketty informs us that current wealth inequality is not as great in most economies as it was about 1900. Figure 13.5 (Piketty 10.6) above is borrowed from his web site. Wealth was more unequally distributed in Old World Europe than New World America a century ago, but this relativity has since been reversed. A great transformation in the wealth holding pattern of societies in the twentieth century took the form of the emergence of a 'patrimonial middle class', by which he means the emergence of substantial wealth holdings on the part of that 40% of the population below the top decile. This development is noteworthy, but Piketty warns that the top percentiles of the wealth distribution may be on their way to controlling a share of all wealth close to their share in the early years of the twentieth century. He illustrates that two elements are critical to this prediction; first is the rate of growth in the economy relative to the return on capital, the second is inheritances.

To illustrate: Imagine an economy with very low growth, and where the owners of capital obtain an annual return of say 5%. If the owners merely maintain their capital intact and consume the remainder of this 5%, then the pattern of wealth holding will continue to be stable. However, if the holders of wealth can reinvest from their return an amount more than is necessary to replace depreciation then their wealth will grow. And if labour income in this economy is almost static on account of low overall growth, then wealth holders will secure a larger share of the economic pie. In contrast, if economic growth is significant, then labour income may grow in line with income from capital and inequality may remain stable. This summarizes Piketty's famous $(r-g)$ law – inequality depends upon the difference between the return on wealth and the growth rate of the economy. This potential for an ever-expanding degree of inequality is magnified when the stock of capital in the economy is large.

Consider now the role of inheritances. That is to say, do individuals leave large or small inheritances when they die, and how concentrated are such inheritances? If individual wealth accumulation patterns are generated by a desire to save for retirement and old age – during which time individuals decumulate by spending their assets – such motivation should result in small bequests being left to following generations. In contrast, if individuals who are in a position to do so save and accumulate, not just for their old age, but because they have dynastic preferences, or if they take pleasure simply from the ownership of wealth, or even if they are very cautious about running down their wealth in their old age, then we should see substantial inheritances passed on to the

sons and daughters of these individuals, thereby perpetuating, and perhaps exacerbating, the inequality of wealth holding in the economy.

Piketty shows that in fact individuals who save substantial amounts tend to leave large bequests; that is they do not save purely for life-cycle motives. In modern economies the annual amount of bequests and gifts from parents to children falls in the range of 10% to 15% of annual GDP. This may grow in future decades, and since wealth is highly concentrated, these bequests in turn are concentrated among a small number of the following generation – inequality is transmitted from one generation to the next.

As a final observation, if we consider the distribution of income and wealth together, particularly at the very top end, we can see readily that a growing concentration of *income* among the top 1% should ultimately translate itself into greater *wealth* inequality. This is because top earners can save more easily than lower earners. To compound matters, if individuals who inherit wealth also tend to inherit more human capital from their parents than others, the concentration of income and wealth may become yet stronger.

The study of distributional issues in economics has probably received too little attention in the modern era. Yet it is vitally important both in terms of the well-being of the individuals who constitute an economy and in terms of adherence to social norms. Given that utility declines with additions to income and wealth, transfers from those at the top to those at the bottom have the potential to increase total utility in the economy. Furthermore, an economy in which justice is seen to prevail—in the form of avoiding excessive inequality—is more likely to achieve a higher degree of social coherence than one where inequality is large.

This page titled [13.8: Wealth and capitalism](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.9: Key Terms

Human capital is the stock of expertise accumulated by a worker that determines future productivity and earnings.

Age-earnings profiles define the pattern of earnings over time for individuals with different characteristics.

Education premium: the difference in earnings between the more and less highly educated.

On-the-job training improves human capital through work experience.

Firm-specific skills raise a worker's productivity in a particular firm.

General skills enhance productivity in many jobs or firms.

Signalling is the decision to undertake an action in order to reveal information.

Screening is the process of obtaining information by observing differences in behaviour.

Discrimination implies an earnings differential that is attributable to a trait other than human capital.

Lorenz curve describes the cumulative percentage of the income distribution going to different quantiles of the population.

Gini index: a measure of how far the Lorenz curve lies from the line of equality. Its maximum value is one; its minimum value is zero.

This page titled [13.9: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

13.10: Exercises for Chapter 13

EXERCISE 13.1

Georgina is contemplating entering the job market after graduating from high school. Her future lifespan is divided into two phases: An initial one during which she may go to university, and a second when she will work. Since dollars today are worth more than dollars in the future she discounts the future by 20%, that is the value today of that future income is the income divided by 1.2. By going to university and then working she will earn (i) -\$60,000; (ii) \$600,000. The negative value implies that she will incur costs in educating herself in the first period. In contrast, if she decides to work for both periods she will earn \$30,000 in the first period and \$480,000 in the second.

- If her objective is to maximize her lifetime earnings, should she go to university or enter the job market immediately?
- If instead of discounting the future at the rate of 20%, she discounts it at the rate of 50%, what should she do?

EXERCISE 13.2

Imagine that you have the following data on the income distribution for two economies.

	Quintile share of total income	
First quintile	4.1	3.0
Second quintile	9.6	9.0
Third quintile	15.3	17.0
Fourth quintile	23.8	29.0
Fifth quintile	47.2	42.0
Total	100	100

- On graph paper, or in a spreadsheet program, plot the Lorenz curves corresponding to the two sets of quintile shares. You must first compute the cumulative shares as we did for Figure 13.3.
- Can you say, from a visual analysis, which distribution is more equal?

EXERCISE 13.3

The distribution of income in the economy is given in the table below. The first numerical column represents the dollars earned by each quintile. Since the numbers add to 100 you can equally think of the dollar values as shares of the total pie. In this economy the government changes the distribution by levying taxes and distributing benefits.

Quintile	Gross income \$m	Taxes \$m	Benefits \$m
First	4	0	9
Second	11	1	6
Third	19	3	5
Fourth	26	7	3
Fifth	40	15	3
Total	100	26	26

- Plot the Lorenz curve for gross income to scale.
- Now subtract the taxes paid and add the benefits received by each quintile. Check that the total income is still \$100. Calculate the cumulative income shares and plot the resulting Lorenz curve. Can you see that taxes and benefits reduce inequality?

EXERCISE 13.4

Consider two individuals, each facing a 45 year horizon at the age of 20. Ivan decides to work immediately and his earnings path takes the following form: $\text{Earnings} = 20,000 + 1,000t - 10t^2$, where the t is time, and it takes on values from 1 to 25, reflecting the working lifespan.

- In a spreadsheet enter values 1... 25 in the first column and then compute the value of earnings in each of the 25 years in the second column using the earnings equation.
- John decides to study some more and only earns a part-time salary in his first few years. He hopes that the additional earnings in future years will compensate for that. His function is given by $10,000 + 2,000t - 12t^2$. In the same spreadsheet compute his annual earnings for 25 years.
- Plot the two earnings functions you have computed using the 'charts' feature of Excel. Does your graph indicate that John passes Ivan between year 10 and year 11?

EXERCISE 13.5

In the short run one half of the labour force has high skills and one half low skills (in terms of Figure 13.2 this means that the short-run supply curve is vertical at 0.5). The relative demand for the high-skill workers is given by $W = 40 \times (1 - f)$, where W is the wage premium and f is the fraction that is skilled. The premium is measured in percent and f has a maximum value of 1. The W function thus has vertical and horizontal intercepts of $\{40, 1\}$.

- Illustrate the supply and demand curves graphically, and illustrate the skill premium going to the high-skill workers in the short run by determining the value of W when $f=0.5$.
- If demand increases to $W = 60 \times (1 - f)$ what is the new premium? Illustrate your answer graphically.

EXERCISE 13.6

Consider the foregoing problem in a long-run context, when the fraction of the labour force that is high-skilled is more elastic with respect to the premium. Let this long-run relative supply function be $W = 40 \times f$.

- Graph this long-run supply function and verify that it goes through the same initial equilibrium as in Exercise 13.5.
 - Illustrate the long run and short run on the same diagram.
 - What is the numerical value of the premium in the long run after the increase in demand? Illustrate graphically.
1. Foley, K. and D. Green, 2015, "Why more education will not solve rising inequality (and may make it worse)" , *Institute for Research in Public Policy*, Montreal, Canada.

This page titled [13.10: Exercises for Chapter 13](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

SECTION OVERVIEW

Unit 6: Government and Trade

Governments play a major role in virtually every economy. They account for one third or more of national product. In addition to providing a legal and constitutional framework and implementing the law, governments can moderate and influence the operation of markets, provide public goods and furnish missing information. Governments distribute and redistribute widely and supply major services such as health and education. Modern economies could not function without a substantial role for their governments. These roles are explored and developed in Chapter 14.

The final chapter looks outwards. Canada is an open economy, with a high percentage of its production imported and exported. We explore the theory of absolute and comparative advantage and illustrate the potential for consumer gains that trade brings. We analyze barriers to trade such as tariffs and quotas and end with an overview of the World's major trading groups and institutions.

14: Government

- 14.1: Market Failure
- 14.2: Fiscal federalism- Taxing and spending
- 14.3: Federal-provincial fiscal relations
- 14.4: Government-to-individual transfers
- 14.5: Regulation and competition policy
- 14.6: Key Terms
- 14.7: Exercises for Chapter 14

15: International trade

- 15.1: Trade in our daily lives
- 15.2: Canada in the world economy
- 15.3: The gains from trade- Comparative advantage
- 15.4: Returns to scale and dynamic gains from trade
- 15.5: Trade barriers- Tariffs, subsidies and quotas
- 15.6: The politics of protection
- 15.7: Institutions governing trade
- 15.8: Key Terms
- 15.9: Exercises for Chapter 15

This page titled [Unit 6: Government and Trade](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14: Government

Chapter 14: Government

In this chapter we will explore:

14.1	Market failure and the role of government
14.2	Fiscal federalism in Canada
14.3	Federal-provincial relations and powers
14.4	Redistribution to individuals
14.5	Regulatory activity and competition policy

Governments have a profound impact on economies. The economies of Scandinavia are very different from those in North America. North and South Korea are night and day, even though they were identical several decades ago. Canada and Argentina were very similar in the early decades of the twentieth century. Both had abundant space, natural resources and migrants. Today Canada is one of the most prosperous economies in the world, while Argentina struggles as a middle-income economy.

Governments are not peripheral to the marketplace. They can facilitate or hinder the operation of markets. They can ignore poverty or implement policies to support those on low income and those who are incapacitated. Governments can treat the economy as their fiefdoms, as has been the case for decades in many underdeveloped economies. By assuming the role of warlords, local governors inhibit economic development, because the fruits of investment and labour are subject to capture by the ruling power elite.

In Canada we take for granted the existence of a generally benign government that serves the economy, rather than one which expects the economy to serve it. The separation of powers, the existence of a constitution, property rights, a police force, and a free press are all crucial ingredients in the mix that ferments economic development and a healthy society.

The analysis of government is worthy not just of a full course, but a full program of study. Accordingly, our objective in this chapter must be limited. We begin by describing the various ways in which markets may be inadequate and what government can do to remedy these deficiencies. Next we describe the size and scope of government in Canada, and define the sources of government revenues. On the expenditure side, we emphasize the redistributive and transfer roles that are played by Canadian governments. Tax revenues, particularly at the federal level, go predominantly to transfers rather than the provision of goods and services. Finally we examine how governments seek to control, limit and generally influence the marketplace: How do governments foster the operation of markets in Canada? How do they attempt to limit monopolies and cartels? How do they attempt to encourage the entry of new producers and generally promote a market structure that is conducive to competition, economic growth and consumer well-being?

14.1 Market failure

Markets are fine institutions when all of the conditions for their efficient operation are in place. In Chapter 5 we explored the meaning of efficient resource allocation, by developing the concepts of consumer and producer surpluses. But, while we have emphasized the benefits of efficient resource allocation in a market economy, there are many situations where markets deliver inefficient outcomes. Several problems beset the operation of markets. The principal sources of *market failure* are: *Externalities*, *public goods*, *asymmetric information*, and the *concentration of power*. In addition markets may produce outcomes that are *unfavourable* to certain groups – perhaps those on low incomes. The circumstances described here lead to what is termed market failure.

Market failure defines outcomes in which the allocation of resources is not efficient.

Externalities

A negative externality is one resulting, perhaps, from the polluting activity of a producer, or the emission of greenhouse gases into the atmosphere. A positive externality is one where the activity of one individual confers a benefit on others. An example here is where individuals choose to get immunized against a particular illness. As more people become immune, the lower is the probability that the illness can propagate itself through hosts, and therefore the greater the benefits to those not immunized.

Solutions to these market failures come in several forms: Government taxes and subsidies, or quota systems that place limits on the production of products generating externalities. Such solutions were explored in Chapter 5. Taxes on gasoline discourage its use and therefore reduce the emission of poisons into the atmosphere. Taxes on cigarettes and alcohol lower the consumption of goods that may place an additional demand on our publicly-funded health system. The provision of free, or low-cost, immunization against specific diseases to children benefits the whole population.

These measures attempt to *compensate for the absence of a market* in certain activities. Producers may not wish to pay for the right to emit pollutants, and consequently if the government steps in to counter such an externality, the government is effectively implementing a solution to the missing market.

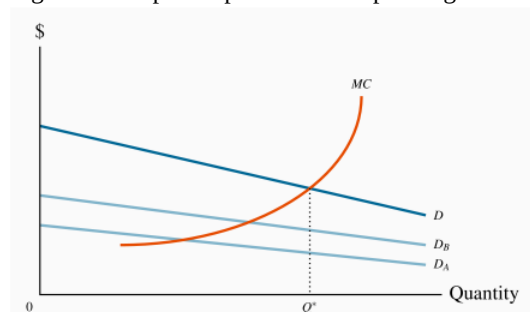
Public goods

Public goods are sometimes called collective consumption goods, on account of their non-rivalrous and non-excludability characteristics. For example, if the government meteorological office provides daily forecasts over the nation's airwaves, it is no more expensive to supply that information to one million than to one hundred individuals in the same region. Its provision to one is not rivalrous with its provision to others – in contrast to private goods that cannot be 'consumed' simultaneously by more than one individual. In addition, it may be difficult to exclude certain individuals from receiving the information.

Public goods are non-rivalrous, in that they can be consumed simultaneously by more than one individual; additionally they may have a non-excludability characteristic.

Examples of such goods and services abound: Highways (up to their congestion point), street lighting, information on trans-fats and tobacco, or public defence provision. Such goods pose a problem for private markets: If it is difficult to exclude individuals from their consumption, then potential private suppliers will likely be deterred from supplying them because the suppliers cannot generate revenue from *free-riders*. Governments therefore normally supply such goods and services. But how much should governments supply? An answer is provided with the help of Figure 14.1.

Figure 14.1 Optimal provision of a public good



The total demand for the public good D is the vertical sum of the individual demands D_A and D_B . The optimal provision is where the MC equals the aggregate marginal valuation, as defined by the demand curve D . At the optimum Q^* , each individual is supplied the same amount of the public good.

This is a supply-demand diagram with a difference. The supply side is conventional, with the MC of production representing the supply curve. An efficient use of the economy's resources, we already know, dictates that an amount should be produced so that the cost at the margin equals the benefit to consumers at the margin. In contrast to the total market demand for private goods, which is obtained by summing individual demands horizontally, the demand for public goods is obtained by summing individual demands *vertically*.

Figure 14.1 depicts an economy with just two individuals whose demands for street lighting are given by D_A and D_B . These demands reveal the value each individual places on the various output levels of the public good, measured on the x-axis. However, since each individual can consume the public good *simultaneously*, the aggregate value of any output produced is the *sum of each individual valuation*. The valuation in the market of any quantity produced is therefore the vertical sum of the individual demands. D is the vertical sum of D_A and D_B , and the optimal output is Q^* . At this equilibrium each individual consumes the same quantity of street lighting, and the MC of the last unit supplied equals the value placed upon it by society – both individuals. Note that this 'optimal' supply depends upon the income distribution, as we have stated several times to date. A different distribution of income may give rise to different demands D_A and D_B , and therefore a different 'optimal' output.

Efficient supply of public goods is where the marginal cost equals the sum of individual marginal valuations, and each individual consumes the same quantity.

Application Box 14.1 Are Wikipedia, Google and MOOCs public goods?

Wikipedia is one of the largest on-line sources of free information in the world. It is an encyclopedia that functions in multiple languages and that furnishes information on millions of topics. It is freely accessible, and is maintained and expanded by its users. Google is the most frequently used search engine on the World Wide Web. It provides information to millions of users simultaneously on every subject imaginable. But it is not quite free of charge; when the user searches she supplies information on herself that can be used profitably by Google in its advertising. MOOCs are 'monster open online courses' offered by numerous universities, frequently for no charge to the student. Are these services public goods in the sense we have described?

Very few goods and services are pure public goods, some have the major characteristics of public goods nonetheless. In this general sense, Google, Wikipedia and MOOCs have public good characteristics. Wikipedia is funded by philanthropic contributions, and its users expand its range by posting information on its servers. Google is funded from advertising revenue. MOOCs are funded by university budgets.

A pure public good is available to additional users at zero marginal cost. This condition is essentially met by these services since their server capacity rarely reaches its limit. Nonetheless, they are constantly adding server capacity, and in that sense cannot furnish their services to an unlimited number of additional users at no additional cost.

Knowledge is perhaps the ultimate public good; Wikipedia, Google and MOOCs all disseminate knowledge, knowledge which has been developed through the millennia by philosophers, scientists, artists, teachers, research laboratories and universities.

A challenge in providing the optimal amount of government-supplied public goods is to know the value that users may place upon them – how can the demand curves D_A and D_B , be ascertained, for example, in Figure 14.1? In contrast to markets for private goods, where consumer demands are essentially revealed through the process of purchase, the demands for public goods may have to be uncovered by means of surveys that are designed so as to elicit the true valuations that users place upon different amounts of a public good. A second challenge relates to the pricing and funding of public goods: For example, should highway lighting be funded from general tax revenue, or should drivers pay for it? These are complexities that are beyond our scope of our current inquiry.

Asymmetric information

Markets for information abound in the modern economy. Governments frequently supply information on account of its public good characteristics. But the problem of *asymmetric information* poses additional challenges. Asymmetric information is where at least one party in an economic relationship has less than full information. This situation characterizes many interactions: Life-insurance companies do not have perfect information on the lifestyle and health of their clients; used vehicle buyers may not know the history of the vehicles they are buying.

Asymmetric information is where at least one party in an economic relationship has less than full information and has a different amount of information from another party.

Asymmetric information can lead to two kinds of problems. The first is adverse selection. For example, can the life-insurance company be sure that it is not insuring only the lives of people who are high risk and likely to die young? If primarily high-risk people buy such insurance then the insurance company must set its premiums accordingly: The company is getting an adverse selection rather than a random selection of clients. Frequently governments decide to run universal compulsory-membership insurance plans (auto or health are examples in Canada) precisely because they may not wish to charge higher rates to higher-risk individuals.

Adverse selection occurs when incomplete or asymmetric information describes an economic relationship.

A related problem is moral hazard. If an individual does not face the full consequences of his actions, his behaviour may be influenced: if a homeowner has a fully insured home he may be less security conscious than an owner who does not.

In Chapter 7 we described how US mortgage providers lent large sums to borrowers with uncertain incomes in the early years of the new millennium. The individuals responsible for the lending were being rewarded on the basis of the amount lent, not the safety of the loan. Nor were the lenders responsible for loans that were not repaid. This 'sub-prime mortgage crisis' was certainly a case of moral hazard.

Moral hazard may characterize behaviour where the costs of certain activities are not incurred by those undertaking them.

Solutions to these problems do not always involve the government, but in critical situations do. For example, the government requires most professional societies and orders to ensure that their members are trained, accredited and capable. Whether for a

medical doctor, a plumber or an engineer, a license or certificate of competence is a signal that the work and advice of these professionals is *bona fide*. Equally, the government sets *standards* so that individuals do not have to incur the cost of ascertaining the quality of their purchases – bicycle helmets must satisfy specific crash norms; so too must air-bags in automobiles.

These situations differ from those where solutions to the information problem can be dealt with reasonably well in the market place. For example, with the advent of buyer and seller rating on *Airbnb*, a potential renter can learn of the quality of the accommodation he is considering, and the letor can assess the potential renter.

Concentration of power

Monopolistic and imperfectly-competitive market structures can give rise to inefficient outcomes, in the sense that the value placed on the last unit of output does not equal the cost at the margin. This arises because the supplier uses his market power in order to maximize profits by limiting output and selling at a higher price.

What can governments do about such power concentrations? Every developed economy has a body similar to Canada's *Competition Bureau*. Such regulatory bodies are charged with seeing that the interests of the consumer, and the economy more broadly, are represented in the market place. Interventions, regulatory procedures and efforts to prevent the abuse of market power come in a variety of forms. These measures are examined in Section 14.5.

Unfavourable market outcomes

Even if governments successfully address the problems posed by the market failures described above, there is nothing to guarantee that market-driven outcomes will be 'fair', or accord with the prevailing notions of justice or equity. The marketplace generates many low-paying jobs, unemployment and poverty. The concentration of economic power has led to the growth in income and wealth inequality in many economies. Governments, to varying degrees, attempt to moderate these outcomes through a variety of social programs and transfers that are discussed in Section 14.4.

14.2 Fiscal federalism: Taxing and spending

Canada is a federal state, in which the federal, provincial and municipal governments exercise different powers and responsibilities. In contrast, most European states are unitary and power is not devolved to their regions to the same degree as in Canada or the US or Australia. Federalism confers several advantages over a unitary form of government where an economy is geographically extensive, or where identifiable differences distinguish one region from another: Regions can adopt different policies in response to the expression of different preferences by their respective voters; smaller governments may be better at experimentation and the introduction of new policies than large governments; political representatives are 'closer' to their constituents.

Despite these advantages, the existence of an additional level of government creates a tension between these levels. Such tension is evident in every federation, and federal and provincial governments argue over the appropriate division of taxation powers and revenue-raising power in general. For example, how should the royalties and taxes from oil and gas deposits offshore be distributed – to the federal government or a provincial government?

In Canada, the federal government collects more in tax revenue than it expends on its own programs. This is a feature of most federations. The provinces simultaneously face a shortfall in their own revenues relative to their program expenditure requirements. The federal government therefore redistributes, or transfers, funds to the provinces so that the latter can perform their constitutionally-assigned roles in the economy. The fact that the federal government bridges this fiscal gap gives it a degree of power over the provinces. This influence is commonly termed federal spending power.

Spending power of a federal government arises when the federal government can influence lower level governments due to its financial rather than constitutional power.

The principal revenue sources for financing federal government activity are given in Figure 14.1 for the fiscal year 2015-16, and the expenditure of these revenues is broken down in Figure 14.2. Further details are accessible in the Department of Finance's 'fiscal reference tables' at www.fin.gc.ca/frt-trf/2019/frt-trf-19-eng.asp. Total revenues for that fiscal year amounted to \$333.2b, and expenditures to \$346.2b.

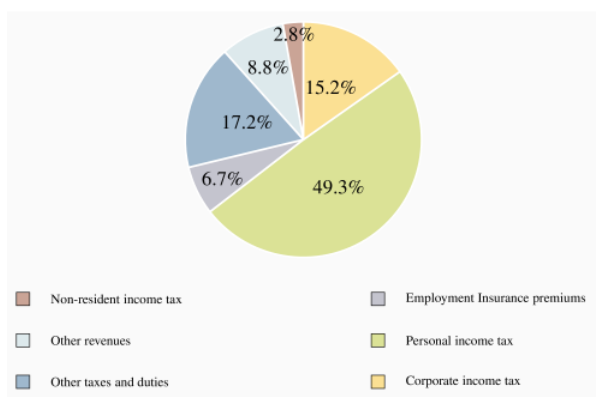


Figure 14.1 Federal Government Revenues 2015–16

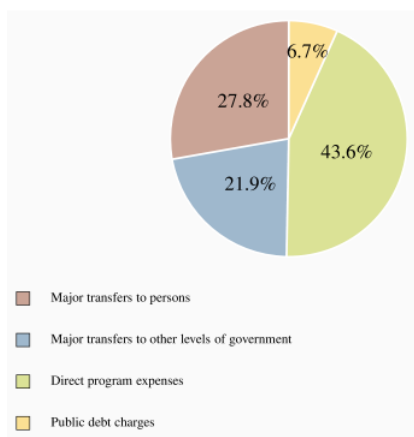


Figure 14.2 Federal expenditures 2015–16

The federal and provincial governments each transfer these revenues to individuals and other levels of government, supply goods and services directly, and also pay interest on accumulated borrowings – the national debt or provincial debt.

Provincial and local governments supply more goods and services than the federal government – health care, drug insurance, education and welfare are the responsibility of provincial and municipal governments. In contrast, national defence, the provision of main traffic arteries, Corrections Canada and a variety of transfer programs to individuals – such as Employment Insurance, Old Age Security and the Canada Pension Plan – are federally funded. The greater part of federal revenues goes towards *transfers* to individuals and provincial governments, as opposed to the supply of goods and services.

14.3 Federal-provincial fiscal relations

The federal government transfers revenue to the provinces using three main programs: Equalization, the Canada Social Transfer and the Canada Health Transfer. Each of these has a different objective. Equalization aims to reduce fiscal disparities among the provinces; The Canada Social Transfer (CST) is for educational and Social Assistance ('welfare') expenditures; The Canada Health Transfer (CHT) performs the same function for health.

Equalization and Territorial Funding

Canada's provinces receive unconditional funding through Canada's Equalization program, whereas the Territories receive federal funding through a separate mechanism - the Territorial Funding Formula.

"Parliament and the Government of Canada are committed to the principle of making equalization payments to ensure that provincial governments have sufficient revenues to provide reasonably comparable levels of public service at reasonably comparable levels of taxation."

This statement, from Section 36(2) of the Constitution Act of 1982, defines the purpose of Equalization. Equalization payments are unconditional – receiving provinces are free to spend the funds on public services according to their own priorities, or even use the revenue to reduce their provincial taxes. Payments are calculated according to a formula that ensures those provinces with revenue-raising ability, or fiscal capacity, below a threshold or 'standard' receive payments from the federal government to bring their capacity up to that standard.

Equalization has gone through very many changes in the several decades of its existence. Its current rules and regulations reflect the 2006 recommendations of a federal Expert Panel. The fiscal capacity of a province is measured by its ability to raise revenues from five major sources: Personal and business income taxes, sales taxes, property taxes, and natural resources. This ability is then compared to the ability of all of the provinces combined to raise revenue; if a difference or shortfall exists, the federal government transfers revenue accordingly, with the amount determined by both the population of the province and the magnitude of its per-person shortfall. Data on annual transfers for Equalization, the Territorial Funding Formula, and the CHT and CST are available at fin.gc.ca/fedprov/mtp-eng.asp

The program transferred \$19.8b to the provinces in 2019-20. The reciprocity status of some provinces varies from year to year. Variation in energy prices and energy-based government revenues are the principal cause of this. British Columbia, Alberta, Saskatchewan and Ontario tend to receive little or zero. Manitoba, Quebec and the Atlantic Provinces have been the major recipient provinces. Quebec receives the largest single amount – approximately two thirds of the total allocation, both on account of its population size and the fact that it has a lower than average fiscal capacity. The Territories received a total of \$3.9b in 2019-20 through the Territorial Funding Formula.

The Canada Social Transfer and the Canada Health Transfer

The CST is a block transfer to provinces in support of post-secondary education, Social Assistance and social services more generally. The CST came into effect in 2004. Prior to that date it was integrated with the health component of federal transfers in a program titled the Canada Health and Social Transfer (CHST). The objective of the separation was to increase the transparency and accountability of federal support for health while continuing to provide funding for other objectives. The CHT is the other part of the unbundled CHST: It provides funding to the provinces for their health expenditures.

The CST and CHT funding comes in two parts: A cash transfer and tax transfer. A tax transfer essentially provides the same support as a cash transfer of equal value, but comes in a different form. In 1977 the federal government agreed with provincial governments to reduce federal personal and corporate tax rates in order to permit the provincial governments to increase the corresponding provincial rates. The net effect was that the federal government got less tax revenue and the provinces got more. And to this day, the federal and provincial governments keep a record of the implied tax transfers that arise from this long-ago agreement. This is the tax transfer component of the CST and the CHT.

The CST support is allocated to provinces and territories on an equal per-capita basis to ensure equal support for all Canadians regardless of their place of residence. The CHT is distributed likewise, and it requires the provinces to abide by the federally-legislated *Canada Health Act*, which demands that provincial health coverage be comprehensive, universal, portable, accessible and publicly administered.

The CHT transfer amounted to \$40.4b and the CST amounted to \$14.6b, in cash, for the year 2019-20. Health care and health expenditures are a core issue of policy at all levels of government on account of the envisaged growth in these expenditures that will inevitably accompany the aging of the baby-boomers.

14.4 Government-to-individual transfers

Many Canadians take pride in Canada's extensive 'social safety net' that aims to protect individuals from misfortune and the reduction of income in old age. Others believe it is too generous. While it is more supportive than the safety net in the US, the Canadian safety net is no more protective than the nets of the developed economies of the European Union. The extent of such support depends in large measure upon the degree to which governments are willing to impose, and individuals are willing to pay, higher or lower tax rates. The major elements of this umbrella of programs are the following.

The *Canada and Quebec Pension Plans (C/QPP)* are funded from the contributions of workers and their employers. Contributions form 9.9% of an individual's earnings up to the maximum pensionable earnings (MPE) figure of \$57,400 in 2019. The contributions are shared equally by employer and employee. The Canada and Quebec components of the plan operate similarly, but are managed separately. The contribution rate to the QPP stands at 10.65%. Contributions to the plans from workers and their employers are largely transferred immediately to retired workers. Part of the contributions is invested in a fund. The objective of

the plans is to ensure that some income is saved for retirement. Many individuals are not very good at planning – they constantly postpone the decision to save, so the state steps in and requires them to save. An individual contributing throughout a full-time working lifecycle can expect an annual pension of about \$14,000 in 2019. The Plans provide a maximum payout of 25% of maximum insurable earnings. The objective is to provide a minimum level of retirement income, not an income that will see individuals live in great comfort.

The C/QPP plans have contributed greatly to the reduction of poverty among the elderly since their introduction in the mid-sixties. The aging of the baby-boom generation – that very large cohort born in the late forties through to the early sixties – means that the percentage of the population in the post-65 age group has begun to increase. To meet this changing demographic, the federal and provincial governments reshaped the plans in the late nineties in order to put them on a sound financial footing – primarily by increasing contributions, that in turn will enable the build-up of a CPP 'fund' that will support the aged in the following decades.

A number of recent studies in Canada on the retirement savings practices of Canadians have proposed that households on average are not saving a sufficient amount for their retirement; many households may thus see a notable decline in their incomes upon retirement. In response to this finding, the federal government agreed with the provinces in June 2016, to add a supplement to the CPP. The new federal provisions, which will be phased in over the period 2019-2025, envisage an increase in contributions that will ultimately lead to a maximum replacement rate of 33% of MPE as opposed to the current goal of 25%. However, full benefits will be experienced only by individuals contributing for their complete lifecycle, meaning that full implementation will take about four decades.

Details of the CPP and the 2016 enhancements are to be found at www.fin.gc.ca/n16/data/16-113_3-eng.asp.

Old Age Security (OAS), the *Guaranteed Income Supplement (GIS)* and the *Spousal Allowance (SPA)* together form the second support leg for the retired. OAS is a payment made automatically to individuals once they attain the age of 65. The GIS is an additional payment made only to those on very low incomes – for example, individuals who have little income from their C/QPP or private pension plans. The SPA, which is payable to the spouse or survivor of an OAS recipient, accounts for a small part of the sums disbursed. As of 2019 the maximum annual OAS payment stood at \$7,360. The federal government in 2016 reversed a plan that would have seen the eligible age for receipt of OAS move to 67.

The payments for these plans come from the general tax revenues of the federal government. Unlike the C/QPP, the benefits received are not related to the contributions that an individual makes over the working lifecycle. This program has also had a substantial impact on poverty reduction among the elderly.

Employment Insurance (EI) and *Social Assistance (SA)* are designed to support, respectively, the unemployed and those with no other source of income. Welfare is the common term used to describe SA. Expenditures on EI and SA are strongly cyclical. At the trough of an economic cycle the real value of expenditures on these programs greatly exceeds expenditures at the peak of the cycle. Unemployment in Canada rose above 8% in 2009, and payments to the unemployed and those on welfare reflected this dire state. The strongly cyclical pattern of the cost of these programs reflects the importance of a healthy job market: Macroeconomic conditions have a major impact on social program expenditures.

EI is funded by contributions from employees and their employers. For each dollar contributed by the employee, the employer contributes \$1.4. Premiums are paid on earned income up to a maximum insurable earnings (MIE) of \$53,100 in 2019. The contribution rate for employees stood at 1.62% of MIE in 2017. EI contributions and pay-outs form part of the federal government's general revenues and expenditures. There is no separate 'fund' for this program. However it is expected to operate on a break-even basis over the longer term. To reflect this, the contribution rate fluctuates with a view to maintaining a balance between payouts and revenues over a seven-year planning period.

EI is called an insurance program, but in reality it is much more than that. Certain groups systematically use the program more than others – those in seasonal jobs, those in rural areas and those in the Atlantic Provinces, for example. Accordingly, using the terminology of Chapter 7, it is not everywhere an actuarially 'fair' insurance program. Benefits payable to unemployed individuals may also depend on their family size, in addition to their work history. While most payments go in the form of 'regular' benefits to unemployed individuals, the EI program also sponsors employee retraining, family benefits that cover maternity and paternity leave, and some other specific target programs for the unemployed.

Social Assistance is provided to individuals who are in serious need of financial support – having no income and few assets. Provincial governments administer SA, although the cost of the program is partly covered by federal transfers through the Canada Social Transfer. The nineteen nineties witnessed a substantial tightening of regulations virtually across the whole of Canada. Access to SA benefits is now more difficult, and benefits have fallen in real terms since the late nineteen eighties.

Welfare dependence peaked in Canada in 1994, when 3.1 million individuals were dependent upon support. As of 2019, the total is approximately half of this, on account of more stringent access conditions, reduced benefit levels and an improved job market. Some groups in Canada believe that benefits should be higher, others believe that making welfare too generous provides young individuals with the wrong incentives in life, and may lead them to neglect schooling and skill development.

Workers Compensation supports workers injured on the job. Worker/employer contributions and general tax revenue form the sources of program revenue, and the mix varies from province-to-province. In contrast to the macro-economy-induced swings in expenditures that characterize SA and EI since the early nineties, expenditures on Worker's Compensation have remained relatively constant.

Canada Child Benefit: The major remaining pillar in Canada's social safety net is the group of payments and tax credits aimed at supporting children: The Canada Child Tax Benefit (CCTB), the Universal Child Care Benefit and the National Child Benefit Supplement were repackaged in 2016 under the title Canada Child Benefit. Child support has evolved and been enriched over the last two decades, partly with the objective of reducing poverty among households with children, and partly with a view to helping parents receiving social assistance to transition back to the labour market. As of 2016, the federal government provides an annual payment to families with children. For each child under the age of 6 the payment is \$6,639 and for each child aged 6-17 the payment is \$5,602. Since these payment are primarily intended for households with low and middle incomes, the amounts are progressively clawed back once the household income reaches a threshold of \$30,000.

Application Box 14.2 Government debts and deficits

Canada's expenditure and tax policies in the nineteen seventies and eighties led to the accumulation of large government debts, as a result of running fiscal deficits. By the mid-nineties the combined federal and provincial debt reached 100% of GDP, with the federal debt accounting for the larger share. This ratio was perilously high: Interest payments absorbed a large fraction of annual government revenues, which in turn limited the ability of the government to embark on new programs or enrich existing ones. Canada's debt rating on international financial markets declined.

In 1995, Finance Minister Paul Martin addressed this problem, and over the following years program spending was pared back. Ultimately, the economy expanded and by the end of the decade the annual deficits at the federal level were eliminated.

As of 2007 the ratio of combined federal and provincial debts stood at just over 60% of GDP. However, the Great Recession of 2008 and following years saw all levels of government experience deficits, with the result that this ratio of combined debt to GDP rose again. Growth in recent years has seen that ratio fall. As of 2018-19, federal interest payments on its debt account for about 7% of its revenues (this figure stood at 28% in the early nineties). At the time of writing, interest rates are low in developed economies and so the interest costs of government debt are low. Low borrowing costs are a reason why some people favor large government spending in the form of infrastructure projects. Those who are fiscally more conservative fear rising rates in the future. The recessionary impacts of the coronavirus pandemic of 2020 will add greatly to accumulated debt, particularly at the federal level.

Debts can be measured in more than a single manner. One measure of debt is the value of all federal government bonds and financial liabilities outstanding. As of 2019-20 this value was approximately \$700b. In addition to this, the federal government has outstanding liabilities to the pensions of its retired employees, and not all of these liabilities have been covered by the contributions of those employees into their pension plans. The federal government also owns assets, both financial and physical - such as office buildings. Hence these assets offset the financial liabilities. To assess the total debt picture of the Canadian economy we need to add provincial and local government debts to the federal debts, and then consider the annual interest costs of this total. It turns out that the interest costs are just above 2% of GDP.

Source: Government of Canada, Fiscal Reference Tables: www.fin.gc.ca/frt-trf/2019/frt-trf-19-eng.asp

14.5 Regulation and competition policy

Goals and objectives

The goals of competition policy are relatively uniform across developed economies: The promotion of domestic competition; the development of new ideas, new products and new enterprises; the promotion of efficiency in the resource-allocation sense; the development of manufacturing and service industries that can compete internationally.

In addition to these economic objectives, governments and citizens frown upon monopolies or monopoly practices if they lead to an undue *concentration of political power*. Such power can lead to a concentration of wealth and influence in the hands of an elite.

Canada's regulatory body is the *Competition Bureau*, whose activity is governed primarily by the *Competition Act* of 1986. This act replaced the *Combines Investigation Act*. The *Competition Tribunal* acts as an adjudication body, and is composed of judges and non-judicial members. This tribunal can issue orders on the maintenance of competition in the marketplace. Canada has had anti-combines legislation since 1889, and the act of 1986 is the most recent form of such legislation and policy. The Competition Act does not forbid monopolies, but it does rule as unlawful the *abuse* of monopoly power. Canada's competition legislation is aimed at anti-competitive practices, and a full description of its activities is to be found on its website at www.competitionbureau.gc.ca. Let us examine some of these proscribed policies.

Anti-competitive practices

Anti-competitive practices may either limit entry into a sector of the economy or force existing competitors out. In either case they lead to a reduction in competition.

Mergers may turn competitive firms into a single organization with excessive market power. The customary justification for mergers is that they permit the merged firms to achieve scale economies that would otherwise be impossible. Such scale economies may in turn result in lower prices in the domestic or international market to the benefit of the consumer, but may alternatively reduce competition and result in higher prices. Equally important in this era of global competition is the impact of a merger on a firm's ability to compete internationally.

Mergers can be of the horizontal type (e.g. two manufacturers of pre-mixed concrete merge) or vertical type (a concrete manufacturer merges with a cement manufacturer). In a market with few suppliers mergers have the potential to reduce domestic competition.

Cartels aim to restrict output and thereby increase profits. These formations are almost universally illegal in individual national economies.

While cartels are one means of increasing prices, price discrimination is another, as we saw when studying monopoly behaviour. For example, if a concrete manufacturer makes her product available to large builders at a lower price than to small-scale builders – perhaps because the large builder has more bargaining power – then the small builder is at a competitive disadvantage in the construction business. If the small firm is forced out of the construction business as a consequence, then competition in this sector is reduced.

We introduced the concept of predatory pricing in Chapter 11. Predatory pricing is a practice that is aimed at driving out competition by artificially reducing the price of one product sold by a supplier. For example, a dominant nationwide transporter could reduce price on a particular route where competition comes from a strictly local competitor. By 'subsidizing' this route from profits on other routes, the dominant firm could undercut the local firm and drive it out of the market.

Predatory pricing is a practice that is aimed at driving out competition by artificially reducing the price of one product sold by a supplier.

Suppliers may also refuse to deal. If the local supplier of pre-mixed concrete refuses to sell the product to a local construction firm, then the ability of such a downstream firm to operate and compete may be compromised. This practice is similar to that of exclusive sales and tied sales. An exclusive sale might involve a large vegetable wholesaler forcing her retail clients to buy *only* from this supplier. Such a practice might hurt the local grower of aubergines or zucchini, and also may prevent the retailer from obtaining *some* of her vegetables at a lower price or at a higher quality elsewhere. A tied sale is one where the purchaser must agree to purchase a *bundle* of goods from a supplier.

Refusal to deal: an illegal practice where a supplier refuses to sell to a purchaser.

Exclusive sale: where a retailer is obliged (perhaps illegally) to purchase all wholesale products from a single supplier only.

Tied sale: one where the purchaser must agree to purchase a bundle of goods from a supplier.

Resale price maintenance involves the producer requiring a retailer to sell a product at a specified price. This practice can hurt consumers since they cannot 'shop around'. In Canada, we frequently encounter a 'manufacturer's suggested retail price' for autos and durable goods. But since these prices are not *required*, the practice conforms to the law.

Resale price maintenance is an illegal practice wherein a producer requires sellers to maintain a specified price.

Bid rigging is an illegal practice in which normally competitive bidders conspire to fix the awarding of contracts or sales. For example, two builders, who consider bidding on construction projects, may decide that one will bid seriously for project X and the other will bid seriously on project Y. In this way they conspire to reduce competition in order to make more profit.

Bid rigging is an illegal practice in which bidders (buyers) conspire to set prices in their own interest.

Deception and dishonesty in promoting products can either short-change the consumer or give one supplier an unfair advantage over other suppliers.

Enforcement

The Competition Act is enforced through the Competition Bureau in a variety of ways. Decisions on acceptable business practices are frequently reached through study and letters of agreement between the Bureau and businesses. In some cases, where laws appear to have been violated, criminal proceedings may follow.

Regulation, deregulation and privatization

The last three decades have witnessed a significant degree of privatization and deregulation in Canada, most notably in the transportation, communication and energy sectors. Modern deregulation in the US began with the passage of the *Airline Deregulation Act* of 1978, and was pursued with great energy under the Reagan administration in the eighties. The Economic Council of Canada produced an influential report in 1981, titled "Reforming Regulation," on the impact of regulation and possible deregulation of specific sectors. The Economic Council proposed that regulation in some sectors was inhibiting competition, entry and innovation. As a consequence, the interests of the consumer were in danger of becoming secondary to the interests of the suppliers.

Telecommunications provision, in the era when the telephone was the main form of such communication, was traditionally viewed as a natural monopoly. The Canadian Radio and Telecommunications Commission (CRTC) regulated its rates. The industry has developed dramatically in the last two decades with the introduction of satellite-facilitated communication, the internet, multi-purpose cable networks, cell phones and service integration.

Transportation, in virtually all forms, has been deregulated in Canada since the nineteen eighties. Railways were originally required to subsidize the transportation of grain under the *Crow's Nest Pass* rate structure. But the subsidization of particular markets requires an excessive rate elsewhere, and if the latter markets become subject to competition then a competitive system cannot function. This structure, along with many other anomalies, was changed with the passage of the *Canada Transportation Act* in 1996.

Trucking, historically, has been regulated by individual provinces. Entry was heavily controlled prior to the federal *National Transportation Act* of 1987, and subsequent legislation introduced by a number of provinces, have made for easier entry and a more competitive rate structure.

Deregulation of the airline industry in the US in the late seventies had a considerable influence on thinking and practice in Canada. The Economic Council report of 1981 recommended in favour of easier entry and greater fare competition. These policies were reflected in the 1987 National Transportation Act. Most economists are favourable to deregulation and freedom to enter, and the US experience indicated that cost reductions and increased efficiency could follow. In 1995 an agreement was reached between the US and Canada that provided full freedom for Canadian carriers to move passengers to any US city, and freedom for US carriers to do likewise, subject to a phase-in provision.

The National Energy Board regulates the development and transmission of oil and natural gas. But earlier powers of the Board, involving the regulation of product prices, were eliminated in 1986, and controls on oil exports were also eliminated.

Agriculture remains a highly controlled area of the economy. Supply 'management', which is really supply restriction, and therefore 'price maintenance', characterizes grain, dairy, poultry and other products. Management is primarily through provincial marketing boards.

The role of the sharing economy

The arrival of universal access to the internet has seen the emergence of what is known as the Sharing economy throughout the world. This expression is used to describe commercial activities that, in the first place, are internet-based. Second, suppliers in the sharing economy use resources in the market place that were initially aimed at a different purpose. *Airbnb* and *Uber* are good examples of companies in sectors of the economy where sharing is possible. In *Uber's* case, the 'ride-share' drivers initially purchased their vehicles for private use, and subsequently redirected them to commercial use. *Airbnb* is a communication corporation that enables the owners of spare home capacity to sell the use of that capacity to short-term renters. With the maturation of such corporations, the concept of 'initial' and 'secondary' use becomes blurred.

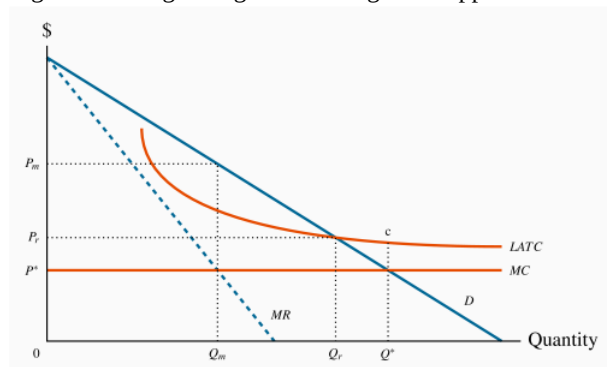
Sharing economy: involves enterprises that are internet based, and that use production resources that have use outside of the marketplace.

The importance of the sharing economy is that it provides an additional source of competition to established suppliers, and therefore limits the market power of the latter. At the same time, the emergence of the sharing economy poses a new set of regulatory challenges: If traditional taxis are required to purchase operating permits (medallions), and the ride-share drivers do not require such permits, is there a reasonable degree of competition in the market, and if not what is the appropriate solution? Should the medallion requirement be abolished, or should ride-share drivers be required to purchase one? In the case of *Airbnb*, the suppliers operate outside of the traditional 'hotel' market. In general they do not charge sales taxes or face any union labour agreements. What is the appropriate response from governments? And how should the sharing economy be taxed?

Price regulation

Regulating monopolistic sectors of the economy is one means of reducing their market power. In Chapter 11 it was proposed that indefinitely decreasing production costs in an industry means that the industry might be considered as a 'natural' monopoly: Higher output can be produced at lower cost with fewer firms. Hence, a single supplier has the potential to supply the market at a lower unit cost; unless, that is, such a single supplier uses his monopoly power. To illustrate how the consumer side may benefit from this production structure through regulation, consider Figure 14.2. For simplicity suppose that long-run marginal costs are constant and that average costs are downward sloping due to an initial fixed cost. The profit-maximizing (monopoly) output is where $MR=MC$ at Q_m and is sold at the price P_m . This output is inefficient because the willingness of buyers to pay for additional units of output exceeds the additional cost. On this criterion the efficient output is Q^* . But $LATC$ exceeds price at Q^* , and therefore it is not feasible for a producer.

Figure 14.2 Regulating a decreasing-cost supplier



The profit-maximizing output is Q_m , where $MR=MC$ and price is P_m . This output is inefficient because marginal benefit is greater than MC . Q^* is the efficient output, but results in losses because $LATC > P$ at that output. A regulated price that covers costs is where $LATC=DQ_r$. This is closer to the efficient output Q^* than the monopoly output Q_m .

One solution is for the regulating body to set a price-quantity combination of P_r and Q_r , where price equals average cost and therefore generates a normal rate of profit. This output level is still lower than the efficient output level Q^* , but is more efficient than the profit-maximizing output Q_m . It is more efficient in the sense that it is closer to the efficient output Q^* . A problem with such a strategy is that it may induce lax management: If producers are allowed to charge an average-cost price, then there is a reduced incentive for them to keep strict control of their costs in the absence of competition in the marketplace.

A second solution to the declining average cost phenomenon is to implement what is called a two-part tariff. This means that customers pay an 'entry fee' in order to be able to purchase the good. For example, in many jurisdictions hydro or natural gas subscribers may pay a fixed charge per month for their supply line and supply guarantee, and then pay an additional charge that varies with quantity. In this way it is possible for the supplier to charge a price per unit of output that is closer to marginal cost and still make a profit, than under an average cost pricing formula. In terms of Figure 14.2, the total value of entry fees, or fixed components of the pricing, would have to cover the difference between MC and $LATC$ times the output supplied. In Figure 14.2 this implies that if the efficient output Q^* is purchased at a price equal to the MC the producer loses the amount $(c-MC)$ on each unit sold. The access fees would therefore have to cover at least this value.

Such a solution is appropriate when fixed costs are high and marginal costs are low. This situation is particularly relevant in the modern market for telecommunications: The cost to suppliers of marginal access to their networks, whether it be for internet, phone or TV, is negligible compared to the cost of maintaining the network and installing capacity.

Two-part tariff: involves an access fee and a per unit of quantity fee.

Finally, a word of caution: Nobel Laureate George Stigler has argued that there is a danger of regulators becoming too close to the regulated, and that the relationship can evolve to a point where the regulator may protect the regulated firms. In contrast, Professor Philippon of New York University argues that regulators are not regulating sufficiently in the US: they have permitted an excessive number of mergers that have, in turn, reduced competition.

Key Terms

Market failure defines outcomes in which the allocation of resources is not efficient.

Public goods are non-rivalrous, in that they can be consumed simultaneously by more than one individual; additionally they may have a non-excludability characteristic.

Efficient supply of public goods is where the marginal cost equals the sum of individual marginal valuations, and each individual consumes the same quantity.

Asymmetric information is where at least one party in an economic relationship has less than full information and has a different amount of information from another party.

Adverse selection occurs when incomplete or asymmetric information describes an economic relationship.

Moral hazard may characterize behaviour where the costs of certain activities are not incurred by those undertaking them.

Spending power of a federal government arises when the federal government can influence lower level governments due to its financial rather than constitutional power.

Predatory pricing is a practice that is aimed at driving out competition by artificially reducing the price of one product sold by a supplier.

Refusal to deal: an illegal practice where a supplier refuses to sell to a purchaser.

Exclusive sale: where a retailer is obliged (perhaps illegally) to purchase all wholesale products from a single supplier only.

Tied sale: one where the purchaser must agree to purchase a bundle of goods from the one supplier.

Resale price maintenance is an illegal practice wherein a producer requires sellers to maintain a specified price.

Bid rigging is an illegal practice in which bidders (buyers) conspire to set prices in their own interest.

Sharing economy: involves enterprises that are internet based, and that use production resources that have use outside of the marketplace.

Two-part tariff: involves an access fee and a per unit of quantity fee.

Exercises for Chapter 14

EXERCISE 14.1

An economy is composed of two individuals, whose demands for a public good – street lighting – are given by $P=12-(1/2)Q$ and $P=8-(1/3)Q$.

1. Graph these demands on a diagram, for values of $Q = 1, \dots, 24$.
2. Graph the total demand for this public good by summing the demands vertically, specifying the numerical value of each intercept.
3. Let the marginal cost of providing the good be \$5 per unit. Illustrate graphically the efficient supply of the public good (Q^*) in this economy.
4. Illustrate graphically the area that represents the total value to the consumers of the amount Q^* .

EXERCISE 14.2

In Exercise 14.1, suppose a new citizen joins the economy, and her demand for the public good is given by $P=10-(5/12)Q$.

1. Add this individual's demand curve to the graphic for the above question and graph the new total demand curve, specifying the intercept values.
2. Illustrate the area on your graph that represents the new total value to the three citizens of the optimal amount supplied.

3. Illustrate graphically the net value to society of the new Q^* – the total value minus the total cost.

EXERCISE 14.3

An industry that is characterized by a decreasing cost structure has a demand curve given by $P=100-Q$ and the marginal revenue curve by $MR=100-2Q$. The marginal cost is $MC=4$, and average cost is $AC=4+188/Q$.

1. Graph this cost and demand structure. [*Hint*: This graph is similar to Figure 14.2.]
2. Illustrate the efficient output and the monopoly output for the industry.
3. Illustrate on the graph the price the monopolist would charge if he were unregulated.

EXERCISE 14.4

Optional: In Question 14.3, suppose the government decides to regulate the behaviour of the supplier, in the interests of the consumer.

1. Illustrate graphically the price and output that would emerge if the supplier were regulated so that his allowable price equalled average cost.
2. Is this greater or less than the efficient output?
3. Compute the AC and P that would be charged with this regulation.
4. Illustrate graphically the deadweight loss associated with the regulated price and compare it with the deadweight loss under monopoly.

EXERCISE 14.5

Optional: As an alternative to regulating the supplier such that price covers average total cost, suppose that a two part tariff were used to generate revenue. This scheme involves charging the MC for each unit that is purchased and in addition charging each buyer in the market a fixed cost that is independent of the amount he purchases. If an efficient output is supplied in the market, illustrate graphically the total revenue to be obtained from the component covering a price per unit of the good supplied, and the component covering fixed cost.

This page titled [14: Government](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14.1: Market Failure

Markets are fine institutions when all of the conditions for their efficient operation are in place. In Chapter 5 we explored the meaning of efficient resource allocation, by developing the concepts of consumer and producer surpluses. But, while we have emphasized the benefits of efficient resource allocation in a market economy, there are many situations where markets deliver inefficient outcomes. Several problems beset the operation of markets. The principal sources of *market failure* are: *Externalities*, *public goods*, *asymmetric information*, and the *concentration of power*. In addition markets may produce outcomes that are *unfavourable* to certain groups – perhaps those on low incomes. The circumstances described here lead to what is termed market failure.

Market failure defines outcomes in which the allocation of resources is not efficient.

Externalities

A negative externality is one resulting, perhaps, from the polluting activity of a producer, or the emission of greenhouse gases into the atmosphere. A positive externality is one where the activity of one individual confers a benefit on others. An example here is where individuals choose to get immunized against a particular illness. As more people become immune, the lower is the probability that the illness can propagate itself through hosts, and therefore the greater the benefits to those not immunized.

Solutions to these market failures come in several forms: Government taxes and subsidies, or quota systems that place limits on the production of products generating externalities. Such solutions were explored in Chapter 5. Taxes on gasoline discourage its use and therefore reduce the emission of poisons into the atmosphere. Taxes on cigarettes and alcohol lower the consumption of goods that may place an additional demand on our publicly-funded health system. The provision of free, or low-cost, immunization against specific diseases to children benefits the whole population.

These measures attempt to *compensate for the absence of a market* in certain activities. Producers may not wish to pay for the right to emit pollutants, and consequently if the government steps in to counter such an externality, the government is effectively implementing a solution to the missing market.

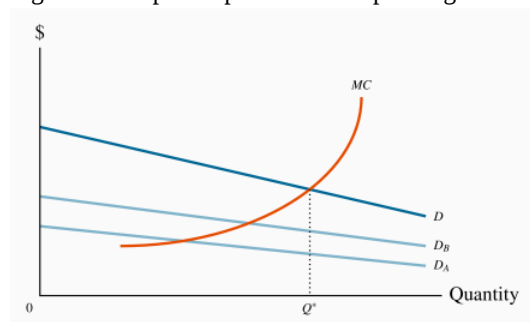
Public goods

Public goods are sometimes called collective consumption goods, on account of their non-rivalrous and non-excludability characteristics. For example, if the government meteorological office provides daily forecasts over the nation's airwaves, it is no more expensive to supply that information to one million than to one hundred individuals in the same region. Its provision to one is not rivalrous with its provision to others – in contrast to private goods that cannot be 'consumed' simultaneously by more than one individual. In addition, it may be difficult to exclude certain individuals from receiving the information.

Public goods are non-rivalrous, in that they can be consumed simultaneously by more than one individual; additionally they may have a non-excludability characteristic.

Examples of such goods and services abound: Highways (up to their congestion point), street lighting, information on trans-fats and tobacco, or public defence provision. Such goods pose a problem for private markets: If it is difficult to exclude individuals from their consumption, then potential private suppliers will likely be deterred from supplying them because the suppliers cannot generate revenue from *free-riders*. Governments therefore normally supply such goods and services. But how much should governments supply? An answer is provided with the help of Figure 14.1.

Figure 14.1 Optimal provision of a public good



The total demand for the public good D is the vertical sum of the individual demands D_A and D_B . The optimal provision is where the MC equals the aggregate marginal valuation, as defined by the demand curve D . At the optimum Q^* , each individual is supplied the same amount of the public good.

This is a supply-demand diagram with a difference. The supply side is conventional, with the MC of production representing the supply curve. An efficient use of the economy's resources, we already know, dictates that an amount should be produced so that the cost at the margin equals the benefit to consumers at the margin. In contrast to the total market demand for private goods, which is obtained by summing individual demands horizontally, the demand for public goods is obtained by summing individual demands *vertically*.

Figure 14.1 depicts an economy with just two individuals whose demands for street lighting are given by D_A and D_B . These demands reveal the value each individual places on the various output levels of the public good, measured on the x -axis. However, since each individual can consume the public good *simultaneously*, the aggregate value of any output produced is the *sum of each individual valuation*. The valuation in the market of any quantity produced is therefore the vertical sum of the individual demands. D is the vertical sum of D_A and D_B , and the optimal output is Q^* . At this equilibrium each individual consumes the same quantity of street lighting, and the MC of the last unit supplied equals the value placed upon it by society – both individuals. Note that this 'optimal' supply depends upon the income distribution, as we have stated several times to date. A different distribution of income may give rise to different demands D_A and D_B , and therefore a different 'optimal' output.

Efficient supply of public goods is where the marginal cost equals the sum of individual marginal valuations, and each individual consumes the same quantity.

Application Box 14.1 Are Wikipedia, Google and MOOCs public goods?

Wikipedia is one of the largest on-line sources of free information in the world. It is an encyclopedia that functions in multiple languages and that furnishes information on millions of topics. It is freely accessible, and is maintained and expanded by its users. Google is the most frequently used search engine on the World Wide Web. It provides information to millions of users simultaneously on every subject imaginable. But it is not quite free of charge; when the user searches she supplies information on herself that can be used profitably by Google in its advertising. MOOCs are 'monster open online courses' offered by numerous universities, frequently for no charge to the student. Are these services public goods in the sense we have described?

Very few goods and services are pure public goods, some have the major characteristics of public goods nonetheless. In this general sense, Google, Wikipedia and MOOCs have public good characteristics. Wikipedia is funded by philanthropic contributions, and its users expand its range by posting information on its servers. Google is funded from advertising revenue. MOOCs are funded by university budgets.

A pure public good is available to additional users at zero marginal cost. This condition is essentially met by these services since their server capacity rarely reaches its limit. Nonetheless, they are constantly adding server capacity, and in that sense cannot furnish their services to an unlimited number of additional users at no additional cost.

Knowledge is perhaps the ultimate public good; Wikipedia, Google and MOOCs all disseminate knowledge, knowledge which has been developed through the millennia by philosophers, scientists, artists, teachers, research laboratories and universities.

A challenge in providing the optimal amount of government-supplied public goods is to know the value that users may place upon them – how can the demand curves D_A and D_B , be ascertained, for example, in Figure 14.1? In contrast to markets for private goods, where consumer demands are essentially revealed through the process of purchase, the demands for public goods may have to be uncovered by means of surveys that are designed so as to elicit the true valuations that users place upon different amounts of a public good. A second challenge relates to the pricing and funding of public goods: For example, should highway lighting be funded from general tax revenue, or should drivers pay for it? These are complexities that are beyond our scope of our current inquiry.

Asymmetric information

Markets for information abound in the modern economy. Governments frequently supply information on account of its public good characteristics. But the problem of *asymmetric information* poses additional challenges. Asymmetric information is where at least one party in an economic relationship has less than full information. This situation characterizes many interactions: Life-insurance companies do not have perfect information on the lifestyle and health of their clients; used vehicle buyers may not know the history of the vehicles they are buying.

Asymmetric information is where at least one party in an economic relationship has less than full information and has a different amount of information from another party.

Asymmetric information can lead to two kinds of problems. The first is adverse selection. For example, can the life-insurance company be sure that it is not insuring only the lives of people who are high risk and likely to die young? If primarily high-risk people buy such insurance then the insurance company must set its premiums accordingly: The company is getting an adverse selection rather than a random selection of clients. Frequently governments decide to run universal compulsory-membership insurance plans (auto or health are examples in Canada) precisely because they may not wish to charge higher rates to higher-risk individuals.

Adverse selection occurs when incomplete or asymmetric information describes an economic relationship.

A related problem is moral hazard. If an individual does not face the full consequences of his actions, his behaviour may be influenced: if a homeowner has a fully insured home he may be less security conscious than an owner who does not.

In Chapter 7 we described how US mortgage providers lent large sums to borrowers with uncertain incomes in the early years of the new millennium. The individuals responsible for the lending were being rewarded on the basis of the amount lent, not the safety of the loan. Nor were the lenders responsible for loans that were not repaid. This 'sub-prime mortgage crisis' was certainly a case of moral hazard.

Moral hazard may characterize behaviour where the costs of certain activities are not incurred by those undertaking them.

Solutions to these problems do not always involve the government, but in critical situations do. For example, the government requires most professional societies and orders to ensure that their members are trained, accredited and capable. Whether for a medical doctor, a plumber or an engineer, a license or certificate of competence is a signal that the work and advice of these professionals is *bona fide*. Equally, the government sets *standards* so that individuals do not have to incur the cost of ascertaining the quality of their purchases – bicycle helmets must satisfy specific crash norms; so too must air-bags in automobiles.

These situations differ from those where solutions to the information problem can be dealt with reasonably well in the market place. For example, with the advent of buyer and seller rating on *Airbnb*, a potential renter can learn of the quality of the accommodation he is considering, and the letor can assess the potential renter.

Concentration of power

Monopolistic and imperfectly-competitive market structures can give rise to inefficient outcomes, in the sense that the value placed on the last unit of output does not equal the cost at the margin. This arises because the supplier uses his market power in order to maximize profits by limiting output and selling at a higher price.

What can governments do about such power concentrations? Every developed economy has a body similar to Canada's *Competition Bureau*. Such regulatory bodies are charged with seeing that the interests of the consumer, and the economy more broadly, are represented in the market place. Interventions, regulatory procedures and efforts to prevent the abuse of market power come in a variety of forms. These measures are examined in Section 14.5.

Unfavourable market outcomes

Even if governments successfully address the problems posed by the market failures described above, there is nothing to guarantee that market-driven outcomes will be 'fair', or accord with the prevailing notions of justice or equity. The marketplace generates many low-paying jobs, unemployment and poverty. The concentration of economic power has led to the growth in income and wealth inequality in many economies. Governments, to varying degrees, attempt to moderate these outcomes through a variety of social programs and transfers that are discussed in Section 14.4.

This page titled [14.1: Market Failure](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis](#) and [Ian Irvine](#) ([Lyryx](#)) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14.2: Fiscal federalism- Taxing and spending

Canada is a federal state, in which the federal, provincial and municipal governments exercise different powers and responsibilities. In contrast, most European states are unitary and power is not devolved to their regions to the same degree as in Canada or the US or Australia. Federalism confers several advantages over a unitary form of government where an economy is geographically extensive, or where identifiable differences distinguish one region from another: Regions can adopt different policies in response to the expression of different preferences by their respective voters; smaller governments may be better at experimentation and the introduction of new policies than large governments; political representatives are 'closer' to their constituents.

Despite these advantages, the existence of an additional level of government creates a tension between these levels. Such tension is evident in every federation, and federal and provincial governments argue over the appropriate division of taxation powers and revenue-raising power in general. For example, how should the royalties and taxes from oil and gas deposits offshore be distributed – to the federal government or a provincial government?

In Canada, the federal government collects more in tax revenue than it expends on its own programs. This is a feature of most federations. The provinces simultaneously face a shortfall in their own revenues relative to their program expenditure requirements. The federal government therefore redistributes, or transfers, funds to the provinces so that the latter can perform their constitutionally-assigned roles in the economy. The fact that the federal government bridges this fiscal gap gives it a degree of power over the provinces. This influence is commonly termed federal spending power.

Spending power of a federal government arises when the federal government can influence lower level governments due to its financial rather than constitutional power.

The principal revenue sources for financing federal government activity are given in Figure 14.1 for the fiscal year 2015-16, and the expenditure of these revenues is broken down in Figure 14.2. Further details are accessible in the Department of Finance's 'fiscal reference tables' at <https://www.fin.gc.ca/ftr-trf/2019/ftr-trf-19-eng.asp>. Total revenues for that fiscal year amounted to \$333.2b, and expenditures to \$346.2b.

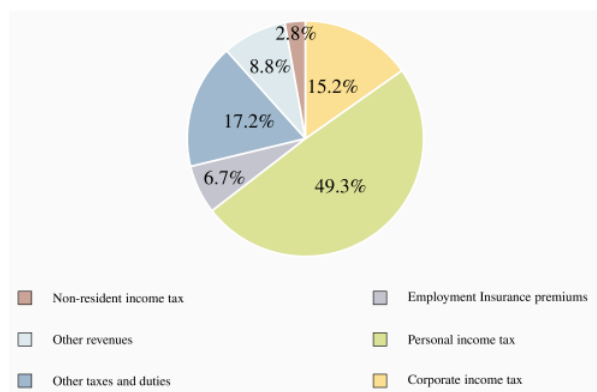


Figure 14.1 Federal Government Revenues 2015–16

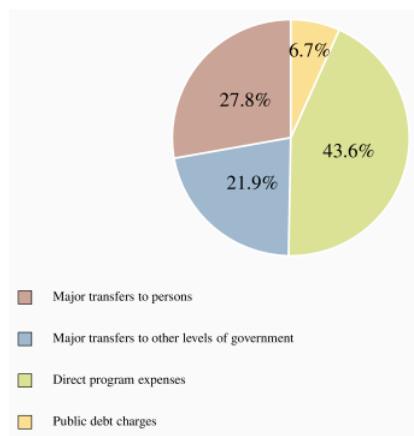


Figure 14.2 Federal expenditures 2015–16

The federal and provincial governments each transfer these revenues to individuals and other levels of government, supply goods and services directly, and also pay interest on accumulated borrowings – the national debt or provincial debt.

Provincial and local governments supply more goods and services than the federal government – health care, drug insurance, education and welfare are the responsibility of provincial and municipal governments. In contrast, national defence, the provision of main traffic arteries, Corrections Canada and a variety of transfer programs to individuals – such as Employment Insurance, Old Age Security and the Canada Pension Plan – are federally funded. The greater part of federal revenues goes towards *transfers* to individuals and provincial governments, as opposed to the supply of goods and services.

This page titled [14.2: Fiscal federalism- Taxing and spending](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14.3: Federal-provincial fiscal relations

The federal government transfers revenue to the provinces using three main programs: Equalization, the Canada Social Transfer and the Canada Health Transfer. Each of these has a different objective. Equalization aims to reduce fiscal disparities among the provinces; The Canada Social Transfer (CST) is for educational and Social Assistance ('welfare') expenditures; The Canada Health Transfer (CHT) performs the same function for health.

Equalization and Territorial Funding

Canada's provinces receive unconditional funding through Canada's Equalization program, whereas the Territories receive federal funding through a separate mechanism - the Territorial Funding Formula.

"Parliament and the Government of Canada are committed to the principle of making equalization payments to ensure that provincial governments have sufficient revenues to provide reasonably comparable levels of public service at reasonably comparable levels of taxation."

This statement, from Section 36(2) of the Constitution Act of 1982, defines the purpose of Equalization. Equalization payments are unconditional – receiving provinces are free to spend the funds on public services according to their own priorities, or even use the revenue to reduce their provincial taxes. Payments are calculated according to a formula that ensures those provinces with revenue-raising ability, or fiscal capacity, below a threshold or 'standard' receive payments from the federal government to bring their capacity up to that standard.

Equalization has gone through very many changes in the several decades of its existence. Its current rules and regulations reflect the 2006 recommendations of a federal Expert Panel. The fiscal capacity of a province is measured by its ability to raise revenues from five major sources: Personal and business income taxes, sales taxes, property taxes, and natural resources. This ability is then compared to the ability of all of the provinces combined to raise revenue; if a difference or shortfall exists, the federal government transfers revenue accordingly, with the amount determined by both the population of the province and the magnitude of its per-person shortfall. Data on annual transfers for Equalization, the Territorial Funding Formula, and the CHT and CST are available at fin.gc.ca/fedprov/mtp-eng.asp

The program transferred \$19.8b to the provinces in 2019-20. The reciprocity status of some provinces varies from year to year. Variation in energy prices and energy-based government revenues are the principal cause of this. British Columbia, Alberta, Saskatchewan and Ontario tend to receive little or zero. Manitoba, Quebec and the Atlantic Provinces have been the major recipient provinces. Quebec receives the largest single amount – approximately two thirds of the total allocation, both on account of its population size and the fact that it has a lower than average fiscal capacity. The Territories received a total of \$3.9b in 2019-20 through the Territorial Funding Formula.

The Canada Social Transfer and the Canada Health Transfer

The CST is a block transfer to provinces in support of post-secondary education, Social Assistance and social services more generally. The CST came into effect in 2004. Prior to that date it was integrated with the health component of federal transfers in a program titled the Canada Health and Social Transfer (CHST). The objective of the separation was to increase the transparency and accountability of federal support for health while continuing to provide funding for other objectives. The CHT is the other part of the unbundled CHST: It provides funding to the provinces for their health expenditures.

The CST and CHT funding comes in two parts: A cash transfer and tax transfer. A tax transfer essentially provides the same support as a cash transfer of equal value, but comes in a different form. In 1977 the federal government agreed with provincial governments to reduce federal personal and corporate tax rates in order to permit the provincial governments to increase the corresponding provincial rates. The net effect was that the federal government got less tax revenue and the provinces got more. And to this day, the federal and provincial governments keep a record of the implied tax transfers that arise from this long-ago agreement. This is the tax transfer component of the CST and the CHT.

The CST support is allocated to provinces and territories on an equal per-capita basis to ensure equal support for all Canadians regardless of their place of residence. The CHT is distributed likewise, and it requires the provinces to abide by the federally-

legislated *Canada Health Act*, which demands that provincial health coverage be comprehensive, universal, portable, accessible and publicly administered.

The CHT transfer amounted to \$40.4b and the CST amounted to \$14.6b, in cash, for the year 2019-20. Health care and health expenditures are a core issue of policy at all levels of government on account of the envisaged growth in these expenditures that will inevitably accompany the aging of the baby-boomers.

This page titled [14.3: Federal-provincial fiscal relations](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14.4: Government-to-individual transfers

Many Canadians take pride in Canada's extensive 'social safety net' that aims to protect individuals from misfortune and the reduction of income in old age. Others believe it is too generous. While it is more supportive than the safety net in the US, the Canadian safety net is no more protective than the nets of the developed economies of the European Union. The extent of such support depends in large measure upon the degree to which governments are willing to impose, and individuals are willing to pay, higher or lower tax rates. The major elements of this umbrella of programs are the following.

The *Canada and Quebec Pension Plans (C/QPP)* are funded from the contributions of workers and their employers. Contributions form 9.9% of an individual's earnings up to the maximum pensionable earnings (MPE) figure of \$57,400 in 2019. The contributions are shared equally by employer and employee. The Canada and Quebec components of the plan operate similarly, but are managed separately. The contribution rate to the QPP stands at 10.65%. Contributions to the plans from workers and their employers are largely transferred immediately to retired workers. Part of the contributions is invested in a fund. The objective of the plans is to ensure that some income is saved for retirement. Many individuals are not very good at planning – they constantly postpone the decision to save, so the state steps in and requires them to save. An individual contributing throughout a full-time working lifecycle can expect an annual pension of about \$14,000 in 2019. The Plans provide a maximum payout of 25% of maximum insurable earnings. The objective is to provide a minimum level of retirement income, not an income that will see individuals live in great comfort.

The C/QPP plans have contributed greatly to the reduction of poverty among the elderly since their introduction in the mid-sixties. The aging of the baby-boom generation – that very large cohort born in the late forties through to the early sixties – means that the percentage of the population in the post-65 age group has begun to increase. To meet this changing demographic, the federal and provincial governments reshaped the plans in the late nineties in order to put them on a sound financial footing – primarily by increasing contributions, that in turn will enable the build-up of a CPP 'fund' that will support the aged in the following decades.

A number of recent studies in Canada on the retirement savings practices of Canadians have proposed that households on average are not saving a sufficient amount for their retirement; many households may thus see a notable decline in their incomes upon retirement. In response to this finding, the federal government agreed with the provinces in June 2016, to add a supplement to the CPP. The new federal provisions, which will be phased in over the period 2019-2025, envisage an increase in contributions that will ultimately lead to a maximum replacement rate of 33% of MPE as opposed to the current goal of 25%. However, full benefits will be experienced only by individuals contributing for their complete lifecycle, meaning that full implementation will take about four decades.

Details of the CPP and the 2016 enhancements are to be found at http://www.fin.gc.ca/n16/data/16-113_3-eng.asp.

Old Age Security (OAS), the *Guaranteed Income Supplement (GIS)* and the *Spousal Allowance (SPA)* together form the second support leg for the retired. OAS is a payment made automatically to individuals once they attain the age of 65. The GIS is an additional payment made only to those on very low incomes – for example, individuals who have little income from their C/QPP or private pension plans. The SPA, which is payable to the spouse or survivor of an OAS recipient, accounts for a small part of the sums disbursed. As of 2019 the maximum annual OAS payment stood at \$7,360. The federal government in 2016 reversed a plan that would have seen the eligible age for receipt of OAS move to 67.

The payments for these plans come from the general tax revenues of the federal government. Unlike the C/QPP, the benefits received are not related to the contributions that an individual makes over the working lifecycle. This program has also had a substantial impact on poverty reduction among the elderly.

Employment Insurance (EI) and *Social Assistance (SA)* are designed to support, respectively, the unemployed and those with no other source of income. Welfare is the common term used to describe SA. Expenditures on EI and SA are strongly cyclical. At the trough of an economic cycle the real value of expenditures on these programs greatly exceeds expenditures at the peak of the cycle. Unemployment in Canada rose above 8% in 2009, and payments to the unemployed and those on welfare reflected this dire state. The strongly cyclical pattern of the cost of these programs reflects the importance of a healthy job market: Macroeconomic conditions have a major impact on social program expenditures.

EI is funded by contributions from employees and their employers. For each dollar contributed by the employee, the employer contributes \$1.4. Premiums are paid on earned income up to a maximum insurable earnings (MIE) of \$53,100 in 2019. The contribution rate for employees stood at 1.62% of MIE in 2017. EI contributions and pay-outs form part of the federal government's general revenues and expenditures. There is no separate 'fund' for this program. However it is expected to operate on

a break-even basis over the longer term. To reflect this, the contribution rate fluctuates with a view to maintaining a balance between payouts and revenues over a seven-year planning period.

EI is called an insurance program, but in reality it is much more than that. Certain groups systematically use the program more than others – those in seasonal jobs, those in rural areas and those in the Atlantic Provinces, for example. Accordingly, using the terminology of Chapter 7, it is not everywhere an actuarially 'fair' insurance program. Benefits payable to unemployed individuals may also depend on their family size, in addition to their work history. While most payments go in the form of 'regular' benefits to unemployed individuals, the EI program also sponsors employee retraining, family benefits that cover maternity and paternity leave, and some other specific target programs for the unemployed.

Social Assistance is provided to individuals who are in serious need of financial support – having no income and few assets. Provincial governments administer SA, although the cost of the program is partly covered by federal transfers through the Canada Social Transfer. The nineteen nineties witnessed a substantial tightening of regulations virtually across the whole of Canada. Access to SA benefits is now more difficult, and benefits have fallen in real terms since the late nineteen eighties.

Welfare dependence peaked in Canada in 1994, when 3.1 million individuals were dependent upon support. As of 2019, the total is approximately half of this, on account of more stringent access conditions, reduced benefit levels and an improved job market. Some groups in Canada believe that benefits should be higher, others believe that making welfare too generous provides young individuals with the wrong incentives in life, and may lead them to neglect schooling and skill development.

Workers Compensation supports workers injured on the job. Worker/employer contributions and general tax revenue form the sources of program revenue, and the mix varies from province-to-province. In contrast to the macro-economy-induced swings in expenditures that characterize SA and EI since the early nineties, expenditures on Worker's Compensation have remained relatively constant.

Canada Child Benefit: The major remaining pillar in Canada's social safety net is the group of payments and tax credits aimed at supporting children: The Canada Child Tax Benefit (CCTB), the Universal Child Care Benefit and the National Child Benefit Supplement were repackaged in 2016 under the title Canada Child Benefit. Child support has evolved and been enriched over the last two decades, partly with the objective of reducing poverty among households with children, and partly with a view to helping parents receiving social assistance to transition back to the labour market. As of 2016, the federal government provides an annual payment to families with children. For each child under the age of 6 the payment is \$6,639 and for each child aged 6-17 the payment is \$5,602. Since these payment are primarily intended for households with low and middle incomes, the amounts are progressively clawed back once the household income reaches a threshold of \$30,000.

Application Box 14.2 Government debts and deficits

Canada's expenditure and tax policies in the nineteen seventies and eighties led to the accumulation of large government debts, as a result of running fiscal deficits. By the mid-nineties the combined federal and provincial debt reached 100% of GDP, with the federal debt accounting for the larger share. This ratio was perilously high: Interest payments absorbed a large fraction of annual government revenues, which in turn limited the ability of the government to embark on new programs or enrich existing ones. Canada's debt rating on international financial markets declined.

In 1995, Finance Minister Paul Martin addressed this problem, and over the following years program spending was pared back. Ultimately, the economy expanded and by the end of the decade the annual deficits at the federal level were eliminated.

As of 2007 the ratio of combined federal and provincial debts stood at just over 60% of GDP. However, the Great Recession of 2008 and following years saw all levels of government experience deficits, with the result that this ratio of combined debt to GDP rose again. Growth in recent years has seen that ratio fall. As of 2018-19, federal interest payments on its debt account for about 7% of its revenues (this figure stood at 28% in the early nineties). At the time of writing, interest rates are low in developed economies and so the interest costs of government debt are low. Low borrowing costs are a reason why some people favor large government spending in the form of infrastructure projects. Those who are fiscally more conservative fear rising rates in the future. The recessionary impacts of the coronavirus pandemic of 2020 will add greatly to accumulated debt, particularly at the federal level.

Debts can be measured in more than a single manner. One measure of debt is the value of all federal government bonds and financial liabilities outstanding. As of 2019-20 this value was approximately \$700b. In addition to this, the federal government has outstanding liabilities to the pensions of its retired employees, and not all of these liabilities have been covered by the contributions of those employees into their pension plans. The federal government also owns assets, both financial and physical - such as office buildings. Hence these assets offset the financial liabilities. To assess the total debt picture of the Canadian economy we need to

add provincial and local government debts to the federal debts, and then consider the annual interest costs of this total. It turns out that the interest costs are just above 2% of GDP.

Source: Government of Canada, Fiscal Reference Tables: <https://www.fin.gc.ca/frt-trf/2019/frt-trf-19-eng.asp>

This page titled [14.4: Government-to-individual transfers](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14.5: Regulation and competition policy

Goals and objectives

The goals of competition policy are relatively uniform across developed economies: The promotion of domestic competition; the development of new ideas, new products and new enterprises; the promotion of efficiency in the resource-allocation sense; the development of manufacturing and service industries that can compete internationally.

In addition to these economic objectives, governments and citizens frown upon monopolies or monopoly practices if they lead to an undue *concentration of political power*. Such power can lead to a concentration of wealth and influence in the hands of an elite.

Canada's regulatory body is the *Competition Bureau*, whose activity is governed primarily by the *Competition Act* of 1986. This act replaced the *Combines Investigation Act*. The *Competition Tribunal* acts as an adjudication body, and is composed of judges and non-judicial members. This tribunal can issue orders on the maintenance of competition in the marketplace. Canada has had anti-combines legislation since 1889, and the act of 1986 is the most recent form of such legislation and policy. The Competition Act does not forbid monopolies, but it does rule as unlawful the *abuse* of monopoly power. Canada's competition legislation is aimed at anti-competitive practices, and a full description of its activities is to be found on its website at www.competitionbureau.gc.ca. Let us examine some of these proscribed policies.

Anti-competitive practices

Anti-competitive practices may either limit entry into a sector of the economy or force existing competitors out. In either case they lead to a reduction in competition.

Mergers may turn competitive firms into a single organization with excessive market power. The customary justification for mergers is that they permit the merged firms to achieve scale economies that would otherwise be impossible. Such scale economies may in turn result in lower prices in the domestic or international market to the benefit of the consumer, but may alternatively reduce competition and result in higher prices. Equally important in this era of global competition is the impact of a merger on a firm's ability to compete internationally.

Mergers can be of the horizontal type (e.g. two manufacturers of pre-mixed concrete merge) or vertical type (a concrete manufacturer merges with a cement manufacturer). In a market with few suppliers mergers have the potential to reduce domestic competition.

Cartels aim to restrict output and thereby increase profits. These formations are almost universally illegal in individual national economies.

While cartels are one means of increasing prices, price discrimination is another, as we saw when studying monopoly behaviour. For example, if a concrete manufacturer makes her product available to large builders at a lower price than to small-scale builders – perhaps because the large builder has more bargaining power – then the small builder is at a competitive disadvantage in the construction business. If the small firm is forced out of the construction business as a consequence, then competition in this sector is reduced.

We introduced the concept of predatory pricing in Chapter 11. Predatory pricing is a practice that is aimed at driving out competition by artificially reducing the price of one product sold by a supplier. For example, a dominant nationwide transporter could reduce price on a particular route where competition comes from a strictly local competitor. By 'subsidizing' this route from profits on other routes, the dominant firm could undercut the local firm and drive it out of the market.

Predatory pricing is a practice that is aimed at driving out competition by artificially reducing the price of one product sold by a supplier.

Suppliers may also refuse to deal. If the local supplier of pre-mixed concrete refuses to sell the product to a local construction firm, then the ability of such a downstream firm to operate and compete may be compromised. This practice is similar to that of exclusive sales and tied sales. An exclusive sale might involve a large vegetable wholesaler forcing her retail clients to buy *only* from this supplier. Such a practice might hurt the local grower of aubergines or zucchini, and also may prevent the retailer from obtaining *some* of her vegetables at a lower price or at a higher quality elsewhere. A tied sale is one where the purchaser must agree to purchase a *bundle* of goods from a supplier.

Refusal to deal: an illegal practice where a supplier refuses to sell to a purchaser.

Exclusive sale: where a retailer is obliged (perhaps illegally) to purchase all wholesale products from a single supplier only.

Tied sale: one where the purchaser must agree to purchase a bundle of goods from a supplier.

Resale price maintenance involves the producer requiring a retailer to sell a product at a specified price. This practice can hurt consumers since they cannot 'shop around'. In Canada, we frequently encounter a 'manufacturer's suggested retail price' for autos and durable goods. But since these prices are not *required*, the practice conforms to the law.

Resale price maintenance is an illegal practice wherein a producer requires sellers to maintain a specified price.

Bid rigging is an illegal practice in which normally competitive bidders conspire to fix the awarding of contracts or sales. For example, two builders, who consider bidding on construction projects, may decide that one will bid seriously for project X and the other will bid seriously on project Y. In this way they conspire to reduce competition in order to make more profit.

Bid rigging is an illegal practice in which bidders (buyers) conspire to set prices in their own interest.

Deception and dishonesty in promoting products can either short-change the consumer or give one supplier an unfair advantage over other suppliers.

Enforcement

The Competition Act is enforced through the Competition Bureau in a variety of ways. Decisions on acceptable business practices are frequently reached through study and letters of agreement between the Bureau and businesses. In some cases, where laws appear to have been violated, criminal proceedings may follow.

Regulation, deregulation and privatization

The last three decades have witnessed a significant degree of privatization and deregulation in Canada, most notably in the transportation, communication and energy sectors. Modern deregulation in the US began with the passage of the *Airline Deregulation Act* of 1978, and was pursued with great energy under the Reagan administration in the eighties. The Economic Council of Canada produced an influential report in 1981, titled "Reforming Regulation," on the impact of regulation and possible deregulation of specific sectors. The Economic Council proposed that regulation in some sectors was inhibiting competition, entry and innovation. As a consequence, the interests of the consumer were in danger of becoming secondary to the interests of the suppliers.

Telecommunications provision, in the era when the telephone was the main form of such communication, was traditionally viewed as a natural monopoly. The Canadian Radio and Telecommunications Commission (CRTC) regulated its rates. The industry has developed dramatically in the last two decades with the introduction of satellite-facilitated communication, the internet, multi-purpose cable networks, cell phones and service integration.

Transportation, in virtually all forms, has been deregulated in Canada since the nineteen eighties. Railways were originally required to subsidize the transportation of grain under the *Crow's Nest Pass* rate structure. But the subsidization of particular markets requires an excessive rate elsewhere, and if the latter markets become subject to competition then a competitive system cannot function. This structure, along with many other anomalies, was changed with the passage of the *Canada Transportation Act* in 1996.

Trucking, historically, has been regulated by individual provinces. Entry was heavily controlled prior to the federal *National Transportation Act* of 1987, and subsequent legislation introduced by a number of provinces, have made for easier entry and a more competitive rate structure.

Deregulation of the airline industry in the US in the late seventies had a considerable influence on thinking and practice in Canada. The Economic Council report of 1981 recommended in favour of easier entry and greater fare competition. These policies were reflected in the 1987 National Transportation Act. Most economists are favourable to deregulation and freedom to enter, and the US experience indicated that cost reductions and increased efficiency could follow. In 1995 an agreement was reached between the US and Canada that provided full freedom for Canadian carriers to move passengers to any US city, and freedom for US carriers to do likewise, subject to a phase-in provision.

The National Energy Board regulates the development and transmission of oil and natural gas. But earlier powers of the Board, involving the regulation of product prices, were eliminated in 1986, and controls on oil exports were also eliminated.

Agriculture remains a highly controlled area of the economy. Supply 'management', which is really supply restriction, and therefore 'price maintenance', characterizes grain, dairy, poultry and other products. Management is primarily through provincial marketing

boards.

The role of the sharing economy

The arrival of universal access to the internet has seen the emergence of what is known as the Sharing economy throughout the world. This expression is used to describe commercial activities that, in the first place, are internet-based. Second, suppliers in the sharing economy use resources in the market place that were initially aimed at a different purpose. *Airbnb* and *Uber* are good examples of companies in sectors of the economy where sharing is possible. In *Uber's* case, the 'ride-share' drivers initially purchased their vehicles for private use, and subsequently redirected them to commercial use. *Airbnb* is a communication corporation that enables the owners of spare home capacity to sell the use of that capacity to short-term renters. With the maturation of such corporations, the concept of 'initial' and 'secondary' use becomes blurred.

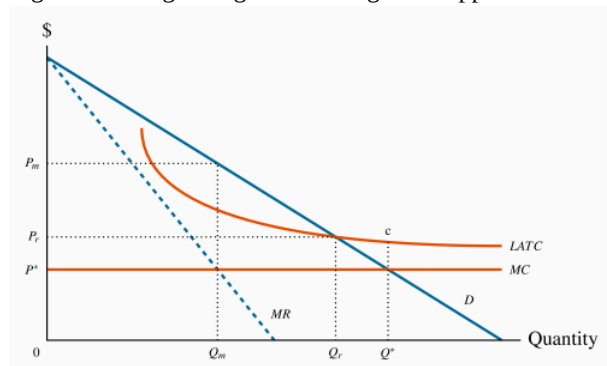
Sharing economy: involves enterprises that are internet based, and that use production resources that have use outside of the marketplace.

The importance of the sharing economy is that it provides an additional source of competition to established suppliers, and therefore limits the market power of the latter. At the same time, the emergence of the sharing economy poses a new set of regulatory challenges: If traditional taxis are required to purchase operating permits (medallions), and the ride-share drivers do not require such permits, is there a reasonable degree of competition in the market, and if not what is the appropriate solution? Should the medallion requirement be abolished, or should ride-share drivers be required to purchase one? In the case of *Airbnb*, the suppliers operate outside of the traditional 'hotel' market. In general they do not charge sales taxes or face any union labour agreements. What is the appropriate response from governments? And how should the sharing economy be taxed?

Price regulation

Regulating monopolistic sectors of the economy is one means of reducing their market power. In Chapter 11 it was proposed that indefinitely decreasing production costs in an industry means that the industry might be considered as a 'natural' monopoly: Higher output can be produced at lower cost with fewer firms. Hence, a single supplier has the potential to supply the market at a lower unit cost; unless, that is, such a single supplier uses his monopoly power. To illustrate how the consumer side may benefit from this production structure through regulation, consider Figure 14.2. For simplicity suppose that long-run marginal costs are constant and that average costs are downward sloping due to an initial fixed cost. The profit-maximizing (monopoly) output is where $MR=MC$ at Q_m and is sold at the price P_m . This output is inefficient because the willingness of buyers to pay for additional units of output exceeds the additional cost. On this criterion the efficient output is Q^* . But $LATC$ exceeds price at Q^* , and therefore it is not feasible for a producer.

Figure 14.2 Regulating a decreasing-cost supplier



The profit-maximizing output is Q_m , where $MR=MC$ and price is P_m . This output is inefficient because marginal benefit is greater than MC . Q^* is the efficient output, but results in losses because $LATC > P$ at that output. A regulated price that covers costs is where $LATC=DQ_r$. This is closer to the efficient output Q^* than the monopoly output Q_m .

One solution is for the regulating body to set a price-quantity combination of P_r , and Q_r , where price equals average cost and therefore generates a normal rate of profit. This output level is still lower than the efficient output level Q^* , but is more efficient than the profit-maximizing output Q_m . It is more efficient in the sense that it is closer to the efficient output Q^* . A problem with such a strategy is that it may induce lax management: If producers are allowed to charge an average-cost price, then there is a reduced incentive for them to keep strict control of their costs in the absence of competition in the marketplace.

A second solution to the declining average cost phenomenon is to implement what is called a two-part tariff. This means that customers pay an 'entry fee' in order to be able to purchase the good. For example, in many jurisdictions hydro or natural gas subscribers may pay a fixed charge per month for their supply line and supply guarantee, and then pay an additional charge that varies with quantity. In this way it is possible for the supplier to charge a price per unit of output that is closer to marginal cost and still make a profit, than under an average cost pricing formula. In terms of Figure 14.2, the total value of entry fees, or fixed components of the pricing, would have to cover the difference between MC and $LATC$ times the output supplied. In Figure 14.2 this implies that if the efficient output Q^* is purchased at a price equal to the MC the producer loses the amount $(c-MC)$ on each unit sold. The access fees would therefore have to cover at least this value.

Such a solution is appropriate when fixed costs are high and marginal costs are low. This situation is particularly relevant in the modern market for telecommunications: The cost to suppliers of marginal access to their networks, whether it be for internet, phone or TV, is negligible compared to the cost of maintaining the network and installing capacity.

Two-part tariff: involves an access fee and a per unit of quantity fee.

Finally, a word of caution: Nobel Laureate George Stigler has argued that there is a danger of regulators becoming too close to the regulated, and that the relationship can evolve to a point where the regulator may protect the regulated firms. In contrast, Professor Philippon of New York University argues that regulators are not regulating sufficiently in the US: they have permitted an excessive number of mergers that have, in turn, reduced competition.

This page titled [14.5: Regulation and competition policy](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14.6: Key Terms

Market failure defines outcomes in which the allocation of resources is not efficient.

Public goods are non-rivalrous, in that they can be consumed simultaneously by more than one individual; additionally they may have a non-excludability characteristic.

Efficient supply of public goods is where the marginal cost equals the sum of individual marginal valuations, and each individual consumes the same quantity.

Asymmetric information is where at least one party in an economic relationship has less than full information and has a different amount of information from another party.

Adverse selection occurs when incomplete or asymmetric information describes an economic relationship.

Moral hazard may characterize behaviour where the costs of certain activities are not incurred by those undertaking them.

Spending power of a federal government arises when the federal government can influence lower level governments due to its financial rather than constitutional power.

Predatory pricing is a practice that is aimed at driving out competition by artificially reducing the price of one product sold by a supplier.

Refusal to deal: an illegal practice where a supplier refuses to sell to a purchaser.

Exclusive sale: where a retailer is obliged (perhaps illegally) to purchase all wholesale products from a single supplier only.

Tied sale: one where the purchaser must agree to purchase a bundle of goods from the one supplier.

Resale price maintenance is an illegal practice wherein a producer requires sellers to maintain a specified price.

Bid rigging is an illegal practice in which bidders (buyers) conspire to set prices in their own interest.

Sharing economy: involves enterprises that are internet based, and that use production resources that have use outside of the marketplace.

Two-part tariff: involves an access fee and a per unit of quantity fee.

This page titled [14.6: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

14.7: Exercises for Chapter 14

EXERCISE 14.1

An economy is composed of two individuals, whose demands for a public good – street lighting – are given by $P=12-(1/2)Q$ and $P=8-(1/3)Q$.

- Graph these demands on a diagram, for values of $Q = 1, \dots, 24$.
- Graph the total demand for this public good by summing the demands vertically, specifying the numerical value of each intercept.
- Let the marginal cost of providing the good be \$5 per unit. Illustrate graphically the efficient supply of the public good (Q^*) in this economy.
- Illustrate graphically the area that represents the total value to the consumers of the amount Q^* .

EXERCISE 14.2

In Exercise 14.1, suppose a new citizen joins the economy, and her demand for the public good is given by $P=10-(5/12)Q$.

- Add this individual's demand curve to the graphic for the above question and graph the new total demand curve, specifying the intercept values.
- Illustrate the area on your graph that represents the new total value to the three citizens of the optimal amount supplied.
- Illustrate graphically the net value to society of the new Q^* – the total value minus the total cost.

EXERCISE 14.3

An industry that is characterized by a decreasing cost structure has a demand curve given by $P=100-Q$ and the marginal revenue curve by $MR=100-2Q$. The marginal cost is $MC=4$, and average cost is $AC=4+188/Q$.

- Graph this cost and demand structure. [*Hint: This graph is similar to Figure 14.2.*]
- Illustrate the efficient output and the monopoly output for the industry.
- Illustrate on the graph the price the monopolist would charge if he were unregulated.

EXERCISE 14.4

Optional: In Question 14.3, suppose the government decides to regulate the behaviour of the supplier, in the interests of the consumer.

- Illustrate graphically the price and output that would emerge if the supplier were regulated so that his allowable price equalled average cost.
- Is this greater or less than the efficient output?
- Compute the AC and P that would be charged with this regulation.
- Illustrate graphically the deadweight loss associated with the regulated price and compare it with the deadweight loss under monopoly.

EXERCISE 14.5

Optional: As an alternative to regulating the supplier such that price covers average total cost, suppose that a two part tariff were used to generate revenue. This scheme involves charging the MC for each unit that is purchased and in addition charging each buyer in the market a fixed cost that is independent of the amount he purchases. If an efficient output is supplied in the market, illustrate graphically the total revenue to be obtained from the component covering a price per unit of the good supplied, and the component covering fixed cost.

This page titled [14.7: Exercises for Chapter 14](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15: International trade

Chapter 15: International trade

In this chapter we will explore:

15.1	Trade in our daily lives
15.2	Canada in the world economy
15.3	Gains from trade: Comparative advantage
15.4	Returns to scale and dynamic gains
15.5	Trade barriers: Tariffs, subsidies and quotas
15.6	The politics of protection
15.7	Institutions governing trade

15.1 Trade in our daily lives

Virtually every economy in the modern world trades with other economies – they are what we call 'open' economies. Evidence of such openness is everywhere evident in our daily life. The world eats Canadian wheat; China exports manufactured goods to almost anywhere we can think of; and Canadians take their holidays in Florida.

As consumers we value the choice and variety of products that trade offers. We benefit from lower prices than would prevail in a world of protectionism. At the same time there is a constant chorus of voices calling for protection from international competition: Manufacturers are threatened by production in Asia; farmers fight against the imports of poultry, beef, and dairy products; even the service sector is concerned about offshore competition from call centres and designers. In this world of competing views it is vital to understand how trade has the potential to improve the well-being of economies.

This chapter examines the theory of international trade, trade flows, and trade policy: Who trades with whom, in what commodities, and why. In general, countries trade with one another because they can buy foreign products at a lower price than it costs to make them at home. International trade reflects specialization and exchange, which in turn improve living standards. It is cost differences between countries rather than technological differences that drive trade: In principle, Canada could supply Toronto with olives and oranges grown in Nunavut greenhouses, but it makes more sense to import them from Greece, Florida or Mexico.

Trade between Canada and other countries differs from trade between provinces. By definition, international trade involves jumping a border, whereas most trade within Canada does not. Internal borders are present in some instances – for example when it comes to recognizing professional qualifications acquired out-of-province. In the second instance, international trade may involve different currencies. When Canadians trade with Europeans the trade is accompanied by financial transactions involving Canadian dollars and Euros. A Canadian buyer of French wine pays in Canadian dollars, but the French vineyard worker is paid in euros. Exchange rates are one factor in determining national competitiveness in international markets. Evidently, not every international trade requires currency trades at the same time – most members of the European Union use the Euro. Indeed a common currency was seen as a means of facilitating trade between member nations of the EU, and thus a means of integrating the constituent economies more effectively.

It is important at the outset to emphasize that while trade has the potential to improve aggregate well being in the trading economies, this does not mean that every citizen will benefit; some will gain others will lose out. Buyers usually gain as a result of having a wider array of products to purchase at lower prices. Successful producers may benefit from production efficiencies associated with accessing a global supply chain, while others may be squeezed by international competition. The former may export more and employ more workers, the latter may contract and lay off employees.

15.2 Canada in the world economy

World trade has grown rapidly since the end of World War II, indicating that trade has become ever more important to national economies. Canada has been no exception. Canada signed the Free Trade Agreement with the US in 1989, and this agreement was expanded in 1994 when Mexico was included under the North America Free Trade Agreement (NAFTA). Imports and exports rose dramatically, from approximately one quarter to forty percent of GDP. Canada is now what is termed a very 'open' economy – one

where trade forms a large fraction of total production. In early 2017, Canada and the EU signed a trade agreement - the Comprehensive Economic and Trade Agreement (CETA). Under this agreement tariffs will be phased out or reduced in most areas of trade over a several-year period. Canada also signed the Comprehensive and Progressive Agreement on Trans-Pacific Trade in 2018 that has the same objective of reducing trade barriers between 11 Pacific-Rim member states.

Smaller economies are typically more open than large economies—Belgium and the Netherlands depend upon trade more than the United States. This is because large economies tend to have a sufficient variety of resources to supply much of an individual country's needs. The European Union is similar, in population terms, to the United States, but it is composed of many distinct economies. Some European economies are equal in size to individual American states. But trade between California and New York is not international, whereas trade between Italy and the Spain is.

Because our economy is increasingly open to international trade, events in the world economy affect our daily lives much more than in the past. The conditions in international markets for basic commodities and energy affect all nations, both importers and exporters. For example, the prices of primary commodities on world markets increased dramatically in the latter part of the 2000s. Higher prices for grains, oil, and fertilizers on world markets brought enormous benefits to Canada, particularly the Western provinces, which produce these commodities. In contrast, by early 2015, many of these prices dropped dramatically and Canadian producers suffered as a consequence.

The service sector accounts for more of our GDP than the manufacturing sector. As incomes grow, the demand for health, education, leisure, financial services, tourism, etc., dominates the demand for physical products. Technically, the income elasticity demand for the former group exceeds the income elasticity of demand for the latter. Internationally, while trade in services is growing rapidly, it still forms a relatively small part of total world trade. Trade in goods—merchandise trade—remains dominant, partly because many countries import unfinished goods, add some value, and re-export them. Even though the value added from such import-export activity may make just a small contribution to GDP, the gross flows of imports and exports can still be large relative to GDP. The transition from agriculture to manufacturing and then to services has been underway in developed economies for over a century. This transition has been facilitated in recent decades by the communications revolution and globalization. Globalization has seen a rapid shift in merchandise production from the developed to the developing world.

Table 15.1 shows the patterns of Canadian merchandise trade in 2018. The US is Canada's major trading partner, buying almost three quarters of our exports and supplying almost two thirds of imports. Table 15.2 details exports and imports by type. Although exports of resource-based products account for only about 40 percent of total merchandise exports, Canada is still viewed as a resource-based economy. This is in part because manufactures account for almost 80 percent of US and European merchandise exports and about 60 percent of Canadian exports. Nevertheless, Canada has important strength in machinery, equipment, and automotive products.

Table 15.1 Canada's Merchandise Trade Patterns 2018

Country	Exports to	Imports from
United States	73.9	64.4
European Union	7.9	10.5
China	5.0	7.6
Mexico	1.6	3.4
Others	11.6	14.1
Total	100	100
Dollar Total	585,255.7	607,205.4

Source: Adapted from Statistics Canada Table 12-10-0011-01

Table 15.2 Canadian Trade by Merchandise Type 2017

Sector	Exports	Imports
Farm, fishing, and intermediate food products	6.5	3.0
Energy products	20.5	5.5

Metal ores and non-metallic minerals	3.8	2.3
Metal and non-metallic mineral products	11.9	7.6
Basic and industrial chemical, plastic and rubber products	6.6	8.5
Forestry products and building and packaging materials	8.5	4.4
Industrial machinery, equipment and parts	5.3	9.7
Electronic and electrical equipment and parts	3.5	11.8
Motor vehicles and parts	16.3	20.0
Aircraft and other transportation equipment and parts	3.7	3.7
Consumer Goods	12.3	21.9
Special transactions trade	1.1	1.6
Total	100	100
Total dollar value in millions	500,892.6	561,425.9

Source: Adapted from Statistics Canada Table 12-10-0002-01

15.3 The gains from trade: Comparative advantage

In the opening chapter of this text we emphasized the importance of opportunity cost and differing efficiencies in the production process as a means of generating benefits to individuals through trade in the marketplace. The simple example we developed illustrated that, where individuals differ in their efficiency levels, benefits can accrue to each individual as a result of specializing and trading. In that example it was assumed that individual A had an absolute advantage in producing one product and that individual Z had an absolute advantage in producing the second good. This set-up could equally well be applied to two economies that have different efficiencies and are considering trade, with the objective of increasing their consumption possibilities. Technically, we could replace Amanda and Zoe with Argentina and Zambia, and nothing in the analysis would have to change in order to illustrate that consumption gains could be attained by both Argentina and Zambia as a result of specialization and trade.

Remember: The opportunity cost of a good is the quantity of another good or service given up in order to have one more unit of the good in question.

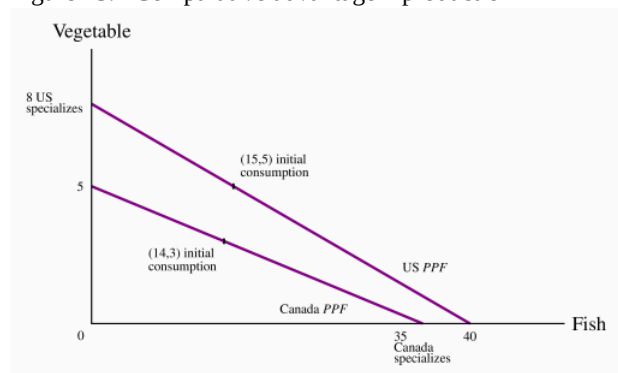
So, let us now consider two economies with differing production capabilities, as illustrated in Figures 15.1 and 15.2. In this instance it is assumed that one economy has an absolute advantage in both goods, but the degree of that advantage is greater in one good than the other. In international trade language, there exists a comparative advantage as well as an absolute advantage. It is frequently a surprise to students that this situation has the capacity to yield consumption advantages to each economy, even though one is absolutely more efficient in producing both of the goods. This is termed the principle of comparative advantage, and it states that even if one country has an absolute advantage in producing both goods, gains to specialization and trade still materialize, provided the opportunity cost of producing the goods differs between economies. This is a remarkable result, and much less intuitive than the principle of absolute advantage. We explore it with the help of the example developed in Figures 15.1 and 15.2.

Principle of comparative advantage states that even if one country has an absolute advantage in producing both goods, gains to specialization and trade still materialize, provided the opportunity cost of producing the goods differs between economies.

We will name these two imaginary economies the US and Canada. Their production possibilities are defined by the *PPFs* in Figure 15.1. Canada can produce 5 units of *V* or 35 units of *F*, or any combination defined by the line joining these points. With the same resources the US can produce 8*V* or 40*F*, or any combination defined by its *PPF*¹. With no trade, Canadians and Americans consume a combination of the goods defined by some point on their respective *PPFs*. The opportunity cost of a unit of *V* in Canada is $\frac{7}{5}$ *F* (the slope of Canada's *PPF* is $\frac{5}{35} = \frac{1}{7}$). In the US the opportunity cost of one unit of *V* is $\frac{5}{8}$ *F* (slope is $\frac{8}{40} = \frac{1}{5}$). In this set-

up the US is more efficient in producing V than F relative to Canada, as reflected by the opportunity costs. Hence we say that the *US has a comparative advantage in the production of V and that Canada has therefore a comparative advantage in producing F .*

Figure 15.1 Comparative advantage – production



Canada specializes completely in Fish at 35, where it has a comparative advantage. Similarly, the US specializes in Vegetable at 8. They trade at a rate of 1:6. The US trades 3V to Canada in return for 18F.

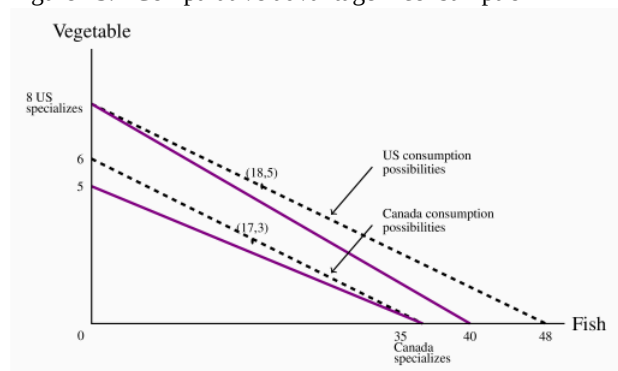
Prior to trade each economy is producing all of the goods it consumes. This no-trade state is termed autarky.

Autarky denotes the no-trade situation.

The gains from trade

We now permit each economy to specialize in producing where it has a comparative advantage. So Canada specializes completely by producing 35F and the US produces 8V. Having done this the economies must now agree on the terms of trade. The terms of trade define the rate at which the two goods will trade post-specialization. Let us suppose that a bargaining process leads to agreement that one unit of V will trade for six units of F . Such a trading rate, one that lies between the opportunity costs of each economy, benefits both economies. This exchange rate lies between Canada's opportunity cost of 1:7 and the US opportunity cost of 1:5. By specializing in F , Canada can now obtain an additional unit of V by sacrificing six units of F , whereas pre-trade it had to sacrifice seven units of F for a unit of V . Technically, by specializing in F and trading at a rate of 1:6 Canada's consumption possibilities have expanded and are given by the consumption possibility frontier (CPF) illustrated in Figure 15.2. The consumption possibility frontier defines what an economy can consume after production specialization and trade.

Figure 15.2 Comparative advantage – consumption



Post specialization the economies trade 1V for 6F. Total production is 35F plus 8V. Hence one consumption possibility would be (18,5) for the US and (17,3) for Canada. Here Canada exchanges 18F in return for 3V.

The US also experiences an improved set of consumption possibilities. By specializing in V and trading at a rate of 1:6 its CPF lies outside its PPF and this enables it to consume more than in the pre-specialization state, where its CPF was defined by its PPF.

Evidently, the US and Canada CPFs are parallel since they trade with each other at the same rate: If Canada exports six units of F for every unit of V that it imports from the US, then the US must import the same six units of F for each unit of V it exports to Canada. The remarkable outcome here is that, even though one economy is more efficient in producing each good, specialization still leads to gains for both economies. The gain is illustrated by the fact that each economy's consumption possibilities lie outside of its production possibilities².

Terms of trade define the rate at which the goods trade internationally.

Consumption possibility frontier defines what an economy can consume after production specialization and trade.

Comparative advantage and factor endowments

A traditional statement of why comparative advantage arises is that economies have different endowments of the factors of production – land, capital and labour endowments differ. A land endowment that facilitates the harvesting of grain (Saskatchewan) or the growing of fruit (California) may be innate to an economy. We say that wheat production is *land intensive*, that aluminum production is *power intensive*, that research and development is *skill intensive*, that auto manufacture is *capital intensive*, that apparel is *labour intensive*. Consequently, if a country is well endowed with some particular factors of production, it is to be expected that it will specialize in producing goods that use those inputs. A relatively abundant supply or endowment of one factor of production tends to make the cost of using that factor relatively cheap: It is relatively less expensive to produce clothing in China and wheat in Canada than the other way around. This explains why Canada's Prairies produce wheat, why Quebec produces aluminum, why Asia produces apparel. But endowments can evolve.

How can we explain why Switzerland specializes in watches, precision instruments, and medical equipment, while Vietnam specializes in rice, tourism and manufactured goods and components? Evidently, Switzerland made a decision to educate its population and invest in the capital required to produce these goods. It was not naturally endowed with these skills, in the same way that Greece is endowed with sun or Saskatchewan is endowed with fertile flat land.

While we have demonstrated the principle of comparative advantage using a two-good example (since we are constrained by the geometry of two dimensions), the conclusions carry over to the case of many goods. Furthermore, the principle has many applications. For example, if one person in the household is more efficient at doing all household chores than another, there are still gains to specialization provided the efficiency differences are not all identical. This is the principle of comparative advantage at work in a microcosm.

Application Box 15.1 The one hundred mile diet

In 2005 two young British Columbians embarked on what has famously become known as the 'one hundred mile diet'—a challenge to eat and drink only products grown within this distance of their home. They succeeded in doing this for a whole year, wrote a book on their experience and went on to produce a TV series. They were convinced that such a project is good for humanity, partly because they wrapped up ideas on organic farming and environmentally friendly practices in the same message.

Reflect now on the implications of this superficially attractive program: If North Americans were to espouse this diet, it would effectively result in the closing down of the midwest of the Continent. From Saskatchewan to Kansas, we are endowed with grain-producing land that is the envy of the planet. But since most of this terrain is not within 100 miles of any big cities, these deluded advocates are proposing that we close up the production of grains and cereals exactly in those locations where such production is extraordinarily efficient. Should we sacrifice grains and cereals completely in this hemisphere, or just cultivate them on a hillside close to home, even if the resulting cultivation were to be more labour and fuel intensive? Should we produce olives in greenhouses in Edmonton rather than importing them from the Mediterranean, or simply stop eating them? Should we sacrifice wine and beer in North Battleford because insufficient grapes and hops are grown locally?

Would production in temperate climates really save more energy than the current practice of shipping vegetables and fruits from a distance—particularly when there are returns to scale associated with their distribution? The 'one hundred mile diet' is based on precepts that are contrary to the norms of the gains from trade. In its extreme the philosophy proposes that food exports be halted and that the world's great natural endowments of land, water, and sun be allowed to lie fallow. Where would that leave a hungry world?

Table 15.1 shows the patterns of Canadian merchandise trade in 2008. The United States was and still is Canada's major trading partner, buying almost three quarters of our exports and supplying almost two thirds of Canadian imports. Table 15.2 details exports by type. Although exports of resource-based products account for only about 40 percent of total exports, Canada is now viewed as a resource-based economy. This is in part because manufactured products account for almost 80 percent of US and European exports but only about 60 percent of Canadian exports. Nevertheless, Canada has important export strength in machinery, equipment, and automotive products.

15.4 Returns to scale and dynamic gains from trade

The theory of comparative advantage explains why economies should wish to trade. The theory is based upon the view that economies are 'inherently' different in their production capabilities. But trade is influenced by more than these differences. We will

explore how returns to scale may be exploited to generate benefits from trade, and also how economies might gain from one-another by learning as a result of trading. This learning can increase domestic productivity.

Returns to scale

One of the reasons Canada signed the North America Free trade Agreement (NAFTA) was that economists convinced the Canadian government that a larger market would enable Canadian producers to be even *more efficient* than in the presence of trade barriers. Rather than opening up trade in order to take advantage of existing comparative advantage, it was proposed that efficiencies would actually increase with market size. This argument is easily understood in terms of increasing returns to scale concepts that we developed in Chapter 8. Essentially, economists suggested that there were several sectors of the Canadian economy that were operating on the downward sloping section of their long-run average cost curve.

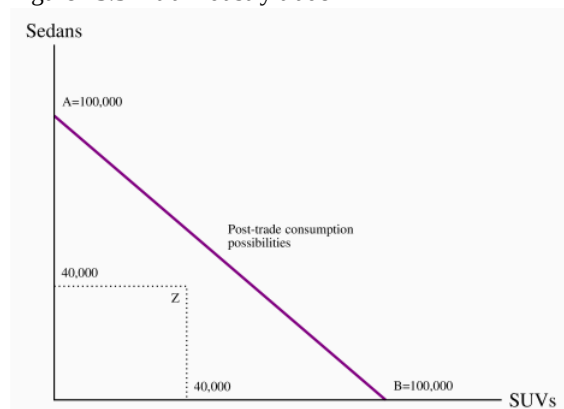
Increasing returns are evident in the world market place as well as the domestic marketplace. Witness the small number of aircraft manufacturers—*Airbus* and *Boeing* are the world's two major manufacturers of large aircraft. Enormous fixed costs—in the form of research, design, and development—or capital outlays frequently result in decreasing unit costs, and the world marketplace can be supplied at a lower cost if some specialization can take place. *Facebook* is the giant in social media. Entertainment streaming companies like *Netflix*, *Amazon*, *Disney* and *Hulu* are few in number because of scale economies. Consider the specific example of automotive trade. In North America, Canadian auto plants produce different vehicle models than their counterparts in the US. Canada exports some *models* of a given manufacturer to the United States and imports other *models*. This is the phenomenon of intra-industry trade and intra-firm trade. How can we explain these patterns?

Intra-industry trade is two-way international trade in products produced within the same industry.

Intra-firm trade is two-way trade in international products produced within the same firm.

In the first instance, intra-industry trade reflects the preference of consumers for a choice of brands; consumers do not all want the same vehicle, or the same software, or the same furnishings. The second element to intra-industry trade is that increasing returns to scale characterize many production processes. Let us see if we can transform the returns to scale ideas developed in earlier chapters into a production possibility framework.

Figure 15.3 Intra industry trade



Hunda can produce either 100,000 of each vehicle or 40,000 of both in each plant. Hence production possibilities are given by the points A, Z, and B. Pre-trade it produces at Z in each economy due to trade barriers. Post-trade it produces at A in one economy and B in the other, and ships the vehicles internationally. Total production increases from 160,000 to 200,000 using the same resources.

Consider the example presented in Figure 15.3, where the hypothetical company *Hunda Motor Corporation* currently has a large assembly plant in *each* of Canada and the US where it produces two types of vehicles; sedans and sports utility vehicles (SUVs). Initially, restrictions on trade in automobiles, in the form of tariffs, between the two countries make it too costly to ship models across the border. Hence Hunda produces both sedans and SUVs in each plant. But for several reasons, *switching between model production is costly and results in reduced output*. Hunda can produce 40,000 vehicles of each type per annum in its plants, but could produce 100,000 of a single model in each plant, using the same amount of capital and labour. This is a situation of increasing returns to scale, and in this instance the scale economies are what make gains from trade possible - as opposed to any innate comparative advantage between the economies. If trade barriers against the shipment of autos across national boundaries can

be eliminated, then Honda can take advantage of scale economies in each plant and increase its total production without using more capital and labour.

As this example implies, an opening up of trade increases the potential market size, and producers who experience increasing returns to scale stand to benefit from an enlarged market because their potential unit costs fall. Returns to scale are not limited to finished goods. Returns to scale characterize the production of many intermediate goods, which are goods used to produce other final goods or services. Manufacturers rarely produce all of the components entering their final products; they have supply chains for components that comprise numerous suppliers. In the automotive industry transmissions, gearboxes and seats are such intermediate goods. If either returns to scale, or comparative advantage, characterize their supply then there are gains to trade in these goods. In the context of the North American Free Trade Agreement (NAFTA), and its most recent form (the US Mexico and Canada Agreement - USMCA), automotive parts as well as automobiles can be shipped free of tariffs across borders provided that they satisfy regional value content (RVC) rules. We observe, for example, transmissions being produced in Ontario and seats in Mexico, and these goods are also shipped freely within North America provided they satisfy the RVC rules. Scale economies characterize transmission manufacture, and comparative advantage characterizes seat production – labour costs are lower in Mexico, and seats are labor intensive.

Content requirements apply to some goods under the NAFTA/USMCA. In the case of vehicles and their components, NAFTA required a 62.5% regional value content - that is, to be imported into one of the NAFTA signatories free of tariffs, 62.5% of the vehicle value had to be attributable to production in one of the three economies. This requirement was designed to prevent the members from using an excessive amount of components from low-wage economies and thereby undermine production of parts and vehicles in North America. The USMCA raises (by the year 2023) the regional value component (RVC) to 75% for most vehicles (70% for heavy trucks), and the RVC to between 65% and 75% on vehicle parts.

Supply chain: denotes the numerous sources for intermediate goods used in producing a final product

Intermediate good: one that is used in the production of final output

Regional value content: requires that a specified percentage of the final value of a product originate in the economies covered in the Agreement.

Dynamic gains from trade

The term dynamic gains denotes the potential for domestic producers to increase productivity as a result of competing with, and learning from, foreign producers.

Dynamic gains: the potential for domestic producers to increase productivity by competing with, and learning from, foreign producers.

Production processes in reality are seldom static. Innovation is constant in the modern world, and innovation is manifested in the form of productivity improvements. An economy's production possibility frontier is determined by its endowments of capital and labour and also the efficiency with which it uses those productive factors. Total factor productivity defines how efficiently the factors of production are combined. Research suggests that in developed economies this productivity increases by about 1% per annum. This means that more output can be produced using the same amounts of capital and labour because production is being carried out more efficiently. In graphical terms, such productivity improvements effectively push out an economy's production possibility frontier by 1% per annum. For economies in the process of development, this productivity growth may be as high as 3% or 4% per annum – for the reason that these economies can observe and learn from economies that are ahead of it technologically.

Freer trade forces domestic firms to compete with foreign firms that may be more productive. Domestic firms that can learn and adapt to competition by becoming more efficient will survive, firms that cannot adapt will not. Inevitably, there will be winners and losers in the production sector of the economy, whereas in the consumption sector most consumers should be winners.

Total factor productivity: a measure of how efficiently the factors of production are combined.

15.5 Trade barriers: Tariffs, subsidies and quotas

Despite the many good arguments favoring free or relatively free trade, we observe numerous trade barriers. These barriers come in several forms. A tariff is a tax on an imported product that is designed to limit trade and generate tax revenue. It is a barrier to trade. An import quota is a limitation on imports; other non-tariff barriers take the form of product content requirements, and subsidies. By raising the domestic price of imports, a tariff helps domestic producers but hurts domestic consumers. Quotas and other non-tariff barriers have similar impacts.

A **tariff** is a tax on an imported product that is designed to limit trade in addition to generating tax revenue.

A **quota** is a quantitative limit on an imported product.

A **trade subsidy** to a domestic manufacturer reduces the domestic cost and limits imports.

Non-tariff barriers, such as product content requirements, limit the gains from trade.

Application Box 15.2 Tariffs – the national policy of J.A. MacDonald

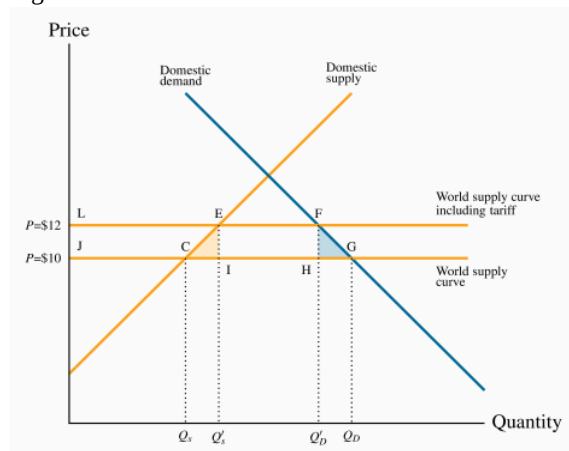
In Canada, tariffs were the main source of government revenues, both before and after Confederation in 1867 and up to World War I. They provided 'incidental protection' for domestic manufacturing. After the 1878 federal election, tariffs were an important part of the National Policy introduced by the government of Sir John A. MacDonald. The broad objective was to create a Canadian nation based on east-west trade and growth.

This National Policy had several dimensions. Initially, to support domestic manufacturing, it increased tariff protection on foreign manufactured goods, but lowered tariffs on raw materials and intermediate goods used in local manufacturing activity. The profitability of domestic manufacturing improved. But on a broader scale, tariff protection, railway promotion, Western settlement, harbour development, and transport subsidies to support the export of Canadian products were intended to support national economic development. Although reciprocity agreements with the United States removed duties on commodities for a time, tariff protection for manufactures was maintained until the GATT negotiations of the post-World War II era.

Tariffs

Figure 15.4 describes how tariffs operate. We can think of this as the Canadian wine market—a market that is heavily taxed in Canada. The world price of Cabernet Sauvignon is, let us say, \$10 per bottle, and this is shown by the horizontal world supply curve at that price. To maintain simplicity, we will neglect any taxes on alcohol here other than the tax represented by the tariff. The international supply curve is horizontal because the domestic market accounts for only a small part of the world demand for wine: we are sufficiently small that international producers can supply us with any amount we wish to buy at the world price. The Canadian demand for this wine is given by the demand curve D , and Canadian suppliers have a supply curve given by S (Canadian Cabernet is assumed to be of the same quality as the imported variety in this example). The effective supply curve in the Canadian market is now BCM . At a price of \$10, Canadian consumers wish to buy Q_D litres, and domestic producers wish to supply Q_S litres. The gap between domestic supply Q_S and domestic demand Q_D is filled by imports. This is the *free trade equilibrium*.

Figure 15.4 Tariffs and trade



At a world price of \$10 the domestic quantity demanded is Q_D . Of this amount Q_S is supplied by domestic producers and the remainder by foreign producers. A tariff increases the world price to \$12. This reduces demand to Q'_D ; the domestic component of supply increases to Q'_S . Of the total loss in consumer surplus (LFGJ), tariff revenue equals EFHI, increased surplus for domestic suppliers equals LECJ, and the deadweight loss is therefore the sum of the triangular areas CEI and HFG.

If the government now imposes a 20 percent tariff on imported wines (or a \$2 per bottle tax), foreign wine sells for \$12 a bottle, inclusive of the tariff. The effective supply curve in the Canadian market becomes BEK . The tariff raises the domestic 'tariff-inclusive' price above the world price, and this shifts the international supply curve of this wine upwards. By raising wine prices in the domestic market, the tariff protects domestic producers by raising the domestic price at which imports become competitive. Those domestic suppliers who were previously not quite competitive at a global price of \$10 are now competitive. The total quantity demanded falls from Q_D to Q'_D at the new equilibrium F . Domestic producers supply the amount Q'_S and imports fall to the

amount $(Q'_D - Q'_S)$. Reduced imports are partly displaced by those domestic producers who can supply at prices between \$10 and \$12. Hence, imports fall both because total consumption falls and because domestic suppliers can displace some imports under the protective tariff; the amount $(Q'_S - Q_S)$.

Since the tariff is a type of tax, its impact in the market depends upon the elasticities of supply and demand, (as illustrated in Chapters 4 and 5). The more elastic is the demand curve, the more a given tariff reduces imports. In contrast, if it is inelastic the quantity of imports declines less.

Costs and benefits of a tariff

The costs of a tariff come from the higher price to consumers, but this is partly offset by the tariff revenue that goes to the government. This tariff revenue is a benefit and can be redistributed to consumers or spent on goods from which consumers derive a benefit. But there are also efficiency costs associated with tariffs—deadweight losses, as we call them. These are the real costs of the tariff, and they arise because the marginal cost of production does not equal the marginal benefit to the consumer. Let us see how these concepts apply with the help of Figure 15.4.

Consumer surplus is the area under the demand curve and above the equilibrium market price. It represents the total amount consumers would have been willing to pay for the product, but did not have to pay, at the equilibrium price. It is a measure of consumer welfare. The tariff raises the market price and reduces this consumer surplus by the amount $LFGJ$. This area measures by how much domestic consumers are worse off as a result of the price increase caused by the tariff. But this is not the net loss for the whole domestic economy, because the government obtains some tax revenue and domestic producers get more revenue and profit.

Government revenue accrues from the domestic sales of imports. On imports of $(Q'_D - Q'_S)$, tax revenue is $EFHI$. Then, domestic producers obtain an additional profit of $LECJ$ —the excess of additional revenue over their cost per additional bottle. If we are not concerned about who gains and who loses, then there is a net loss to the domestic economy equal to the areas CEI and HFG .

The area HFG is the consumer side measure of deadweight loss. At the quantity Q'_D , the production cost of an additional bottle is less than the value placed on it by consumers; and, by not having those additional bottles supplied, consumers forgo a potential gain. The area CEI tells us that when supply by domestic higher-cost producers is increased, and supply of lower-cost foreign producers is reduced, the corresponding resources are not being used efficiently. The sum of the areas CEI and HFG is therefore the total deadweight loss of the tariff.

In the real world we should also be interested in the magnitude of the financial amounts involved here: In particular, how much more do consumers pay with the tariff in place, relative to the additional amounts going to domestic suppliers/corporations? How much tax revenue is generated? How many jobs are created domestically as a result of 'distorting' the market? Regardless of the magnitude of the two deadweight loss areas, which represent the net cost of the tariff, we should be interested in whether the owners of capital gain at the expense of consumers.

Tariffs by country of origin - trade diversion

The imposition of tariffs is governed by the World Trade Organization (WTO). Tariffs are permitted under the WTO rules in specific circumstances: if a particular economy is deemed to be subsidizing exports, and those exports have employment impacts on the destination economy, then a 'retaliatory' tariff may be imposed. A related justification is dumping. It is frequently difficult to prove subsidization or dumping by an exporting economy. An example of such a tariff was one placed by the US on washing machines originating in China in 2016, on the basis of a dumping claim by the United States. The immediate result of this was that the manufacturers located in China switched most of their production to other plants they owned in Vietnam and Thailand. In this particular instance there was virtually no impact on the retail price of washing machines in the US.

Dumping is a predatory practice, based on artificially low costs aimed at driving out domestic producers.

The traditional theory of tariffs described in 15.4 implicitly assumes that production and employment increase in the importing economy as a result of domestic production displacing imported goods. This analysis assumes that the tariff is imposed on a particular commodity, regardless of its economy of origin.

Production subsidies

Figure 15.5 illustrates the effect of a subsidy to a domestic supplier. As in Figure 15.4, the amount Q_D is demanded in the free trade equilibrium and, of this, Q_S is supplied domestically. With a subsidy per unit of output sold, the government can reduce the supply cost of the domestic supplier, thereby shifting the supply curve downward from S to S' . In this illustration, the total quantity demanded remains at Q_D , but the domestic share increases to Q'_S .

price P_{dom} and quantity Q'_D ; the free-trade equilibrium is at G. Of the amount Q'_D , *quota* is supplied by foreign suppliers and the remainder by domestic suppliers. The quota increases the price in the domestic market.

The resulting supply curve yields an equilibrium quantity Q'_D . There are several features to note about this equilibrium. First, the quota pushes the domestic price above the world price (P_{dom} is greater than P) because low-cost international suppliers are partially supplanted by higher-cost domestic suppliers. Second, if the quota is chosen 'appropriately', the same domestic market price could exist under the quota as under the tariff in Figure 15.4. Third, in contrast to the tariff case, the government obtains no tax revenue from the quotas. The higher market price under a quota means that the price per unit received by foreign suppliers is now P_{dom} rather than P . *De facto*, instead of tax revenue being generated in the importing economy, the foreign supplier benefits from a higher price. Fourth, inefficiencies are associated with the equilibrium at Q'_D . These inefficiencies arise because the lower-cost international suppliers are not permitted to supply the amount they would be willing to supply at the quota-induced market equilibrium. In other words, more efficient producers are being squeezed out of the market by quotas that make space for less-efficient producers.

Application Box 15.3 Cheese quota in Canada

In 1978 the federal government set a cheese import quota for Canada at just over 20,000 tonnes. This quota was implemented initially to protect the interests of domestic suppliers. Despite a strong growth in population and income in the intervening decades, the import quota has remained unchanged. The result is a price for cheese that is considerably higher than it would otherwise be. The quotas are owned by individuals and companies who have the right to import cheese. The quotas are also traded among importers, at a price. Importers wishing to import cheese beyond their available quota pay a tariff of about 250 percent. So, while the consumer is the undoubted loser in this game, who gains?

First the suppliers gain, as illustrated in Figure 15.6. Canadian consumers are required to pay high-cost domestic producers who displace lower-cost producers from overseas. Second, the holders of the quotas gain. With the increase in demand for cheese that comes with higher incomes, the domestic cheese price increases over time and this in turn makes an individual quota more valuable.

In the 2018 United States Mexico Canada Agreement, a slight increase in access to Canadian markets was granted in return for a corresponding increase in access to the US market.

15.6 The politics of protection

Objections to imports are frequent and come from many different sectors of the economy. In the face of the gains from trade which we have illustrated in this chapter, why do we observe such strong opposition to imported goods and services?

Structural change and technology

In a nutshell the answer is that, while consumers in the aggregate gain from the reduction of trade barriers, and there is a net gain to the economy at large, some *individual sectors of the economy lose out*. Not surprisingly the sectors that will be adversely affected are vociferous in lodging their objections. Sectors of the economy that cannot compete with overseas suppliers generally see a reduction in jobs. This has been the case in the manufacturing sector of the Canadian and US economies in recent decades, as manufacturing and assembly has flown off-shore to Asia and Mexico where labour costs are lower. Domestic job losses are painful, and frequently workers who have spent decades in a particular job find reemployment difficult, and rarely get as high a wage as in their displaced job.

Such job losses are reflected in calls for tariffs on imports from China, for example, in order to 'level the playing field' – that is, to counter the impact of lower wages in China. Of course it is precisely because of lower labour costs in China that the Canadian consumer benefits.

In Canada we deal with such dislocation first by providing unemployment payments to workers, and by furnishing retraining allowances, both coming from Canada's Employment Insurance program. While such support does not guarantee an equally good alternative job, structural changes in the economy, due to both internal and external developments, must be confronted. For example, the information technology revolution made tens of thousands of 'data entry' workers redundant. Should producers have shunned the technological developments which increased their productivity dramatically? If they did, would they be able to compete in world markets?

While job losses feature heavily in protests against technological development and freer trade, most modern economies continue to grow and create more jobs in the service sector than are lost in the manufacturing sector. Developed economies now have many

more workers in service than manufacture. Service jobs are not just composed of low-wage jobs in fast food establishments – 'Mcjobs', they are high paying jobs in the health, education, legal, financial and communications sectors of the economy.

Successful lobbying and concentration

While efforts to protect manufacture have not resulted in significant barriers to imports of manufactures, objections in some specific sectors of the economy seem to be effective worldwide. One sector that stands out is agriculture, where political conditions are conducive to the continuance of protection and what is called 'supply management' – domestic production quotas. The reason for 'successful' supply limitation appears to rest in the geographic concentration of potential beneficiaries of such protection and the scattered beneficiaries of freer trade on the one hand, and the costs and benefits of political organization on the other: Farmers tend to be concentrated in a limited number of rural electoral ridings and hence they can collectively have a major impact on electoral outcomes. Second, the benefits that accrue to trade restriction are heavily concentrated in the economy – keep in mind that about two percent of the population lives on farms, or relies on farming for its income. By contrast the costs on a per person scale are small, and are spread over the whole population. Thus, in terms of the costs of political organization, the incentives for consumers are small, but the incentives for producers are high.

In addition to the differing patterns of costs and benefits, rural communities tend to be more successful in pushing trade restrictions based on a 'way-of-life' argument. By permitting imports that might displace local supply, lobbyists are frequently successful in convincing politicians that long-standing way-of-life traditions would be endangered, even if such 'traditions' are accompanied by monopolies and exceptionally high tariffs.

Valid trade barriers: Infant industries and dumping?

An argument that carries both intellectual and emotional appeal to voters is the 'infant industry' argument. It goes as follows: New ventures and sectors of the economy may require time before that can compete internationally. Scale economies may be involved, for example, and time may be required for producers to expand their scale of operation, at which time costs will have fallen to international (i.e. competitive) levels. In addition, learning-by-doing may be critical in more high-tech sectors and, once again, with the passage of time costs should decline for this reason also.

The problem with this stance is that these 'infants' have insufficient incentive to 'grow up' and become competitive. A protection measure that is initially intended to be temporary can become permanent because of the potential job losses associated with a cessation of the protection to an industry that fails to become internationally competitive. Furthermore, employees and managers in protected sectors have insufficient incentive to make their production competitive if they realize that their government will always be there to protect them.

In contrast to the infant industry argument, economists are more favourable to restrictions that are aimed at preventing dumping.

Dumping may occur either because foreign suppliers choose to sell at artificially low prices (prices below their break-even price for example), or because of surpluses in foreign markets resulting from oversupply. For example, if, as a result of price support in its own market, a foreign government induced oversupply in butter and it chose to sell such butter on world markets at a price well below the going ('competitive') world supply price, such a sale would constitute dumping. Alternatively, an established foreign supplier might choose to enter our domestic market by selling its products at artificially low prices, with a view to driving domestic competition out of the domestic market. Having driven out the domestic competition it would then be in a position to raise prices. This is predatory pricing as explored in the last chapter. Such behaviour differs from a permanently lower price on the part of foreign suppliers. This latter may be welcomed as a gain from trade, whereas the former may generate no gains and serve only to displace domestic labour and capital.

Protectionism in the age of pandemics

The year 2020 will be remembered in history as the year of the coronavirus pandemic. An uncountable number of men and women died all across the globe as a result of contracting COVID-19, the respiratory disorder brought on by an attack of the coronavirus. In the absence of a vaccine, health authorities the world over implemented a twin policy of social distancing and quarantining (or self-isolation). The world economy went into a tailspin, as huge fractions of the labor force were laid off. Trade patterns were disrupted and serious shortages of personal protection equipment (PPE - masks, visors, gowns), ventilators and drugs emerged. The world demand for PPE and ventilators skyrocketed. But the production of PPE was concentrated in China; most western economies did not have the necessary productive capacity to supply even non-pandemic requirements. Bidding wars erupted amongst countries and hospitals as they vied for supply, while domestic producers of some products added to their production capacity.

Following this chaos, we ask if self-sufficiency would not be a better model than open trade. Would a world where each country ensured it had the production capacity to produce these necessities in times of emergency not be superior to one where global supply chains characterize everything from computers to generic drugs? India is a major producer of generic drugs and the components for such drugs. The demand for anti-biotics and pain killers also rocketed upwards with the pandemic.

There is more than one way to plan for a pandemic, and such planning should not involve a generalized move to self-sufficiency on the part of the global economy. One strategy is to build up inventories of PPE and ventilators domestically. This is costly, but for the most part feasible. It does not represent a complete solution because technology changes will make 30-year old ventilators sitting in inventory redundant for the next pandemic. In addition, most medications have a limited shelf life. Hence one solution is to maintain and rotate substantial inventories of emergency equipment using existing supply chains, and benefit from the efficiencies that are built into these chains.

A second option is to maintain excess production capacity on the part of domestic manufacturers of critical pandemic products. Maintaining such capacity should be considered at least partially as a social cost; pandemics ravage societies, not just individuals, and therefore society should undertake part of the cost of insuring against them.

A more general argument against global trade comes in the form of protecting food supplies. In the early 2000s an increase in global cereal prices led some economies to limit exports of specific crops on account of the fact that global demand was pushing prices to a level that low-income consumers could not afford. But such a policy may threaten consumers in other low-income economies whose demands have not changed in a context of reduced supplies. The reality is that world food supply is adequate for world consumption, even in the presence of disruptions. It is also the case that certain economies have huge advantages in producing specific kinds of food. For example, Canada, the US and the Ukraine produce cereals very economically. Mountainous regions are unsuitable for this production. It would benefit no economy for these economies to lower their production of grains to the point where they produced only enough for their own production. By the same reasoning, warmer climates produce fruits, coffee beans, olives etc that cannot easily be produced in many regions suited to wheat. The gains to specialization in the world economy are enormous. Where food shortages occur we frequently encounter the scourges of drought or war or political upheaval, and these conditions inhibit the distribution of foodstuffs.

What about supply chains? If motherboards produced in China are not being exported in sufficient quantities then indeed production of computers in North America will suffer. But to infer from this that North America should decide to produce all of its computer components in North America is illogical. First, in the time of a pandemic, if certain economies in the supply chain are on lockdown, we cannot be sure that the domestic economy would not be on lockdown simultaneously. Second, the cost to moving the production of all computer parts to North America would likely double the cost of computer hardware - including cell-phones. Perhaps a disruption to our supply chains is something we need to bear in extraordinary times. In case it requires emphasis, most producers in supply chains have incentives to produce and sell. If they do not they will die economically.

The energy sector of every economy is impacted with the outbreak of a pandemic. This is because the demand for fuel (primarily oil) declines following policies of social distancing, limits on permissible travel, and the closure of some production facilities that depend upon oil. In North America, as we saw in Chapter 4 earlier, the price of oil declined from US \$60 per barrel to US \$20 in the space of two months in early 2020. Since production costs are higher in both Canada and much of the US than in Saudi Arabia, the North Sea and Russia, producers in North America were squeezed. Many were no longer able to cover their full production costs, and forced to cease drilling and recovering oil. Inevitably there was a clamor for protection. Producers sought tariffs on competing oil: Tariffs would increase the price of cheaper-to-produce foreign oil and enable domestic producers to survive.

While protection might seem like a 'sensible' policy in this instance, the fact is that unilateral tariffs usually invite reprisals, and raise the danger of a trade war with ever-expanding counter protectionism. In contrast to the case of a *shortage* of medical supplies, the energy sector in Canada suffered from a *glut* of world oil supply. The domestic issue is not about the health of consumers (as in the case of medical supplies), it is about the health of producers.

To conclude: a pandemic is a profoundly serious event and such events inflict major costs on all societies. There are no magic bullets in the form of low-cost ideal economic policies to counter viral warfare. The key to policy making is to recognize constraints and recognize an attack as soon as possible. A wholesale move to insulate the domestic economy is ill-conceived. Comparative advantage confers enormous benefits to all nations. Specific policies should take the form of inventory management and excess production capacity in specific sectors of the economy.

15.7 Institutions governing trade

In the nineteenth century, world trade grew rapidly, in part because the leading trading nation at the time—the United Kingdom—pursued a vigorous policy of free trade. In contrast, US tariffs averaged about 50 percent, although they had fallen to around 30 percent by the early 1920s. As the industrial economies went into the Great Depression of the late 1920s and 1930s, there was pressure to protect domestic jobs by keeping out imports. Tariffs in the United States returned to around 50 percent, and the United Kingdom abandoned the policy of free trade that had been pursued for nearly a century. The combination of world recession and increasing tariffs led to a disastrous slump in the volume of world trade, further exacerbated by World War II.

The WTO and GATT

After World War II, there was a collective determination to see world trade restored. Bodies such as the International Monetary Fund and the World Bank were set up, and many countries signed the General Agreement on Tariffs and Trade (GATT), a commitment to reduce tariffs progressively and dismantle trade restrictions.

Under successive rounds of GATT, tariffs fell steadily. By 1960, United States tariffs were only one-fifth of their level at the outbreak of the War. In the United Kingdom, the system of wartime quotas on imports had been dismantled by the mid-1950s, after which tariffs were reduced by nearly half in the ensuing 25 years. Europe as a whole moved toward an enlarged European Union in which tariffs between member countries have been abolished. By the late 1980s, Canada's tariffs had been reduced to about one-quarter of their immediate post-World War II level.

The GATT Secretariat, now called the World Trade Organization (WTO), aims both to dismantle existing protection that reduces efficiency and to extend trade liberalization to more and more countries. Tariff levels throughout the world are now as low as they have ever been, and trade liberalization has been an engine of growth for many economies. The consequence has been a substantial growth in world trade.

NAFTA, the USMCA, the EU, the CETA, and the TPP

In North America, policy since the 1980s has led to a free trade area that covers the flow of trade between Canada, the United States, and Mexico. The Canada/United States free trade agreement (FTA) of 1989 expanded in 1994 to include Mexico in the North American Free Trade Agreement (NAFTA). The objective in both cases was to institute freer trade between these countries in most goods and services. This meant the elimination or reduction of tariffs and non-tariff barriers over a period of years, with a few exceptions in specific products and cultural industries. A critical component of the Agreement was the establishment of a dispute-resolution mechanism, under which disputes would be resolved by a panel of 'judges' nominated from the member economies. Evidence of the success of these agreements is reflected in the fact that Canadian exports have grown to more than 30 percent of GDP, and trade with the United States accounts for the lion's share of Canadian trade flows. NAFTA was updated and replaced in 2018 and the new agreement is termed the United States Mexico Canada Agreement.

The European Union was formed after World War II, with the prime objective of bringing about a greater degree of *political* integration in Europe. Two world wars had laid waste to their economies and social fabric. Closer economic ties and greater trade were seen as the means of achieving this integration. The Union was called the "Common Market" for much of its existence. The Union originally had six member states, and as of 2019 the number is 28, with several other candidate countries in the process of application, most notably Turkey. The European Union (EU) has a secretariat and parliament in Bruxelles. The UK intends to exit the EU as of late 2019.

Canada has concluded a free trade agreement with the European Union that is termed the Comprehensive Economic and Trade Agreement (CETA). It has the objective of implementing free trade between the two negotiating parties, though there remain some exceptions, for example agriculture.

The Comprehensive and Progressive Agreement for Trans-Pacific Partnership (CPTPP) is a trading agreement between Canada and ten other Pacific-Rim economies that came into being in 2018. Negotiations for a Trans Pacific Partnership treaty were complete by 2016. Those negotiations involved 12 Pacific Rim economies including Canada and the United States, but excluding China. The Obama presidency appeared ready to sign the treaty, however the Trump presidency (and also the Democratic candidate for president of the US, Hillary Clinton) decided that the Partnership was not in the interests of the United States and withdrew its affiliation. The remaining 11 economies reached an agreement to implement the partnership in December 2018.

Key Terms

Autarky denotes the no-trade situation.

Principle of comparative advantage states that even if one country has an absolute advantage in producing both goods, gains to specialization and trade still materialize, provided the opportunity cost of producing the goods differs between economies.

Terms of trade define the rate at which goods trade internationally.

Consumption possibility frontier defines what an economy can consume after production specialization and trade.

Intra-industry trade is two-way international trade in products produced within the same industry.

Intra-firm trade is two-way trade in international products produced within the same firm.

Supply chain: denotes the numerous sources for intermediate goods used in producing a final product.

Intermediate good: one that is used in the production of final output.

Content requirement: requires that a specified percentage of the final value of a product originate in the producing economy.

Dynamic gains: the potential for domestic producers to increase productivity by competing with, and learning from, foreign producers.

Total factor productivity: how efficiently the factors of production are combined.

Tariff is a tax on an imported product that is designed to limit trade in addition to generating tax revenue. It is a barrier to trade.

Quota is a quantitative limit on an imported product.

Trade subsidy to a domestic manufacturer reduces the domestic cost and limits imports.

Non-tariff barriers, such as product content requirements, limits the gains from trade.

Dumping is a predatory practice, based on artificial costs aimed at driving out domestic producers.

Exercises for Chapter 15

EXERCISE 15.1

The following table shows the labour input requirements to produce a bushel of wheat and a litre of wine in two countries, Northland and Southland, on the assumption of constant cost production technology – meaning that the production possibility curves in each are straight lines. You can answer this question either by analyzing the table or developing a graph similar to Figure 15.1, assuming each economy has 4 units of labour.

Labour requirements per unit produced		
	Northland	Southland
Per bushel of wheat	1	3
Per litre of wine	2	4

1. Which country has an absolute advantage in the production of both wheat and wine?
2. What is the opportunity cost of wheat in each economy? Of wine?
3. What is the pattern of comparative advantage here?
4. Suppose the country with a comparative advantage in wine reduces wheat production by one bushel and reallocates the labour involved to wine production. How much additional wine does it produce?

EXERCISE 15.2

Canada and the United States can produce two goods, xylophones and yogurt. Each good can be produced with labour alone. Canada requires 60 hours to produce a ton of yogurt and 6 hours to produce a xylophone. The United States requires 40 hours to produce the ton of yogurt and 5 hours to produce a xylophone.

1. Describe the state of absolute advantage between these economies in producing goods.
2. In which good does Canada have a comparative advantage? Does this mean the United States has a comparative advantage in the other good?

3. Draw the production possibility frontier for each economy to scale on a diagram, assuming that each economy has an endowment of 240 hours of labour, and that the PPFs are linear.
4. On the same diagram, draw Canada's consumption possibility frontier on the assumption that it can trade with the United States at the United States' rate of transformation.
5. Draw the US consumption possibility frontier under the assumption that it can trade at Canada's rate of transformation.

EXERCISE 15.3

The domestic demand for bicycles is given by $P=36-0.3Q$. The foreign supply is given by $P=18$ and domestic supply by $P=16+0.4Q$.

1. Illustrate the market equilibrium on a diagram, and illustrate the amounts supplied by domestic and foreign suppliers in equilibrium.
2. If the government now imposes a tariff of \$6 per unit on the foreign good, illustrate the impact geometrically.
3. In the diagram, illustrate the area representing tariff revenue.
4. *Optional:* Compute the price and quantity in equilibrium with free trade, and again in the presence of the tariff.

EXERCISE 15.4

1. In Exercise 15.3, illustrate graphically the deadweight losses associated with the imposition of the tariff.
2. Illustrate on your diagram the additional amount of profit made by the domestic producer as a result of the tariff. [*Hint:* Refer to Figure 15.4 in the text.]

EXERCISE 15.5

The domestic demand for office printers is given by $P=40-0.2Q$. The supply of domestic producers is given by $P=12+0.1Q$, and international supply by $P=20$.

1. Illustrate this market geometrically.
2. If the government gives a production subsidy of \$2 per unit to domestic suppliers in order to increase their competitiveness, illustrate the impact of this on the domestic supply curve.
3. Illustrate geometrically the cost to the government of this scheme.

EXERCISE 15.6

Consider the data underlying Figure 15.1. Suppose, from the initial state of comparative advantage, where Canada specializes in fish and the US in vegetable, we have a technological change in fishing. The US invents the multi-hook fishing line, and as a result can now produce 64 units of fish with the same amount of labour, rather than the 40 units it could produce before the technological change. This technology does not spread to Canada however.

1. Illustrate the new *PPF* for the US in addition to the *PPF* for Canada.
2. What is the new opportunity cost (number of fish) associated with one unit of *V*?
3. Has comparative advantage changed here – which economy should specialize in the production of each good?

EXERCISE 15.7

The following are hypothetical (straight line) production possibilities tables for Canada and the United States. For each line required, plot any two or more points on the line.

Canada					United States				
	A	B	C	D		A	B	C	D
Peaches	0	5	10	15	Peaches	0	10	20	30
Apples	30	20	10	0	Apples	15	10	5	0

1. Plot Canada's production possibilities curve.
2. Plot the United States' production possibilities curve.

3. What is each country's cost ratio of producing peaches and apples?
 4. Which economy should specialize in which product?
 5. Plot the United States' trading possibilities curve (by plotting at least 2 points on the curve) if the actual terms of the trade are 1 apple for 1 peach.
 6. Plot the Canada' trading possibilities curve (by plotting at least 2 points on the curve) if the actual terms of the trade are 1 apple for 1 peach.
 7. Suppose that the optimum product mixes before specialization and trade were B in the United States and C in Canada. What are the gains from specialization and trade?
1. Note that we are considering the *PPFs* to be straight lines rather than concave shapes. The result we illustrate here carries over to that case also, but it is simpler to illustrate with the linear *PPFs*.
 2. To illustrate the gains numerically, let Canada import $3V$ from the US in return for exporting $18F$. Note that this is a trading rate of 1:6. Hence, Canada consumes $3V$ and $17F$ (Canada produced $35F$ and exported $18F$, leaving it with $17F$). It follows that the US consumes $5V$, having exported $3V$ of the $8V$ it produced, and obtained in return $18F$ in imports. The new consumption bundles are illustrated in the figure: $(17,3)$ for Canada and $(18,5)$ for the US. These consumption bundles clearly represent an improvement over the autarky situation. For example, if Canada had wished to consume $3V$ pre-trade, it would only have been able to consume $14F$, whereas with trade it can consume $17F$.

This page titled [15: International trade](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.1: Trade in our daily lives

Virtually every economy in the modern world trades with other economies – they are what we call 'open' economies. Evidence of such openness is everywhere evident in our daily life. The world eats Canadian wheat; China exports manufactured goods to almost anywhere we can think of; and Canadians take their holidays in Florida.

As consumers we value the choice and variety of products that trade offers. We benefit from lower prices than would prevail in a world of protectionism. At the same time there is a constant chorus of voices calling for protection from international competition: Manufacturers are threatened by production in Asia; farmers fight against the imports of poultry, beef, and dairy products; even the service sector is concerned about offshore competition from call centres and designers. In this world of competing views it is vital to understand how trade has the potential to improve the well-being of economies.

This chapter examines the theory of international trade, trade flows, and trade policy: Who trades with whom, in what commodities, and why. In general, countries trade with one another because they can buy foreign products at a lower price than it costs to make them at home. International trade reflects specialization and exchange, which in turn improve living standards. It is cost differences between countries rather than technological differences that drive trade: In principle, Canada could supply Toronto with olives and oranges grown in Nunavut greenhouses, but it makes more sense to import them from Greece, Florida or Mexico.

Trade between Canada and other countries differs from trade between provinces. By definition, international trade involves jumping a border, whereas most trade within Canada does not. Internal borders are present in some instances – for example when it comes to recognizing professional qualifications acquired out-of-province. In the second instance, international trade may involve different currencies. When Canadians trade with Europeans the trade is accompanied by financial transactions involving Canadian dollars and Euros. A Canadian buyer of French wine pays in Canadian dollars, but the French vineyard worker is paid in euros. Exchange rates are one factor in determining national competitiveness in international markets. Evidently, not every international trade requires currency trades at the same time – most members of the European Union use the Euro. Indeed a common currency was seen as a means of facilitating trade between member nations of the EU, and thus a means of integrating the constituent economies more effectively.

It is important at the outset to emphasize that while trade has the potential to improve aggregate well being in the trading economies, this does not mean that every citizen will benefit; some will gain others will lose out. Buyers usually gain as a result of having a wider array of products to purchase at lower prices. Successful producers may benefit from production efficiencies associated with accessing a global supply chain, while others may be squeezed by international competition. The former may export more and employ more workers, the latter may contract and lay off employees.

This page titled [15.1: Trade in our daily lives](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.2: Canada in the world economy

World trade has grown rapidly since the end of World War II, indicating that trade has become ever more important to national economies. Canada has been no exception. Canada signed the Free Trade Agreement with the US in 1989, and this agreement was expanded in 1994 when Mexico was included under the North America Free Trade Agreement (NAFTA). Imports and exports rose dramatically, from approximately one quarter to forty percent of GDP. Canada is now what is termed a very 'open' economy – one where trade forms a large fraction of total production. In early 2017, Canada and the EU signed a trade agreement - the Comprehensive Economic and Trade Agreement (CETA). Under this agreement tariffs will be phased out or reduced in most areas of trade over a several-year period. Canada also signed the Comprehensive and Progressive Agreement on Trans-Pacific Trade in 2018 that has the same objective of reducing trade barriers between 11 Pacific-Rim member states.

Smaller economies are typically more open than large economies—Belgium and the Netherlands depend upon trade more than the United States. This is because large economies tend to have a sufficient variety of resources to supply much of an individual country's needs. The European Union is similar, in population terms, to the United States, but it is composed of many distinct economies. Some European economies are equal in size to individual American states. But trade between California and New York is not international, whereas trade between Italy and the Spain is.

Because our economy is increasingly open to international trade, events in the world economy affect our daily lives much more than in the past. The conditions in international markets for basic commodities and energy affect all nations, both importers and exporters. For example, the prices of primary commodities on world markets increased dramatically in the latter part of the 2000s. Higher prices for grains, oil, and fertilizers on world markets brought enormous benefits to Canada, particularly the Western provinces, which produce these commodities. In contrast, by early 2015, many of these prices dropped dramatically and Canadian producers suffered as a consequence.

The service sector accounts for more of our GDP than the manufacturing sector. As incomes grow, the demand for health, education, leisure, financial services, tourism, etc., dominates the demand for physical products. Technically, the income elasticity demand for the former group exceeds the income elasticity of demand for the latter. Internationally, while trade in services is growing rapidly, it still forms a relatively small part of total world trade. Trade in goods—merchandise trade—remains dominant, partly because many countries import unfinished goods, add some value, and re-export them. Even though the value added from such import-export activity may make just a small contribution to GDP, the gross flows of imports and exports can still be large relative to GDP. The transition from agriculture to manufacturing and then to services has been underway in developed economies for over a century. This transition has been facilitated in recent decades by the communications revolution and globalization. Globalization has seen a rapid shift in merchandise production from the developed to the developing world.

Table 15.1 shows the patterns of Canadian merchandise trade in 2018. The US is Canada's major trading partner, buying almost three quarters of our exports and supplying almost two thirds of imports. Table 15.2 details exports and imports by type. Although exports of resource-based products account for only about 40 percent of total merchandise exports, Canada is still viewed as a resource-based economy. This is in part because manufactures account for almost 80 percent of US and European merchandise exports and about 60 percent of Canadian exports. Nevertheless, Canada has important strength in machinery, equipment, and automotive products.

Table 15.1 Canada's Merchandise Trade Patterns 2018

Country	Exports to	Imports from
United States	73.9	64.4
European Union	7.9	10.5
China	5.0	7.6
Mexico	1.6	3.4
Others	11.6	14.1
Total	100	100
Dollar Total	585,255.7	607,205.4

Source: Adapted from Statistics Canada Table 12-10-0011-01

Table 15.2 Canadian Trade by Merchandise Type 2017

Sector	Exports	Imports
Farm, fishing, and intermediate food products	6.5	3.0
Energy products	20.5	5.5
Metal ores and non-metallic minerals	3.8	2.3
Metal and non-metallic mineral products	11.9	7.6
Basic and industrial chemical, plastic and rubber products	6.6	8.5
Forestry products and building and packaging materials	8.5	4.4
Industrial machinery, equipment and parts	5.3	9.7
Electronic and electrical equipment and parts	3.5	11.8
Motor vehicles and parts	16.3	20.0
Aircraft and other transportation equipment and parts	3.7	3.7
Consumer Goods	12.3	21.9
Special transactions trade	1.1	1.6
Total	100	100
Total dollar value in millions	500,892.6	561,425.9

Source: Adapted from Statistics Canada Table 12-10-0002-01

This page titled [15.2: Canada in the world economy](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.3: The gains from trade- Comparative advantage

In the opening chapter of this text we emphasized the importance of opportunity cost and differing efficiencies in the production process as a means of generating benefits to individuals through trade in the marketplace. The simple example we developed illustrated that, where individuals differ in their efficiency levels, benefits can accrue to each individual as a result of specializing and trading. In that example it was assumed that individual A had an absolute advantage in producing one product and that individual Z had an absolute advantage in producing the second good. This set-up could equally well be applied to two economies that have different efficiencies and are considering trade, with the objective of increasing their consumption possibilities. Technically, we could replace Amanda and Zoe with Argentina and Zambia, and nothing in the analysis would have to change in order to illustrate that consumption gains could be attained by both Argentina and Zambia as a result of specialization and trade.

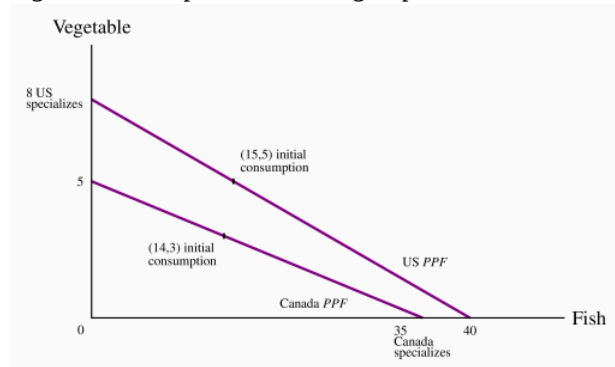
Remember: The opportunity cost of a good is the quantity of another good or service given up in order to have one more unit of the good in question.

So, let us now consider two economies with differing production capabilities, as illustrated in Figures 15.1 and 15.2. In this instance it is assumed that one economy has an absolute advantage in both goods, but the degree of that advantage is greater in one good than the other. In international trade language, there exists a comparative advantage as well as an absolute advantage. It is frequently a surprise to students that this situation has the capacity to yield consumption advantages to each economy, even though one is absolutely more efficient in producing both of the goods. This is termed the principle of comparative advantage, and it states that even if one country has an absolute advantage in producing both goods, gains to specialization and trade still materialize, provided the opportunity cost of producing the goods differs between economies. This is a remarkable result, and much less intuitive than the principle of absolute advantage. We explore it with the help of the example developed in Figures 15.1 and 15.2.

Principle of comparative advantage states that even if one country has an absolute advantage in producing both goods, gains to specialization and trade still materialize, provided the opportunity cost of producing the goods differs between economies.

We will name these two imaginary economies the US and Canada. Their production possibilities are defined by the *PPFs* in Figure 15.1. Canada can produce 5 units of *V* or 35 units of *F*, or any combination defined by the line joining these points. With the same resources the US can produce 8*V* or 40*F*, or any combination defined by its *PPF*¹. With no trade, Canadians and Americans consume a combination of the goods defined by some point on their respective *PPFs*. The opportunity cost of a unit of *V* in Canada is 7*F* (the slope of Canada's *PPF* is $5/35=1/7$). In the US the opportunity cost of one unit of *V* is 5*F* (slope is $8/40=1/5$). In this set-up the US is more efficient in producing *V* than *F* relative to Canada, as reflected by the opportunity costs. Hence we say that the *US has a comparative advantage in the production of V* and that *Canada has therefore a comparative advantage in producing F*.

Figure 15.1 Comparative advantage – production



Canada specializes completely in Fish at 35, where it has a comparative advantage. Similarly, the US specializes in Vegetable at 8. They trade at a rate of 1:6. The US trades 3*V* to Canada in return for 18*F*.

Prior to trade each economy is producing all of the goods it consumes. This no-trade state is termed autarky.

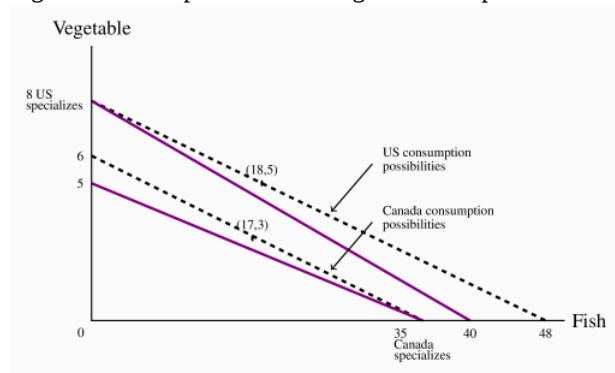
Autarky denotes the no-trade situation.

The gains from trade

We now permit each economy to specialize in producing where it has a comparative advantage. So Canada specializes completely by producing 35*F* and the US produces 8*V*. Having done this the economies must now agree on the terms of trade. The terms of

trade define the rate at which the two goods will trade post-specialization. Let us suppose that a bargaining process leads to agreement that one unit of V will trade for six units of F . Such a trading rate, one that lies between the opportunity costs of each economy, benefits both economies. This exchange rate lies between Canada's opportunity cost of 1:7 and the US opportunity cost of 1:5. By specializing in F , Canada can now obtain an additional unit of V by sacrificing six units of F , whereas pre-trade it had to sacrifice seven units of F for a unit of V . Technically, by specializing in F and trading at a rate of 1:6 Canada's consumption possibilities have expanded and are given by the consumption possibility frontier (CPF) illustrated in Figure 15.2. The consumption possibility frontier defines what an economy can consume after production specialization and trade.

Figure 15.2 Comparative advantage – consumption



Post specialization the economies trade 1V for 6F. Total production is 35F plus 8V. Hence one consumption possibility would be (18,5) for the US and (17,3) for Canada. Here Canada exchanges 18F in return for 3V.

The US also experiences an improved set of consumption possibilities. By specializing in V and trading at a rate of 1:6 its CPF lies outside its PPF and this enables it to consume more than in the pre-specialization state, where its CPF was defined by its PPF .

Evidently, the US and Canada CPF s are parallel since they trade with each other at the same rate: If Canada exports six units of F for every unit of V that it imports from the US, then the US must import the same six units of F for each unit of V it exports to Canada. The remarkable outcome here is that, even though one economy is more efficient in producing each good, specialization still leads to gains for both economies. The gain is illustrated by the fact that each economy's consumption possibilities lie outside of its production possibilities².

Terms of trade define the rate at which the goods trade internationally.

Consumption possibility frontier defines what an economy can consume after production specialization and trade.

Comparative advantage and factor endowments

A traditional statement of why comparative advantage arises is that economies have different endowments of the factors of production – land, capital and labour endowments differ. A land endowment that facilitates the harvesting of grain (Saskatchewan) or the growing of fruit (California) may be innate to an economy. We say that wheat production is *land intensive*, that aluminum production is *power intensive*, that research and development is *skill intensive*, that auto manufacture is *capital intensive*, that apparel is *labour intensive*. Consequently, if a country is well endowed with some particular factors of production, it is to be expected that it will specialize in producing goods that use those inputs. A relatively abundant supply or endowment of one factor of production tends to make the cost of using that factor relatively cheap: It is relatively less expensive to produce clothing in China and wheat in Canada than the other way around. This explains why Canada's Prairies produce wheat, why Quebec produces aluminum, why Asia produces apparel. But endowments can evolve.

How can we explain why Switzerland specializes in watches, precision instruments, and medical equipment, while Vietnam specializes in rice, tourism and manufactured goods and components? Evidently, Switzerland made a decision to educate its population and invest in the capital required to produce these goods. It was not naturally endowed with these skills, in the same way that Greece is endowed with sun or Saskatchewan is endowed with fertile flat land.

While we have demonstrated the principle of comparative advantage using a two-good example (since we are constrained by the geometry of two dimensions), the conclusions carry over to the case of many goods. Furthermore, the principle has many applications. For example, if one person in the household is more efficient at doing all household chores than another, there are still gains to specialization provided the efficiency differences are not all identical. This is the principle of comparative advantage at work in a microcosm.

Application Box 15.1 The one hundred mile diet

In 2005 two young British Columbians embarked on what has famously become known as the 'one hundred mile diet'—a challenge to eat and drink only products grown within this distance of their home. They succeeded in doing this for a whole year, wrote a book on their experience and went on to produce a TV series. They were convinced that such a project is good for humanity, partly because they wrapped up ideas on organic farming and environmentally friendly practices in the same message.

Reflect now on the implications of this superficially attractive program: If North Americans were to espouse this diet, it would effectively result in the closing down of the midwest of the Continent. From Saskatchewan to Kansas, we are endowed with grain-producing land that is the envy of the planet. But since most of this terrain is not within 100 miles of any big cities, these deluded advocates are proposing that we close up the production of grains and cereals exactly in those locations where such production is extraordinarily efficient. Should we sacrifice grains and cereals completely in this hemisphere, or just cultivate them on a hillside close to home, even if the resulting cultivation were to be more labour and fuel intensive? Should we produce olives in greenhouses in Edmonton rather than importing them from the Mediterranean, or simply stop eating them? Should we sacrifice wine and beer in North Battleford because insufficient grapes and hops are grown locally?

Would production in temperate climates really save more energy than the current practice of shipping vegetables and fruits from a distance—particularly when there are returns to scale associated with their distribution? The 'one hundred mile diet' is based on precepts that are contrary to the norms of the gains from trade. In its extreme the philosophy proposes that food exports be halted and that the world's great natural endowments of land, water, and sun be allowed to lie fallow. Where would that leave a hungry world?

Table 15.1 shows the patterns of Canadian merchandise trade in 2008. The United States was and still is Canada's major trading partner, buying almost three quarters of our exports and supplying almost two thirds of Canadian imports. Table 15.2 details exports by type. Although exports of resource-based products account for only about 40 percent of total exports, Canada is now viewed as a resource-based economy. This is in part because manufactured products account for almost 80 percent of US and European exports but only about 60 percent of Canadian exports. Nevertheless, Canada has important export strength in machinery, equipment, and automotive products.

This page titled [15.3: The gains from trade- Comparative advantage](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.4: Returns to scale and dynamic gains from trade

The theory of comparative advantage explains why economies should wish to trade. The theory is based upon the view that economies are 'inherently' different in their production capabilities. But trade is influenced by more than these differences. We will explore how returns to scale may be exploited to generate benefits from trade, and also how economies might gain from one-another by learning as a result of trading. This learning can increase domestic productivity.

Returns to scale

One of the reasons Canada signed the North America Free trade Agreement (NAFTA) was that economists convinced the Canadian government that a larger market would enable Canadian producers to be even *more efficient* than in the presence of trade barriers. Rather than opening up trade in order to take advantage of existing comparative advantage, it was proposed that efficiencies would actually increase with market size. This argument is easily understood in terms of increasing returns to scale concepts that we developed in Chapter 8. Essentially, economists suggested that there were several sectors of the Canadian economy that were operating on the downward sloping section of their long-run average cost curve.

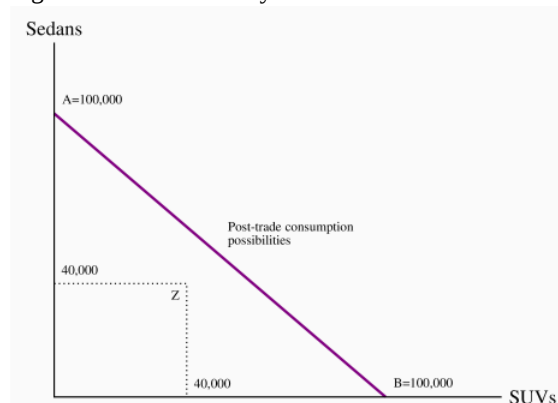
Increasing returns are evident in the world market place as well as the domestic marketplace. Witness the small number of aircraft manufacturers—*Airbus* and *Boeing* are the world's two major manufacturers of large aircraft. Enormous fixed costs—in the form of research, design, and development—or capital outlays frequently result in decreasing unit costs, and the world marketplace can be supplied at a lower cost if some specialization can take place. *Facebook* is the giant in social media. Entertainment streaming companies like *Netflix*, *Amazon*, *Disney* and *Hulu* are few in number because of scale economies. Consider the specific example of automotive trade. In North America, Canadian auto plants produce different vehicle models than their counterparts in the US. Canada exports some *models* of a given manufacturer to the United States and imports other *models*. This is the phenomenon of intra-industry trade and intra-firm trade. How can we explain these patterns?

Intra-industry trade is two-way international trade in products produced within the same industry.

Intra-firm trade is two-way trade in international products produced within the same firm.

In the first instance, intra-industry trade reflects the preference of consumers for a choice of brands; consumers do not all want the same vehicle, or the same software, or the same furnishings. The second element to intra-industry trade is that increasing returns to scale characterize many production processes. Let us see if we can transform the returns to scale ideas developed in earlier chapters into a production possibility framework.

Figure 15.3 Intra industry trade



Hunda can produce either 100,000 of each vehicle or 40,000 of both in each plant. Hence production possibilities are given by the points A, Z, and B. Pre-trade it produces at Z in each economy due to trade barriers. Post-trade it produces at A in one economy and B in the other, and ships the vehicles internationally. Total production increases from 160,000 to 200,000 using the same resources.

Consider the example presented in Figure 15.3, where the hypothetical company *Hunda Motor Corporation* currently has a large assembly plant in *each* of Canada and the US where it produces two types of vehicles; sedans and sports utility vehicles (SUVs). Initially, restrictions on trade in automobiles, in the form of tariffs, between the two countries make it too costly to ship models across the border. Hence Hunda produces both sedans and SUVs in each plant. But for several reasons, *switching between model production is costly and results in reduced output*. Hunda can produce 40,000 vehicles of each type per annum in its plants, but

could produce 100,000 of a single model in each plant, using the same amount of capital and labour. This is a situation of increasing returns to scale, and in this instance the scale economies are what make gains from trade possible - as opposed to any innate comparative advantage between the economies. If trade barriers against the shipment of autos across national boundaries can be eliminated, then Honda can take advantage of scale economies in each plant and increase its total production without using more capital and labour.

As this example implies, an opening up of trade increases the potential market size, and producers who experience increasing returns to scale stand to benefit from an enlarged market because their potential unit costs fall. Returns to scale are not limited to finished goods. Returns to scale characterize the production of many intermediate goods, which are goods used to produce other final goods or services. Manufacturers rarely produce all of the components entering their final products; they have supply chains for components that comprise numerous suppliers. In the automotive industry transmissions, gearboxes and seats are such intermediate goods. If either returns to scale, or comparative advantage, characterize their supply then there are gains to trade in these goods. In the context of the North American Free Trade Agreement (NAFTA), and its most recent form (the US Mexico and Canada Agreement - USMCA), automotive parts as well as automobiles can be shipped free of tariffs across borders provided that they satisfy regional value content (RVC) rules. We observe, for example, transmissions being produced in Ontario and seats in Mexico, and these goods are also shipped freely within North America provided they satisfy the RVC rules. Scale economies characterize transmission manufacture, and comparative advantage characterizes seat production – labour costs are lower in Mexico, and seats are labor intensive.

Content requirements apply to some goods under the NAFTA/USMCA. In the case of vehicles and their components, NAFTA required a 62.5% regional value content - that is, to be imported into one of the NAFTA signatories free of tariffs, 62.5% of the vehicle value had to be attributable to production in one of the three economies. This requirement was designed to prevent the members from using an excessive amount of components from low-wage economies and thereby undermine production of parts and vehicles in North America. The USMCA raises (by the year 2023) the regional value component (RVC) to 75% for most vehicles (70% for heavy trucks), and the RVC to between 65% and 75% on vehicle parts.

Supply chain: denotes the numerous sources for intermediate goods used in producing a final product

Intermediate good: one that is used in the production of final output

Regional value content: requires that a specified percentage of the final value of a product originate in the economies covered in the Agreement.

Dynamic gains from trade

The term dynamic gains denotes the potential for domestic producers to increase productivity as a result of competing with, and learning from, foreign producers.

Dynamic gains: the potential for domestic producers to increase productivity by competing with, and learning from, foreign producers.

Production processes in reality are seldom static. Innovation is constant in the modern world, and innovation is manifested in the form of productivity improvements. An economy's production possibility frontier is determined by its endowments of capital and labour and also the efficiency with which it uses those productive factors. Total factor productivity defines how efficiently the factors of production are combined. Research suggests that in developed economies this productivity increases by about 1% per annum. This means that more output can be produced using the same amounts of capital and labour because production is being carried out more efficiently. In graphical terms, such productivity improvements effectively push out an economy's production possibility frontier by 1% per annum. For economies in the process of development, this productivity growth may be as high as 3% or 4% per annum – for the reason that these economies can observe and learn from economies that are ahead of it technologically.

Freer trade forces domestic firms to compete with foreign firms that may be more productive. Domestic firms that can learn and adapt to competition by becoming more efficient will survive, firms that cannot adapt will not. Inevitably, there will be winners and losers in the production sector of the economy, whereas in the consumption sector most consumers should be winners.

Total factor productivity: a measure of how efficiently the factors of production are combined.

This page titled [15.4: Returns to scale and dynamic gains from trade](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.5: Trade barriers- Tariffs, subsidies and quotas

Despite the many good arguments favoring free or relatively free trade, we observe numerous trade barriers. These barriers come in several forms. A tariff is a tax on an imported product that is designed to limit trade and generate tax revenue. It is a barrier to trade. An import quota is a limitation on imports; other non-tariff barriers take the form of product content requirements, and subsidies. By raising the domestic price of imports, a tariff helps domestic producers but hurts domestic consumers. Quotas and other non-tariff barriers have similar impacts.

A **tariff** is a tax on an imported product that is designed to limit trade in addition to generating tax revenue.

A **quota** is a quantitative limit on an imported product.

A **trade subsidy** to a domestic manufacturer reduces the domestic cost and limits imports.

Non-tariff barriers, such as product content requirements, limit the gains from trade.

Application Box 15.2 Tariffs – the national policy of J.A. MacDonald

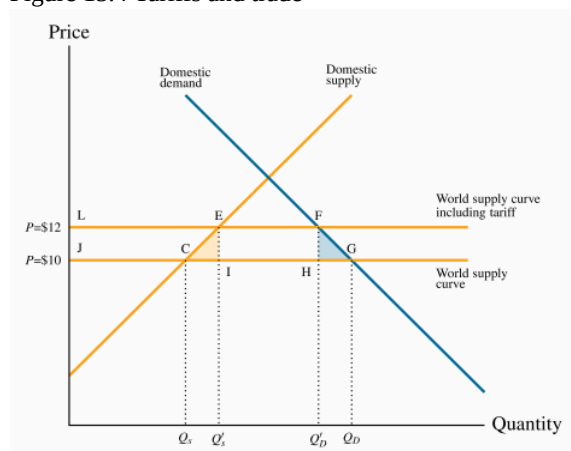
In Canada, tariffs were the main source of government revenues, both before and after Confederation in 1867 and up to World War I. They provided 'incidental protection' for domestic manufacturing. After the 1878 federal election, tariffs were an important part of the National Policy introduced by the government of Sir John A. MacDonald. The broad objective was to create a Canadian nation based on east-west trade and growth.

This National Policy had several dimensions. Initially, to support domestic manufacturing, it increased tariff protection on foreign manufactured goods, but lowered tariffs on raw materials and intermediate goods used in local manufacturing activity. The profitability of domestic manufacturing improved. But on a broader scale, tariff protection, railway promotion, Western settlement, harbour development, and transport subsidies to support the export of Canadian products were intended to support national economic development. Although reciprocity agreements with the United States removed duties on commodities for a time, tariff protection for manufactures was maintained until the GATT negotiations of the post-World War II era.

Tariffs

Figure 15.4 describes how tariffs operate. We can think of this as the Canadian wine market—a market that is heavily taxed in Canada. The world price of Cabernet Sauvignon is, let us say, \$10 per bottle, and this is shown by the horizontal world supply curve at that price. To maintain simplicity, we will neglect any taxes on alcohol here other than the tax represented by the tariff. The international supply curve is horizontal because the domestic market accounts for only a small part of the world demand for wine: we are sufficiently small that international producers can supply us with any amount we wish to buy at the world price. The Canadian demand for this wine is given by the demand curve D , and Canadian suppliers have a supply curve given by S (Canadian Cabernet is assumed to be of the same quality as the imported variety in this example). The effective supply curve in the Canadian market is now BCM . At a price of \$10, Canadian consumers wish to buy Q_D litres, and domestic producers wish to supply Q_S litres. The gap between domestic supply Q_S and domestic demand Q_D is filled by imports. This is the *free trade equilibrium*.

Figure 15.4 Tariffs and trade



At a world price of \$10 the domestic quantity demanded is Q_D . Of this amount Q_S is supplied by domestic producers and the remainder by foreign producers. A tariff increases the world price to \$12. This reduces demand to Q'_D ; the domestic component of

supply increases to Q'_s . Of the total loss in consumer surplus (LFGJ), tariff revenue equals EFHI, increased surplus for domestic suppliers equals LECJ, and the deadweight loss is therefore the sum of the triangular areas CEI and HFG.

If the government now imposes a 20 percent tariff on imported wines (or a \$2 per bottle tax), foreign wine sells for \$12 a bottle, inclusive of the tariff. The effective supply curve in the Canadian market becomes *BEK*. The tariff raises the domestic 'tariff-inclusive' price above the world price, and this shifts the international supply curve of this wine upwards. By raising wine prices in the domestic market, the tariff protects domestic producers by raising the domestic price at which imports become competitive. Those domestic suppliers who were previously not quite competitive at a global price of \$10 are now competitive. The total quantity demanded falls from Q_D to Q'_D at the new equilibrium *F*. Domestic producers supply the amount Q'_s and imports fall to the amount $(Q'_D - Q'_s)$. Reduced imports are partly displaced by those domestic producers who can supply at prices between \$10 and \$12. Hence, imports fall both because total consumption falls and because domestic suppliers can displace some imports under the protective tariff; the amount $(Q'_s - Q_s)$.

Since the tariff is a type of tax, its impact in the market depends upon the elasticities of supply and demand, (as illustrated in Chapters 4 and 5). The more elastic is the demand curve, the more a given tariff reduces imports. In contrast, if it is inelastic the quantity of imports declines less.

Costs and benefits of a tariff

The costs of a tariff come from the higher price to consumers, but this is partly offset by the tariff revenue that goes to the government. This tariff revenue is a benefit and can be redistributed to consumers or spent on goods from which consumers derive a benefit. But there are also efficiency costs associated with tariffs—deadweight losses, as we call them. These are the real costs of the tariff, and they arise because the marginal cost of production does not equal the marginal benefit to the consumer. Let us see how these concepts apply with the help of Figure 15.4.

Consumer surplus is the area under the demand curve and above the equilibrium market price. It represents the total amount consumers would have been willing to pay for the product, but did not have to pay, at the equilibrium price. It is a measure of consumer welfare. The tariff raises the market price and reduces this consumer surplus by the amount *LFGJ*. This area measures by how much domestic consumers are worse off as a result of the price increase caused by the tariff. But this is not the net loss for the whole domestic economy, because the government obtains some tax revenue and domestic producers get more revenue and profit.

Government revenue accrues from the domestic sales of imports. On imports of $(Q'_D - Q'_s)$, tax revenue is *EFHI*. Then, domestic producers obtain an additional profit of *LECJ*—the excess of additional revenue over their cost per additional bottle. If we are not concerned about who gains and who loses, then there is a net loss to the domestic economy equal to the areas *CEI* and *HFG*.

The area *HFG* is the consumer side measure of deadweight loss. At the quantity Q'_D , the production cost of an additional bottle is less than the value placed on it by consumers; and, by not having those additional bottles supplied, consumers forgo a potential gain. The area *CEI* tells us that when supply by domestic higher-cost producers is increased, and supply of lower-cost foreign producers is reduced, the corresponding resources are not being used efficiently. The sum of the areas *CEI* and *HFG* is therefore the total deadweight loss of the tariff.

In the real world we should also be interested in the magnitude of the financial amounts involved here: In particular, how much more do consumers pay with the tariff in place, relative to the additional amounts going to domestic suppliers/corporations? How much tax revenue is generated? How many jobs are created domestically as a result of 'distorting' the market? Regardless of the magnitude of the two deadweight loss areas, which represent the net cost of the tariff, we should be interested in whether the owners of capital gain at the expense of consumers.

Tariffs by country of origin - trade diversion

The imposition of tariffs is governed by the World Trade Organization (WTO). Tariffs are permitted under the WTO rules in specific circumstances: if a particular economy is deemed to be subsidizing exports, and those exports have employment impacts on the destination economy, then a 'retaliatory' tariff may be imposed. A related justification is dumping. It is frequently difficult to prove subsidization or dumping by an exporting economy. An example of such a tariff was one placed by the US on washing machines originating in China in 2016, on the basis of a dumping claim by the United States. The immediate result of this was that the manufacturers located in China switched most of their production to other plants they owned in Vietnam and Thailand. In this particular instance there was virtually no impact on the retail price of washing machines in the US.

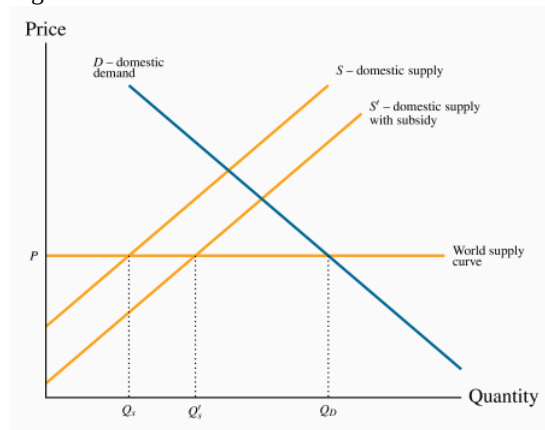
Dumping is a predatory practice, based on artificially low costs aimed at driving out domestic producers.

The traditional theory of tariffs described in 15.4 implicitly assumes that production and employment increase in the importing economy as a result of domestic production displacing imported goods. This analysis assumes that the tariff is imposed on a particular commodity, regardless of its economy of origin.

Production subsidies

Figure 15.5 illustrates the effect of a subsidy to a domestic supplier. As in Figure 15.4, the amount Q_D is demanded in the free trade equilibrium and, of this, Q_S is supplied domestically. With a subsidy per unit of output sold, the government can reduce the supply cost of the domestic supplier, thereby shifting the supply curve downward from S to S' . In this illustration, the total quantity demanded remains at Q_D , but the domestic share increases to Q'_S .

Figure 15.5 Subsidies and trade



With a world supply price of P , a domestic supply curve S , and a domestic demand D , the amount Q_D is purchased. Of this, Q_S is supplied domestically and $(Q_D - Q_S)$ by foreign suppliers. A per-unit subsidy to domestic suppliers shifts their supply curve to S' , and increases their market share to Q'_S .

The new equilibrium represents a misallocation of resources. When domestic output increases from Q_S to Q'_S , a low-cost international producer is being replaced by a higher cost domestic supplier; the domestic supply curve S lies above the international supply curve P in this range of output.

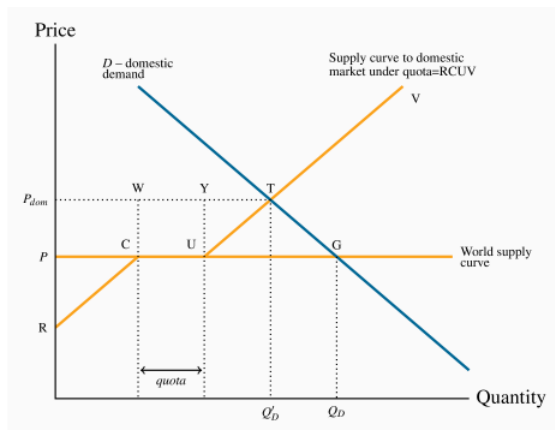
Note that this example deals with a subsidy to domestic suppliers who are selling in the domestic market. It is not a subsidy to domestic producers who are selling in the international market – an export subsidy.

This subsidy comes with a cost to the domestic economy: Taxpayers-at-large must pay higher taxes to support this policy; and each dollar raised in tax itself has a deadweight loss, as we examined in Chapter 5.

Quotas

A quota is a limit placed upon the amount of a good that can be imported. Consider Figure 15.6, where again there is a domestic supply curve coupled with a world price of P . Rather than imposing a tariff, the government imposes a quota that restricts imports to a physical amount denoted by the distance *quota* on the quantity axis. The supply curve facing domestic consumers then has several segments to it. First it has the segment RC , reflecting the fact that domestic suppliers are competitive with world suppliers up to the amount C . Beyond this output, world suppliers can supply at a price of P , whereas domestic suppliers cannot compete at this price. Therefore the supply curve becomes horizontal, but only *up to the amount permitted under the quota*—the quantity CU corresponding to *quota*. Beyond this amount, international supply is not permitted and therefore additional amounts are supplied by the (higher cost) domestic suppliers. Hence the supply curve to domestic buyers becomes the supply curve from the domestic suppliers once again.

Figure 15.6 Quotas and trade



At the world price P , plus a *quota*, the supply curve becomes RCUV. This has three segments: (i) domestic suppliers who can supply below P ; (ii) *quota*; and (iii) domestic suppliers who can only supply at a price above P . The quota equilibrium is at T, with price P_{dom} and quantity Q_D ; the free-trade equilibrium is at G. Of the amount Q_D , *quota* is supplied by foreign suppliers and the remainder by domestic suppliers. The quota increases the price in the domestic market.

The resulting supply curve yields an equilibrium quantity Q_D . There are several features to note about this equilibrium. First, the quota pushes the domestic price above the world price (P_{dom} is greater than P) because low-cost international suppliers are partially supplanted by higher-cost domestic suppliers. Second, if the quota is chosen 'appropriately', the same domestic market price could exist under the quota as under the tariff in Figure 15.4. Third, in contrast to the tariff case, the government obtains no tax revenue from the quotas. The higher market price under a quota means that the price per unit received by foreign suppliers is now P_{dom} rather than P . *De facto*, instead of tax revenue being generated in the importing economy, the foreign supplier benefits from a higher price. Fourth, inefficiencies are associated with the equilibrium at Q_D . These inefficiencies arise because the lower-cost international suppliers are not permitted to supply the amount they would be willing to supply at the quota-induced market equilibrium. In other words, more efficient producers are being squeezed out of the market by quotas that make space for less-efficient producers.

Application Box 15.3 Cheese quota in Canada

In 1978 the federal government set a cheese import quota for Canada at just over 20,000 tonnes. This quota was implemented initially to protect the interests of domestic suppliers. Despite a strong growth in population and income in the intervening decades, the import quota has remained unchanged. The result is a price for cheese that is considerably higher than it would otherwise be. The quotas are owned by individuals and companies who have the right to import cheese. The quotas are also traded among importers, at a price. Importers wishing to import cheese beyond their available quota pay a tariff of about 250 percent. So, while the consumer is the undoubted loser in this game, who gains?

First the suppliers gain, as illustrated in Figure 15.6. Canadian consumers are required to pay high-cost domestic producers who displace lower-cost producers from overseas. Second, the holders of the quotas gain. With the increase in demand for cheese that comes with higher incomes, the domestic cheese price increases over time and this in turn makes an individual quota more valuable.

In the 2018 United States Mexico Canada Agreement, a slight increase in access to Canadian markets was granted in return for a corresponding increase in access to the US market.

This page titled 15.5: Trade barriers- Tariffs, subsidies and quotas is shared under a CC BY-NC-SA 4.0 license and was authored, remixed, and/or curated by Douglas Curtis and Ian Irvine (Lyryx) via source content that was edited to the style and standards of the LibreTexts platform.

15.6: The politics of protection

Objections to imports are frequent and come from many different sectors of the economy. In the face of the gains from trade which we have illustrated in this chapter, why do we observe such strong opposition to imported goods and services?

Structural change and technology

In a nutshell the answer is that, while consumers in the aggregate gain from the reduction of trade barriers, and there is a net gain to the economy at large, some *individual sectors of the economy lose out*. Not surprisingly the sectors that will be adversely affected are vociferous in lodging their objections. Sectors of the economy that cannot compete with overseas suppliers generally see a reduction in jobs. This has been the case in the manufacturing sector of the Canadian and US economies in recent decades, as manufacturing and assembly has flown off-shore to Asia and Mexico where labour costs are lower. Domestic job losses are painful, and frequently workers who have spent decades in a particular job find reemployment difficult, and rarely get as high a wage as in their displaced job.

Such job losses are reflected in calls for tariffs on imports from China, for example, in order to 'level the playing field' – that is, to counter the impact of lower wages in China. Of course it is precisely because of lower labour costs in China that the Canadian consumer benefits.

In Canada we deal with such dislocation first by providing unemployment payments to workers, and by furnishing retraining allowances, both coming from Canada's Employment Insurance program. While such support does not guarantee an equally good alternative job, structural changes in the economy, due to both internal and external developments, must be confronted. For example, the information technology revolution made tens of thousands of 'data entry' workers redundant. Should producers have shunned the technological developments which increased their productivity dramatically? If they did, would they be able to compete in world markets?

While job losses feature heavily in protests against technological development and freer trade, most modern economies continue to grow and create more jobs in the service sector than are lost in the manufacturing sector. Developed economies now have many more workers in service than manufacture. Service jobs are not just composed of low-wage jobs in fast food establishments – 'Mcjobs', they are high paying jobs in the health, education, legal, financial and communications sectors of the economy.

Successful lobbying and concentration

While efforts to protect manufacture have not resulted in significant barriers to imports of manufactures, objections in some specific sectors of the economy seem to be effective worldwide. One sector that stands out is agriculture, where political conditions are conducive to the continuance of protection and what is called 'supply management' – domestic production quotas. The reason for 'successful' supply limitation appears to rest in the geographic concentration of potential beneficiaries of such protection and the scattered beneficiaries of freer trade on the one hand, and the costs and benefits of political organization on the other: Farmers tend to be concentrated in a limited number of rural electoral ridings and hence they can collectively have a major impact on electoral outcomes. Second, the benefits that accrue to trade restriction are heavily concentrated in the economy – keep in mind that about two percent of the population lives on farms, or relies on farming for its income. By contrast the costs on a per person scale are small, and are spread over the whole population. Thus, in terms of the costs of political organization, the incentives for consumers are small, but the incentives for producers are high.

In addition to the differing patterns of costs and benefits, rural communities tend to be more successful in pushing trade restrictions based on a 'way-of-life' argument. By permitting imports that might displace local supply, lobbyists are frequently successful in convincing politicians that long-standing way-of-life traditions would be endangered, even if such 'traditions' are accompanied by monopolies and exceptionally high tariffs.

Valid trade barriers: Infant industries and dumping?

An argument that carries both intellectual and emotional appeal to voters is the 'infant industry' argument. It goes as follows: New ventures and sectors of the economy may require time before that can compete internationally. Scale economies may be involved, for example, and time may be required for producers to expand their scale of operation, at which time costs will have fallen to international (i.e. competitive) levels. In addition, learning-by-doing may be critical in more high-tech sectors and, once again, with the passage of time costs should decline for this reason also.

The problem with this stance is that these 'infants' have insufficient incentive to 'grow up' and become competitive. A protection measure that is initially intended to be temporary can become permanent because of the potential job losses associated with a cessation of the protection to an industry that fails to become internationally competitive. Furthermore, employees and managers in protected sectors have insufficient incentive to make their production competitive if they realize that their government will always be there to protect them.

In contrast to the infant industry argument, economists are more favourable to restrictions that are aimed at preventing dumping.

Dumping may occur either because foreign suppliers choose to sell at artificially low prices (prices below their break-even price for example), or because of surpluses in foreign markets resulting from oversupply. For example, if, as a result of price support in its own market, a foreign government induced oversupply in butter and it chose to sell such butter on world markets at a price well below the going ('competitive') world supply price, such a sale would constitute dumping. Alternatively, an established foreign supplier might choose to enter our domestic market by selling its products at artificially low prices, with a view to driving domestic competition out of the domestic market. Having driven out the domestic competition it would then be in a position to raise prices. This is predatory pricing as explored in the last chapter. Such behaviour differs from a permanently lower price on the part of foreign suppliers. This latter may be welcomed as a gain from trade, whereas the former may generate no gains and serve only to displace domestic labour and capital.

Protectionism in the age of pandemics

The year 2020 will be remembered in history as the year of the coronavirus pandemic. An uncountable number of men and women died all across the globe as a result of contracting COVID-19, the respiratory disorder brought on by an attack of the coronavirus. In the absence of a vaccine, health authorities the world over implemented a twin policy of social distancing and quarantining (or self-isolation). The world economy went into a tailspin, as huge fractions of the labor force were laid off. Trade patterns were disrupted and serious shortages of personal protection equipment (PPE - masks, visors, gowns), ventilators and drugs emerged. The world demand for PPE and ventilators skyrocketed. But the production of PPE was concentrated in China; most western economies did not have the necessary productive capacity to supply even non-pandemic requirements. Bidding wars erupted amongst countries and hospitals as they vied for supply, while domestic producers of some products added to their production capacity.

Following this chaos, we ask if self-sufficiency would not be a better model than open trade. Would a world where each country ensured it had the production capacity to produce these necessities in times of emergency not be superior to one where global supply chains characterize everything from computers to generic drugs? India is a major producer of generic drugs and the components for such drugs. The demand for anti-biotics and pain killers also rocketed upwards with the pandemic.

There is more than one way to plan for a pandemic, and such planning should not involve a generalized move to self-sufficiency on the part of the global economy. One strategy is to build up inventories of PPE and ventilators domestically. This is costly, but for the most part feasible. It does not represent a complete solution because technology changes will make 30-year old ventilators sitting in inventory redundant for the next pandemic. In addition, most medications have a limited shelf life. Hence one solution is to maintain and rotate substantial inventories of emergency equipment using existing supply chains, and benefit from the efficiencies that are built into these chains.

A second option is to maintain excess production capacity on the part of domestic manufacturers of critical pandemic products. Maintaining such capacity should be considered at least partially as a social cost; pandemics ravage societies, not just individuals, and therefore society should undertake part of the cost of insuring against them.

A more general argument against global trade comes in the form of protecting food supplies. In the early 2000s an increase in global cereal prices led some economies to limit exports of specific crops on account of the fact that global demand was pushing prices to a level that low-income consumers could not afford. But such a policy may threaten consumers in other low-income economies whose demands have not changed in a context of reduced supplies. The reality is that world food supply is adequate for world consumption, even in the presence of disruptions. It is also the case that certain economies have huge advantages in producing specific kinds of food. For example, Canada, the US and the Ukraine produce cereals very economically. Mountainous regions are unsuitable for this production. It would benefit no economy for these economies to lower their production of grains to the point where they produced only enough for their own production. By the same reasoning, warmer climates produce fruits, coffee beans, olives etc that cannot easily be produced in many regions suited to wheat. The gains to specialization in the world economy are enormous. Where food shortages occur we frequently encounter the scourges of drought or war or political upheaval, and these conditions inhibit the distribution of foodstuffs.

What about supply chains? If motherboards produced in China are not being exported in sufficient quantities then indeed production of computers in North America will suffer. But to infer from this that North America should decide to produce all of its computer components in North America is illogical. First, in the time of a pandemic, if certain economies in the supply chain are on lockdown, we cannot be sure that the domestic economy would not be on lockdown simultaneously. Second, the cost to moving the production of all computer parts to North America would likely double the cost of computer hardware - including cell-phones. Perhaps a disruption to our supply chains is something we need to bear in extraordinary times. In case it requires emphasis, most producers in supply chains have incentives to produce and sell. If they do not they will die economically.

The energy sector of every economy is impacted with the outbreak of a pandemic. This is because the demand for fuel (primarily oil) declines following policies of social distancing, limits on permissible travel, and the closure of some production facilities that depend upon oil. In North America, as we saw in Chapter 4 earlier, the price of oil declined from US \$60 per barrel to US \$20 in the space of two months in early 2020. Since production costs are higher in both Canada and much of the US than in Saudi Arabia, the North Sea and Russia, producers in North America were squeezed. Many were no longer able to cover their full production costs, and forced to cease drilling and recovering oil. Inevitably there was a clamor for protection. Producers sought tariffs on competing oil: Tariffs would increase the price of cheaper-to-produce foreign oil and enable domestic producers to survive.

While protection might seem like a 'sensible' policy in this instance, the fact is that unilateral tariffs usually invite reprisals, and raise the danger of a trade war with ever-expanding counter protectionism. In contrast to the case of a *shortage* of medical supplies, the energy sector in Canada suffered from a *glut* of world oil supply. The domestic issue is not about the health of consumers (as in the case of medical supplies), it is about the health of producers.

To conclude: a pandemic is a profoundly serious event and such events inflict major costs on all societies. There are no magic bullets in the form of low-cost ideal economic policies to counter viral warfare. The key to policy making is to recognize constraints and recognize an attack as soon as possible. A wholesale move to insulate the domestic economy is ill-conceived. Comparative advantage confers enormous benefits to all nations. Specific policies should take the form of inventory management and excess production capacity in specific sectors of the economy.

This page titled [15.6: The politics of protection](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.7: Institutions governing trade

In the nineteenth century, world trade grew rapidly, in part because the leading trading nation at the time—the United Kingdom—pursued a vigorous policy of free trade. In contrast, US tariffs averaged about 50 percent, although they had fallen to around 30 percent by the early 1920s. As the industrial economies went into the Great Depression of the late 1920s and 1930s, there was pressure to protect domestic jobs by keeping out imports. Tariffs in the United States returned to around 50 percent, and the United Kingdom abandoned the policy of free trade that had been pursued for nearly a century. The combination of world recession and increasing tariffs led to a disastrous slump in the volume of world trade, further exacerbated by World War II.

The WTO and GATT

After World War II, there was a collective determination to see world trade restored. Bodies such as the International Monetary Fund and the World Bank were set up, and many countries signed the General Agreement on Tariffs and Trade (GATT), a commitment to reduce tariffs progressively and dismantle trade restrictions.

Under successive rounds of GATT, tariffs fell steadily. By 1960, United States tariffs were only one-fifth of their level at the outbreak of the War. In the United Kingdom, the system of wartime quotas on imports had been dismantled by the mid-1950s, after which tariffs were reduced by nearly half in the ensuing 25 years. Europe as a whole moved toward an enlarged European Union in which tariffs between member countries have been abolished. By the late 1980s, Canada's tariffs had been reduced to about one-quarter of their immediate post-World War II level.

The GATT Secretariat, now called the World Trade Organization (WTO), aims both to dismantle existing protection that reduces efficiency and to extend trade liberalization to more and more countries. Tariff levels throughout the world are now as low as they have ever been, and trade liberalization has been an engine of growth for many economies. The consequence has been a substantial growth in world trade.

NAFTA, the USMCA, the EU, the CETA, and the TPP

In North America, policy since the 1980s has led to a free trade area that covers the flow of trade between Canada, the United States, and Mexico. The Canada/United States free trade agreement (FTA) of 1989 expanded in 1994 to include Mexico in the North American Free Trade Agreement (NAFTA). The objective in both cases was to institute freer trade between these countries in most goods and services. This meant the elimination or reduction of tariffs and non-tariff barriers over a period of years, with a few exceptions in specific products and cultural industries. A critical component of the Agreement was the establishment of a dispute-resolution mechanism, under which disputes would be resolved by a panel of 'judges' nominated from the member economies. Evidence of the success of these agreements is reflected in the fact that Canadian exports have grown to more than 30 percent of GDP, and trade with the United States accounts for the lion's share of Canadian trade flows. NAFTA was updated and replaced in 2018 and the new agreement is termed the United States Mexico Canada Agreement.

The European Union was formed after World War II, with the prime objective of bringing about a greater degree of *political* integration in Europe. Two world wars had laid waste to their economies and social fabric. Closer economic ties and greater trade were seen as the means of achieving this integration. The Union was called the "Common Market" for much of its existence. The Union originally had six member states, and as of 2019 the number is 28, with several other candidate countries in the process of application, most notably Turkey. The European Union (EU) has a secretariat and parliament in Bruxelles. The UK intends to exit the EU as of late 2019.

Canada has concluded a free trade agreement with the European Union that is termed the Comprehensive Economic and Trade Agreement (CETA). It has the objective of implementing free trade between the two negotiating parties, though there remain some exceptions, for example agriculture.

The Comprehensive and Progressive Agreement for Trans-Pacific Partnership (CPTPP) is a trading agreement between Canada and ten other Pacific-Rim economies that came into being in 2018. Negotiations for a Trans Pacific Partnership treaty were complete by 2016. Those negotiations involved 12 Pacific Rim economies including Canada and the United States, but excluding China. The Obama presidency appeared ready to sign the treaty, however the Trump presidency (and also the Democratic candidate for president of the US, Hillary Clinton) decided that the Partnership was not in the interests of the United States and withdrew its affiliation. The remaining 11 economies reached an agreement to implement the partnership in December 2018.

This page titled [15.7: Institutions governing trade](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.8: Key Terms

Autarky denotes the no-trade situation.

Principle of comparative advantage states that even if one country has an absolute advantage in producing both goods, gains to specialization and trade still materialize, provided the opportunity cost of producing the goods differs between economies.

Terms of trade define the rate at which goods trade internationally.

Consumption possibility frontier defines what an economy can consume after production specialization and trade.

Intra-industry trade is two-way international trade in products produced within the same industry.

Intra-firm trade is two-way trade in international products produced within the same firm.

Supply chain: denotes the numerous sources for intermediate goods used in producing a final product.

Intermediate good: one that is used in the production of final output.

Content requirement: requires that a specified percentage of the final value of a product originate in the producing economy.

Dynamic gains: the potential for domestic producers to increase productivity by competing with, and learning from, foreign producers.

Total factor productivity: how efficiently the factors of production are combined.

Tariff is a tax on an imported product that is designed to limit trade in addition to generating tax revenue. It is a barrier to trade.

Quota is a quantitative limit on an imported product.

Trade subsidy to a domestic manufacturer reduces the domestic cost and limits imports.

Non-tariff barriers, such as product content requirements, limits the gains from trade.

Dumping is a predatory practice, based on artificial costs aimed at driving out domestic producers.

This page titled [15.8: Key Terms](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

15.9: Exercises for Chapter 15

EXERCISE 15.1

The following table shows the labour input requirements to produce a bushel of wheat and a litre of wine in two countries, Northland and Southland, on the assumption of constant cost production technology – meaning that the production possibility curves in each are straight lines. You can answer this question either by analyzing the table or developing a graph similar to Figure 15.1, assuming each economy has 4 units of labour.

Labour requirements per unit produced		
	Northland	Southland
Per bushel of wheat	1	3
Per litre of wine	2	4

- Which country has an absolute advantage in the production of both wheat and wine?
- What is the opportunity cost of wheat in each economy? Of wine?
- What is the pattern of comparative advantage here?
- Suppose the country with a comparative advantage in wine reduces wheat production by one bushel and reallocates the labour involved to wine production. How much additional wine does it produce?

EXERCISE 15.2

Canada and the United States can produce two goods, xylophones and yogurt. Each good can be produced with labour alone. Canada requires 60 hours to produce a ton of yogurt and 6 hours to produce a xylophone. The United States requires 40 hours to produce the ton of yogurt and 5 hours to produce a xylophone.

- Describe the state of absolute advantage between these economies in producing goods.
- In which good does Canada have a comparative advantage? Does this mean the United States has a comparative advantage in the other good?
- Draw the production possibility frontier for each economy to scale on a diagram, assuming that each economy has an endowment of 240 hours of labour, and that the PPFs are linear.
- On the same diagram, draw Canada's consumption possibility frontier on the assumption that it can trade with the United States at the United States' rate of transformation.
- Draw the US consumption possibility frontier under the assumption that it can trade at Canada's rate of transformation.

EXERCISE 15.3

The domestic demand for bicycles is given by $P=36-0.3Q$. The foreign supply is given by $P=18$ and domestic supply by $P=16+0.4Q$.

- Illustrate the market equilibrium on a diagram, and illustrate the amounts supplied by domestic and foreign suppliers in equilibrium.
- If the government now imposes a tariff of \$6 per unit on the foreign good, illustrate the impact geometrically.
- In the diagram, illustrate the area representing tariff revenue.
- Optional:* Compute the price and quantity in equilibrium with free trade, and again in the presence of the tariff.

EXERCISE 15.4

- In Exercise 15.3, illustrate graphically the deadweight losses associated with the imposition of the tariff.
- Illustrate on your diagram the additional amount of profit made by the domestic producer as a result of the tariff. [*Hint:* Refer to Figure 15.4 in the text.]

EXERCISE 15.5

The domestic demand for office printers is given by $P=40-0.2Q$. The supply of domestic producers is given by $P=12+0.1Q$, and international supply by $P=20$.

- Illustrate this market geometrically.
- If the government gives a production subsidy of \$2 per unit to domestic suppliers in order to increase their competitiveness, illustrate the impact of this on the domestic supply curve.
- Illustrate geometrically the cost to the government of this scheme.

EXERCISE 15.6

Consider the data underlying Figure 15.1. Suppose, from the initial state of comparative advantage, where Canada specializes in fish and the US in vegetable, we have a technological change in fishing. The US invents the multi-hook fishing line, and as a result can now produce 64 units of fish with the same amount of labour, rather than the 40 units it could produce before the technological change. This technology does not spread to Canada however.

- Illustrate the new *PPF* for the US in addition to the *PPF* for Canada.
- What is the new opportunity cost (number of fish) associated with one unit of *V*?
- Has comparative advantage changed here – which economy should specialize in the production of each good?

EXERCISE 15.7

The following are hypothetical (straight line) production possibilities tables for Canada and the United States. For each line required, plot any two or more points on the line.

Canada					United States				
	A	B	C	D		A	B	C	D
Peaches	0	5	10	15	Peaches	0	10	20	30
Apples	30	20	10	0	Apples	15	10	5	0

- Plot Canada's production possibilities curve.
 - Plot the United States' production possibilities curve.
 - What is each country's cost ratio of producing peaches and apples?
 - Which economy should specialize in which product?
 - Plot the United States' trading possibilities curve (by plotting at least 2 points on the curve) if the actual terms of the trade are 1 apple for 1 peach.
 - Plot the Canada' trading possibilities curve (by plotting at least 2 points on the curve) if the actual terms of the trade are 1 apple for 1 peach.
 - Suppose that the optimum product mixes before specialization and trade were B in the United States and C in Canada. What are the gains from specialization and trade?
- Note that we are considering the *PPFs* to be straight lines rather than concave shapes. The result we illustrate here carries over to that case also, but it is simpler to illustrate with the linear *PPFs*.
 - To illustrate the gains numerically, let Canada import 3*V* from the US in return for exporting 18*F*. Note that this is a trading rate of 1:6. Hence, Canada consumes 3*V* and 17*F* (Canada produced 35*F* and exported 18*F*, leaving it with 17*F*). It follows that the US consumes 5*V*, having exported 3*V* of the 8*V* it produced, and obtained in return 18*F* in imports. The new consumption bundles are illustrated in the figure: (17,3) for Canada and (18,5) for the US. These consumption bundles clearly represent an improvement over the autarky situation. For example, if Canada had wished to consume 3*V* pre-trade, it would only have been able to consume 14*F*, whereas with trade it can consume 17*F*.

This page titled [15.9: Exercises for Chapter 15](#) is shared under a [CC BY-NC-SA 4.0](#) license and was authored, remixed, and/or curated by [Douglas Curtis and Ian Irvine \(Lyryx\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.

Index

D

dire

Detailed Licensing

Overview

Title: [Principles of Microeconomics \(Curtis and Irvine\)](#)

Webpages: 182

Applicable Restrictions: Noncommercial

All licenses found:

- [CC BY-NC-SA 4.0](#): 98.4% (179 pages)
- [Undeclared](#): 1.6% (3 pages)

By Page

- [Principles of Microeconomics \(Curtis and Irvine\) - CC BY-NC-SA 4.0](#)
 - [Front Matter - CC BY-NC-SA 4.0](#)
 - [TitlePage - CC BY-NC-SA 4.0](#)
 - [InfoPage - CC BY-NC-SA 4.0](#)
 - [Table of Contents - Undeclared](#)
 - [Licensing - Undeclared](#)
 - [Unit 1: The Building Blocks - CC BY-NC-SA 4.0](#)
 - [1: Introduction to key ideas - CC BY-NC-SA 4.0](#)
 - [1.1: What's it all about? - CC BY-NC-SA 4.0](#)
 - [1.2: Understanding through the use of models - CC BY-NC-SA 4.0](#)
 - [1.3: Opportunity cost and the market - CC BY-NC-SA 4.0](#)
 - [1.4: A model of exchange and specialization - CC BY-NC-SA 4.0](#)
 - [1.5: Economy-wide production possibilities - CC BY-NC-SA 4.0](#)
 - [1.6: New Page - CC BY-NC-SA 4.0](#)
 - [1.7: Conclusion - CC BY-NC-SA 4.0](#)
 - [1.8: Key Terms - CC BY-NC-SA 4.0](#)
 - [1.9: Exercises for Chapter 1 - CC BY-NC-SA 4.0](#)
 - [2: Theories, data and beliefs - CC BY-NC-SA 4.0](#)
 - [2.1: Data analysis - CC BY-NC-SA 4.0](#)
 - [2.2: Data, theory and economic models - CC BY-NC-SA 4.0](#)
 - [2.3: Ethics, efficiency and beliefs - CC BY-NC-SA 4.0](#)
 - [2.4: Key Terms - CC BY-NC-SA 4.0](#)
 - [2.5: Exercises for Chapter 2 - CC BY-NC-SA 4.0](#)
 - [3: The classical marketplace – demand and supply - CC BY-NC-SA 4.0](#)
 - [3.1: The marketplace - trading - CC BY-NC-SA 4.0](#)
 - [3.2: The market's building blocks - CC BY-NC-SA 4.0](#)
 - [3.3: Demand and supply curves - CC BY-NC-SA 4.0](#)
 - [3.4: Non-price influences on demand - CC BY-NC-SA 4.0](#)
 - [3.5: Non-price influences on supply - CC BY-NC-SA 4.0](#)
 - [3.6: Simultaneous supply and demand impacts - CC BY-NC-SA 4.0](#)
 - [3.7: Market interventions - governments and interest groups - CC BY-NC-SA 4.0](#)
 - [3.8: Individual and market functions - CC BY-NC-SA 4.0](#)
 - [3.9: Useful techniques - demand and supply equations - CC BY-NC-SA 4.0](#)
 - [3.10: Conclusion - CC BY-NC-SA 4.0](#)
 - [3.11: Key Terms - CC BY-NC-SA 4.0](#)
 - [3.12: Exercises for Chapter 3 - CC BY-NC-SA 4.0](#)
 - [Unit 2: Responsiveness and the Value of Markets - CC BY-NC-SA 4.0](#)
 - [4: Measures of response- Elasticities - CC BY-NC-SA 4.0](#)
 - [4.1: Price responsiveness of demand - CC BY-NC-SA 4.0](#)
 - [4.2: Price elasticities and public policy - CC BY-NC-SA 4.0](#)
 - [4.3: The time horizon and inflation - CC BY-NC-SA 4.0](#)
 - [4.4: Cross-price elasticities - cable or satellite - CC BY-NC-SA 4.0](#)
 - [4.5: The income elasticity of demand - CC BY-NC-SA 4.0](#)
 - [4.6: Elasticity of supply - CC BY-NC-SA 4.0](#)
 - [4.7: Elasticities and tax incidence - CC BY-NC-SA 4.0](#)
 - [4.8: Technical tricks with elasticities - CC BY-NC-SA 4.0](#)
 - [4.9: Key Terms - CC BY-NC-SA 4.0](#)
 - [4.10: Exercises for Chapter 4 - CC BY-NC-SA 4.0](#)
 - [5: Welfare economics and externalities - CC BY-NC-SA 4.0](#)

- 5.1: Equity and efficiency - CC BY-NC-SA 4.0
- 5.2: Consumer and producer surplus - CC BY-NC-SA 4.0
- 5.3: Efficient market outcomes - CC BY-NC-SA 4.0
- 5.4: Taxation, surplus and efficiency - CC BY-NC-SA 4.0
- 5.5: Market failures - externalities - CC BY-NC-SA 4.0
- 5.6: Other market failures - CC BY-NC-SA 4.0
- 5.7: Environmental policy and climate change - CC BY-NC-SA 4.0
- 5.8: Conclusion - CC BY-NC-SA 4.0
- 5.9: Key Terms - CC BY-NC-SA 4.0
- 5.10: Exercises for Chapter 5 - CC BY-NC-SA 4.0
- Unit 3: Decision Making by Consumer and Producers - CC BY-NC-SA 4.0
 - 6: Individual choice - CC BY-NC-SA 4.0
 - 6.1: Rationality - CC BY-NC-SA 4.0
 - 6.2: Choice with measurable utility - CC BY-NC-SA 4.0
 - 6.3: Choice with ordinal utility - CC BY-NC-SA 4.0
 - 6.4: Applications of indifference analysis - CC BY-NC-SA 4.0
 - 6.5: Key Terms - CC BY-NC-SA 4.0
 - 6.6: Exercises for Chapter 6 - CC BY-NC-SA 4.0
 - 7: Firms, investors and capital markets - CC BY-NC-SA 4.0
 - 7.1: Business organization - CC BY-NC-SA 4.0
 - 7.2: Profit - CC BY-NC-SA 4.0
 - 7.3: Risk and the investor - CC BY-NC-SA 4.0
 - 7.4: Risk pooling and diversification - CC BY-NC-SA 4.0
 - 7.5: Conclusion - CC BY-NC-SA 4.0
 - 7.6: Key Terms - CC BY-NC-SA 4.0
 - 7.7: Exercises for Chapter 7 - CC BY-NC-SA 4.0
 - 8: Production and cost - CC BY-NC-SA 4.0
 - 8.1: Efficient production - CC BY-NC-SA 4.0
 - 8.2: The time frame - CC BY-NC-SA 4.0
 - 8.3: Production in the short run - CC BY-NC-SA 4.0
 - 8.4: Costs in the short run - CC BY-NC-SA 4.0
 - 8.5: Fixed costs and sunk costs - CC BY-NC-SA 4.0
 - 8.6: Long-run production and costs - CC BY-NC-SA 4.0
 - 8.7: Technological change- globalization and localization - CC BY-NC-SA 4.0
 - 8.8: Clusters, learning by doing, scope economics - CC BY-NC-SA 4.0
 - 8.9: Conclusion - CC BY-NC-SA 4.0
 - 8.10: Key Terms - CC BY-NC-SA 4.0
 - 8.11: Exercises for Chapter 8 - CC BY-NC-SA 4.0
- Unit 4: Market Structures - CC BY-NC-SA 4.0
 - 9: Perfect competition - CC BY-NC-SA 4.0
 - 9.1: The perfect competition paradigm - CC BY-NC-SA 4.0
 - 9.2: Market characteristics - CC BY-NC-SA 4.0
 - 9.3: The firm's supply decision - CC BY-NC-SA 4.0
 - 9.4: Dynamics- Entry and exit - CC BY-NC-SA 4.0
 - 9.5: Long-run industry supply - CC BY-NC-SA 4.0
 - 9.6: Globalization and technological change - CC BY-NC-SA 4.0
 - 9.7: Efficient resource allocation - CC BY-NC-SA 4.0
 - 9.8: Key Terms - CC BY-NC-SA 4.0
 - 9.9: Exercises for Chapter 9 - CC BY-NC-SA 4.0
 - 10: Monopoly - CC BY-NC-SA 4.0
 - 10.1: Monopolies - CC BY-NC-SA 4.0
 - 10.2: Profit maximizing behaviour - CC BY-NC-SA 4.0
 - 10.3: Long-run choices - CC BY-NC-SA 4.0
 - 10.4: Output inefficiency - CC BY-NC-SA 4.0
 - 10.5: Price discrimination - CC BY-NC-SA 4.0
 - 10.6: Cartels- Acting like a monopolist - CC BY-NC-SA 4.0
 - 10.7: Invention, innovation and rent seeking - CC BY-NC-SA 4.0
 - 10.8: Conclusion - CC BY-NC-SA 4.0
 - 10.9: Key Terms - CC BY-NC-SA 4.0
 - 10.10: Exercises for Chapter 10 - CC BY-NC-SA 4.0
 - 11: Imperfect competition - CC BY-NC-SA 4.0
 - 11.1: The principle ideas - CC BY-NC-SA 4.0
 - 11.2: Imperfect competitors - CC BY-NC-SA 4.0
 - 11.3: Imperfect competitors- measures of structure and market power - CC BY-NC-SA 4.0
 - 11.4: Imperfect competition- monopolistic competition - CC BY-NC-SA 4.0
 - 11.5: Imperfect competition- economies of scope and platforms - CC BY-NC-SA 4.0
 - 11.6: Strategic behaviour- Oligopoly and games - CC BY-NC-SA 4.0
 - 11.7: Strategic behaviour- Duopoly and Cournot games - CC BY-NC-SA 4.0
 - 11.8: Strategic behaviour- Entry, exit and potential competition - CC BY-NC-SA 4.0
 - 11.9: Matching markets- design - CC BY-NC-SA 4.0
 - 11.10: Conclusion - CC BY-NC-SA 4.0

- 11.11: Key Terms - CC BY-NC-SA 4.0
- 11.12: Exercises for Chapter 11 - CC BY-NC-SA 4.0
- Unit 5: The Factors of Production - CC BY-NC-SA 4.0
 - 12: Labour and capital - CC BY-NC-SA 4.0
 - 12.1: Labour - a derived demand - CC BY-NC-SA 4.0
 - 12.2: The supply of labour - CC BY-NC-SA 4.0
 - 12.3: Labour market equilibrium and mobility - CC BY-NC-SA 4.0
 - 12.4: Capital - concepts - CC BY-NC-SA 4.0
 - 12.5: The capital market - CC BY-NC-SA 4.0
 - 12.6: Land - CC BY-NC-SA 4.0
 - 12.7: Key Terms - CC BY-NC-SA 4.0
 - 12.8: Exercises for Chapter 12 - CC BY-NC-SA 4.0
 - 13: Human capital and the income distribution - CC BY-NC-SA 4.0
 - 13.1: Human capital - CC BY-NC-SA 4.0
 - 13.2: Productivity and education - CC BY-NC-SA 4.0
 - 13.3: On-the-job training - CC BY-NC-SA 4.0
 - 13.4: Education as signalling - CC BY-NC-SA 4.0
 - 13.5: Education returns and quality - CC BY-NC-SA 4.0
 - 13.6: Discrimination - CC BY-NC-SA 4.0
 - 13.7: The income distribution - CC BY-NC-SA 4.0
 - 13.8: Wealth and capitalism - CC BY-NC-SA 4.0
 - 13.9: Key Terms - CC BY-NC-SA 4.0
 - 13.10: Exercises for Chapter 13 - CC BY-NC-SA 4.0
- Unit 6: Government and Trade - CC BY-NC-SA 4.0
 - 14: Government - CC BY-NC-SA 4.0
 - 14.1: Market Failure - CC BY-NC-SA 4.0
 - 14.2: Fiscal federalism- Taxing and spending - CC BY-NC-SA 4.0
 - 14.3: Federal-provincial fiscal relations - CC BY-NC-SA 4.0
 - 14.4: Government-to-individual transfers - CC BY-NC-SA 4.0
 - 14.5: Regulation and competition policy - CC BY-NC-SA 4.0
 - 14.6: Key Terms - CC BY-NC-SA 4.0
 - 14.7: Exercises for Chapter 14 - CC BY-NC-SA 4.0
- 15: International trade - CC BY-NC-SA 4.0
 - 15.1: Trade in our daily lives - CC BY-NC-SA 4.0
 - 15.2: Canada in the world economy - CC BY-NC-SA 4.0
 - 15.3: The gains from trade- Comparative advantage - CC BY-NC-SA 4.0
 - 15.4: Returns to scale and dynamic gains from trade - CC BY-NC-SA 4.0
 - 15.5: Trade barriers- Tariffs, subsidies and quotas - CC BY-NC-SA 4.0
 - 15.6: The politics of protection - CC BY-NC-SA 4.0
 - 15.7: Institutions governing trade - CC BY-NC-SA 4.0
 - 15.8: Key Terms - CC BY-NC-SA 4.0
 - 15.9: Exercises for Chapter 15 - CC BY-NC-SA 4.0
- Back Matter - CC BY-NC-SA 4.0
 - Index - CC BY-NC-SA 4.0
 - Glossary - CC BY-NC-SA 4.0
 - Solutions To Exercises - CC BY-NC-SA 4.0
 - 1.1: Chapter 1 Solutions - CC BY-NC-SA 4.0
 - 1.2: Chapter 2 Solutions - CC BY-NC-SA 4.0
 - 1.3: Chapter 3 Solutions - CC BY-NC-SA 4.0
 - 1.4: Chapter 4 Solutions - CC BY-NC-SA 4.0
 - 1.5: Chapter 5 Solutions - CC BY-NC-SA 4.0
 - 1.6: Chapter 6 Solutions - CC BY-NC-SA 4.0
 - 1.7: Chapter 7 Solutions - CC BY-NC-SA 4.0
 - 1.8: Chapter 8 Solutions - CC BY-NC-SA 4.0
 - 1.9: Chapter 9 Solutions - CC BY-NC-SA 4.0
 - 1.10: Chapter 10 Solutions - CC BY-NC-SA 4.0
 - 1.11: Chapter 11 Solutions - CC BY-NC-SA 4.0
 - 1.12: Chapter 12 Solutions - CC BY-NC-SA 4.0
 - 1.13: Chapter 13 Solutions - CC BY-NC-SA 4.0
 - 1.14: Chapter 14 Solutions - CC BY-NC-SA 4.0
 - 1.15: Chapter 15 Solutions - CC BY-NC-SA 4.0
 - Detailed Licensing - Undeclared