

1.2: Chapter 2- Describing Data Using Distributions and Graphs

Key Terms

[bell curve](#)

[bimodal distribution](#)

[bin widths](#)

[box plots](#)

[categorical variables](#)

[frequency polygons](#)

[histogram](#)

[lie factor](#)

[skew](#)

[stem-and-leaf display](#)

[whiskers](#)

Before we can understand our analyses, we must first understand our data. The first step in doing this is using tables, charts, graphs, plots, and other visual tools to see what our data look like.

Graphing Qualitative Variables

When Apple Computer introduced the iMac computer in August 1998, the company wanted to learn whether the iMac was expanding Apple's market share. Was the iMac just attracting previous Macintosh owners? Or was it purchased by newcomers to the computer market and by previous Windows users who were switching over? To find out, 500 iMac customers were interviewed. Each customer was categorized as a previous Macintosh owner, a previous Windows owner, or a new computer purchaser.

This section examines graphical methods for displaying the results of the interviews. We'll learn some general lessons about how to graph data that fall into a small number of categories. A later section will consider how to graph numerical data in which each observation is represented by a number in some range. The key point about the qualitative data that occupy us in the present section is that they do not come with a pre-established ordering (the way numbers are ordered). For example, there is no natural sense in which the category of previous Windows users comes before or after the category of previous Macintosh users. This situation may be contrasted with quantitative data, such as a person's weight. People of one weight are naturally ordered with respect to people of a different weight.

Frequency Tables

All of the graphical methods shown in this section are derived from frequency tables. [Table 2.1](#) shows a frequency table for the results of the iMac study; it shows the frequencies of the various response categories. It also shows the relative frequencies, which are the proportion of responses in each category. For example, the relative frequency for "none" of $.17 = 85/500$.

Table 2.1. Frequency table for the iMac data.

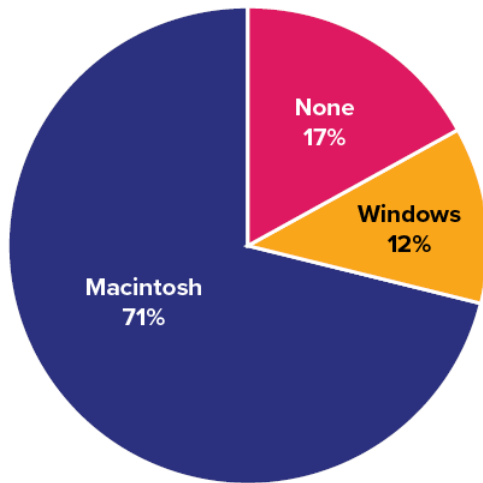
Previous Ownership	Frequency	Relative Frequency
None	85	.17
Windows	60	.12
Macintosh	355	.71
Total	500	1.00

Pie Charts

The pie chart in [Figure 2.1](#) shows the results of the iMac study. In a pie chart, each category is represented by a slice of the pie. The area of the slice is proportional to the percentage of responses in the category. This is simply the relative frequency multiplied by

100. Although most iMac purchasers were Macintosh owners (71%), Apple was encouraged by the 12% of purchasers who were former Windows users, and by the 17% of purchasers who were buying a computer for the first time.

Figure 2.1. Pie chart of iMac purchases illustrating frequencies of previous computer ownership: 71% of purchasers owned a Macintosh before buying their iMac. (“[Mac Pie Chart](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Pie charts are effective for displaying the relative frequencies of a small number of categories. They are not recommended, however, when you have a large number of categories. Pie charts can also be confusing when they are used to compare the outcomes of two different surveys or experiments. In an influential book on the use of graphs, Edward Tufte asserted, “The only worse design than a pie chart is several of them.”¹

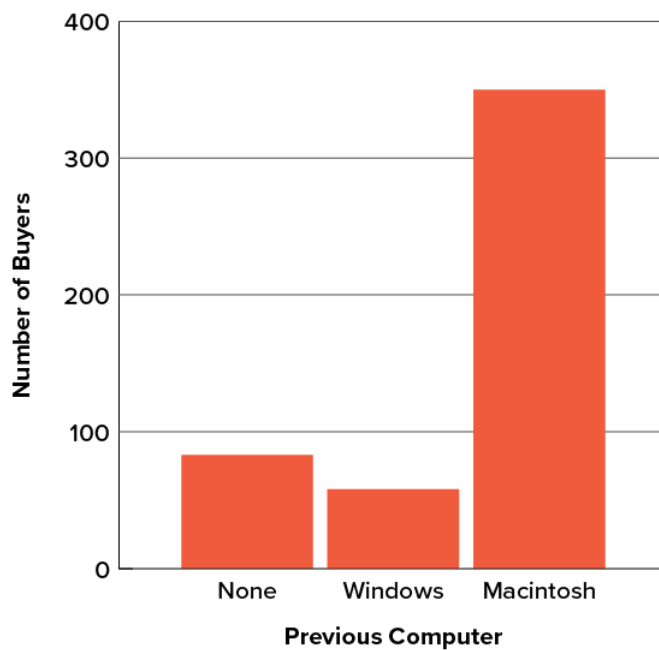
¹ Tufte, E. R. (1983). The visual display of quantitative information (p. 178). Graphics Press.

Here is another important point about pie charts. If they are based on a small number of observations, it can be misleading to label the pie slices with percentages. For example, if just 5 people had been interviewed by Apple Computers, and 3 were former Windows users, it would be misleading to display a pie chart with the Windows slice showing 60%. With so few people interviewed, such a large percentage of Windows users might easily have occurred since chance can cause large errors with small samples. In this case, it is better to alert the user of the pie chart to the actual numbers involved. The slices should therefore be labeled with the actual frequencies observed (e.g., 3) instead of with percentages.

Bar Charts

Bar charts can also be used to represent frequencies of different categories. A bar chart of the iMac purchases is shown in [Figure 2.2](#). Frequencies are shown on the y-axis and the type of computer previously owned is shown on the x-axis. Typically, the y-axis shows the number of observations in each category rather than the percentage of observations in each category as is typical in pie charts.

Figure 2.2. Bar chart of iMac purchases as a function of previous computer ownership. (“[Mac Bar Chart](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

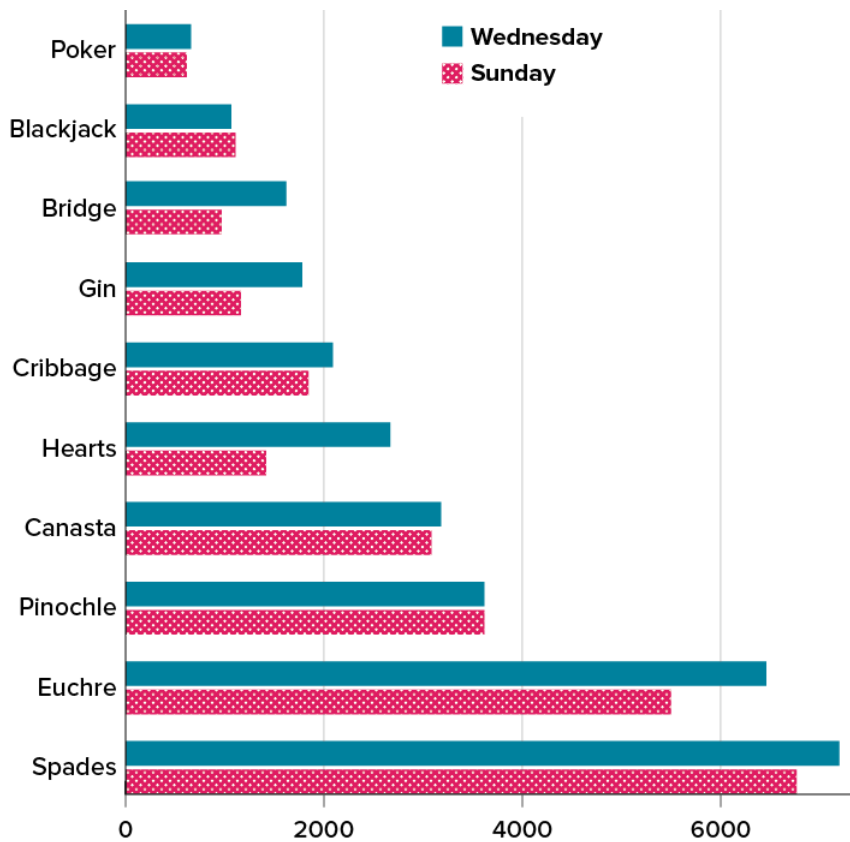


Comparing Distributions

Often we need to compare the results of different surveys, or of different conditions within the same overall survey. In this case, we are comparing the “distributions” of responses between the surveys or conditions. Bar charts are often excellent for illustrating differences between two distributions. [Figure 2.3](#) shows the number of people playing card games at the Yahoo web site on a Sunday and on a Wednesday in the spring of 2001. We see that there were more players overall on Wednesday compared to Sunday. The number of people playing Pinochle was nonetheless the same on these two days. In contrast, there were about twice as many people playing Hearts on Wednesday as on Sunday. Facts like these emerge clearly from a well-designed bar chart.

The bars in [Figure 2.3](#) are oriented horizontally rather than vertically. The horizontal format is useful when you have many categories because there is more room for the category labels. We’ll have [more to say about bar charts](#) when we consider numerical quantities later in this chapter.

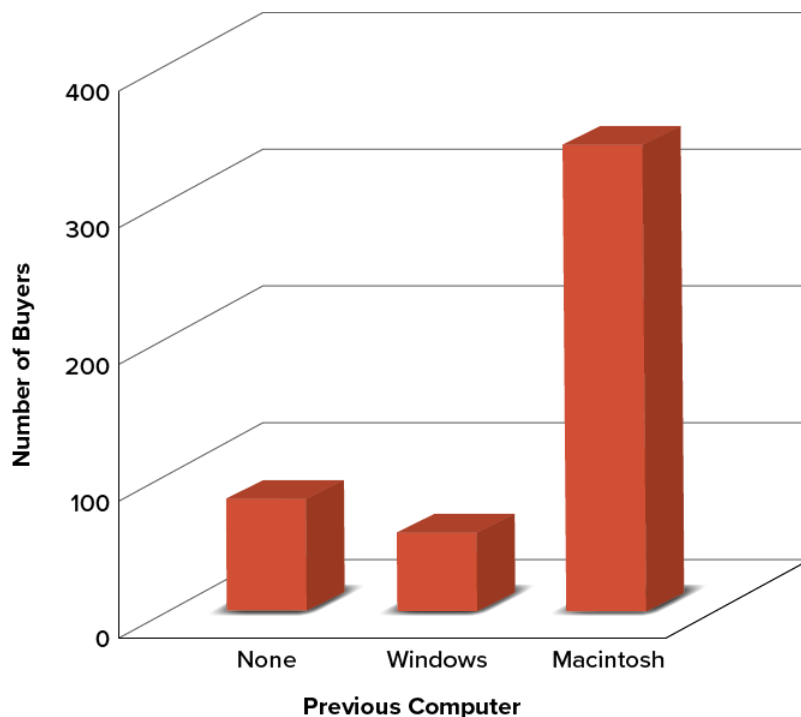
Figure 2.3. A bar chart of the number of people playing different card games on Sunday and Wednesday. (“[Card Game Bar Chart](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Some Graphical Mistakes to Avoid

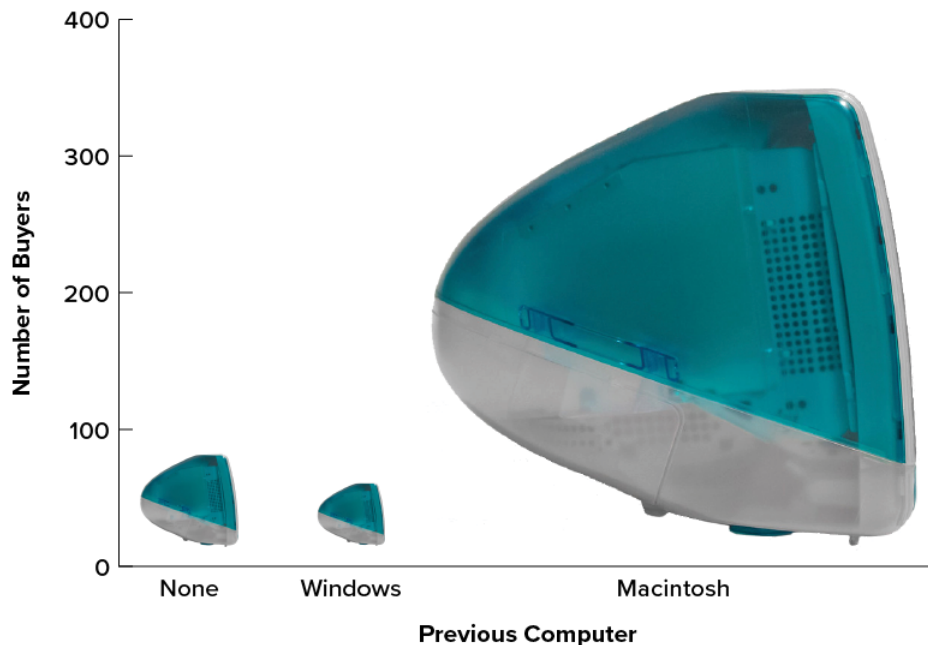
Don't get fancy! People sometimes add features to graphs that don't help to convey their information. For example, three-dimensional bar charts such as the one shown in [Figure 2.4](#) are usually not as effective as their two-dimensional counterparts.

Figure 2.4. A three-dimensional version of [Figure 2.2](#). Charts like this are less effective. ("Mac Bar Chart 3D" by Judy Schmitt is licenced under [CC BY-NC-SA 4.0](#).)



Here is another way that fanciness can lead to trouble. Instead of plain bars, it is tempting to substitute meaningful images. For example, [Figure 2.5](#) presents the iMac data using pictures of computers. The heights of the pictures accurately represent the number of buyers, yet [Figure 2.5](#) is misleading because the viewer's attention will be captured by areas. The areas can exaggerate the size differences between the groups. In terms of percentages, the ratio of previous Macintosh owners to previous Windows owners is about 6 to 1. But the ratio of the two areas in [Figure 2.5](#) is about 35 to 1. A biased person wishing to hide the fact that many Windows owners purchased iMacs would be tempted to use [Figure 2.5](#) instead of [Figure 2.2](#)!

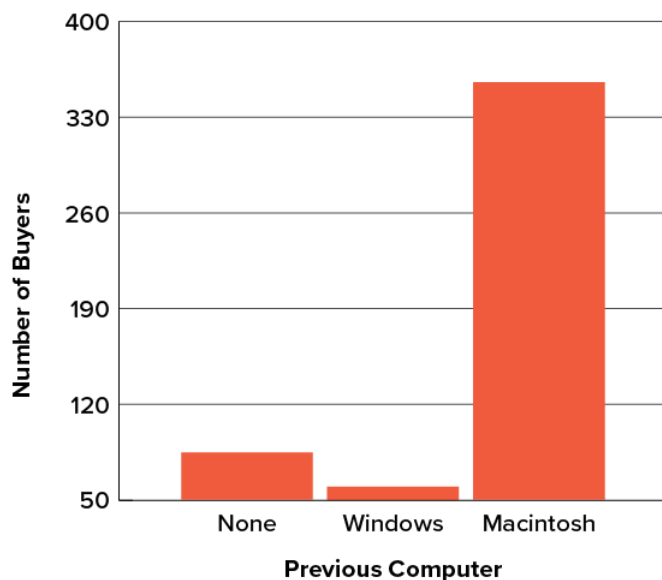
Figure 2.5. A redrawing of [Figure 2.2](#) with a lie factor greater than 8. (“Mac Bar Chart Lie Factor” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#). “Apple iMac G3 (1998)” by albaco/Flickr is licensed under [CC BY-NC-SA 2.0](#); image was brightened and background was removed.)



Edward Tufte coined the term lie factor to refer to the ratio of the size of the effect shown in a graph to the size of the effect shown in the data. He suggests that lie factors greater than 1.05 or less than 0.95 produce unacceptable distortion.

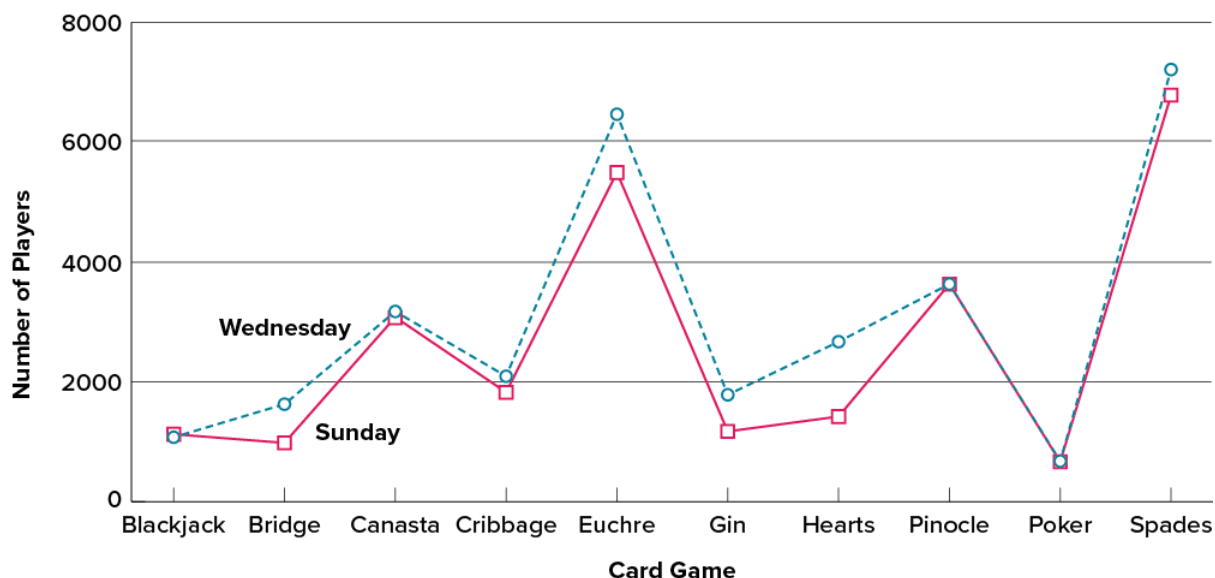
Another distortion in bar charts results from setting the baseline to a value other than zero. The baseline is the bottom of the y-axis, representing the least number of cases that could have occurred in a category. Normally, but not always, this number should be zero. [Figure 2.6](#) shows the iMac data with a baseline of 50. Once again, the differences in areas suggests a different story than the true differences in percentages. The number of Windows-switchers seems minuscule compared to its true value of 12%.

Figure 2.6. A redrawing of [Figure 2.2](#) with a baseline of 50. (“Mac Bar Chart Baseline 50” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Finally, we note that it is a serious mistake to use a line graph when the x-axis contains merely qualitative variables. A line graph is essentially a bar graph with the tops of the bars represented by points joined by lines (the rest of the bar is suppressed). Figure 2.7 inappropriately shows a line graph of the card game data from Yahoo that was presented in Figure 2.3. The drawback to Figure 2.7 is that it gives the false impression that the games are naturally ordered in a numerical way when, in fact, they are ordered alphabetically.

Figure 2.7. A line graph used inappropriately to depict the number of people playing different card games on Sunday and Wednesday. (“Line Chart Inappropriately Used” by Judy Schmitt is licensed under CC BY-NC-SA 4.0.)



Summary

Pie charts and bar charts can both be effective methods of portraying qualitative data. Bar charts are better when there are more than just a few categories and for comparing two or more distributions. Be careful to avoid creating misleading graphs.

Graphing Quantitative Variables

As discussed in the section on variables in Chapter 1, quantitative variables are variables measured on a numeric scale. Height, weight, response time, subjective rating of pain, temperature, and score on an exam are all examples of quantitative variables. Quantitative variables are distinguished from qualitative variables (sometimes called categorical variables or nominal variables), such as favorite color, religion, city of birth, and favorite sport, in which there is no ordering or measuring involved.

There are many types of graphs that can be used to portray distributions of quantitative variables. The upcoming sections cover the following types of graphs: (1) stem-and-leaf displays, (2) histograms, (3) frequency polygons, (4) box plots, (5) bar charts, (6) line graphs, (7) dot plots, and (8) scatter plots (discussed in [Chapter 12](#)). Some graph types, such as stem-and-leaf displays, are best-suited for small to moderate amounts of data, whereas others, such as histograms, are best-suited for large amounts of data. Graph types such as box plots are good at depicting differences between distributions. Scatter plots are used to show the relationship between two variables.

Stem-and-Leaf Displays

A stem-and-leaf display is a graphical method of displaying data. It is particularly useful when your data are not too numerous. In this section, we will explain how to construct and interpret this kind of graph.

As usual, we will start with an example. Consider [Figure 2.8](#), which shows the number of touchdown passes (TD passes) thrown by each of the 31 teams in the National Football League during the 2000 season.

Figure 2.8. Number of touchdown passes. (“[Touchdown Passes Raw Data](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

```
37, 33, 33, 32, 29, 28,
28, 23, 22, 22, 22, 21,
21, 21, 20, 20, 19, 19,
18, 18, 18, 18, 16, 15,
14, 14, 14, 12, 12, 9, 6
```

A stem-and-leaf display of the data is shown in [Figure 2.9](#). The left portion of [Figure 2.9](#) contains the stems. They are the numbers 3, 2, 1, and 0, arranged as a column to the left of the bars. Think of these numbers as 10s digits. A stem of 3, for example, can be used to represent the 10s digit in any of the numbers from 30 to 39. The numbers to the right of the bar are leaves, and they represent the 1s digits. Every leaf in the graph therefore stands for the result of adding the leaf to 10 times its stem.

Figure 2.9. Stem-and-leaf display of the number of touchdown passes. (“[Touchdown Passes Stem and Leaf](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

```
3 | 2337
2 | 001112223889
1 | 2244456888899
0 | 69
```

To make this clear, let us examine [Figure 2.9](#) more closely. In the top row, the four leaves to the right of stem 3 are 2, 3, 3, and 7. Combined with the stem, these leaves represent the numbers 32, 33, 33, and 37, which are the numbers of TD passes for the first four teams in [Figure 2.8](#). The next row has a stem of 2 and 12 leaves. Together, they represent 12 data points, namely, two occurrences of 20 TD passes, three occurrences of 21 TD passes, three occurrences of 22 TD passes, one occurrence of 23 TD passes, two occurrences of 28 TD passes, and one occurrence of 29 TD passes. We leave it to you to figure out what the third row represents. The fourth row has a stem of 0 and two leaves. It stands for the last two entries in [Figure 2.8](#), namely 9 TD passes and 6 TD passes. (The latter two numbers may be thought of as 09 and 06.)

One purpose of a stem-and-leaf display is to clarify the shape of the distribution. You can see many facts about TD passes more easily in [Figure 2.9](#) than in [Figure 2.8](#). For example, by looking at the stems and the shape of the plot, you can tell that most of the teams had between 10 and 29 passing TDs, with a few having more and a few having less. The precise numbers of TD passes can be determined by examining the leaves.

We can make our figure even more revealing by splitting each stem into two parts. [Figure 2.10](#) shows how to do this. The top row is reserved for numbers from 35 to 39 and holds only the 37 TD passes made by the first team in [Figure 2.8](#). The second row is reserved for the numbers from 30 to 34 and holds the 32, 33, and 33 TD passes made by the next three teams in the table. You can see for yourself what the other rows represent.

Figure 2.10. Stem-and-leaf display with the stems split in two. (“[Touchdown Passes Split Stem and Leaf](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

3		7
3		233
2		889
2		001112223
1		56888899
1		22444
0		69

Figure 2.10 is more revealing than Figure 2.9 because the latter figure lumps too many values into a single row. Whether you should split stems in a display depends on the exact form of your data. If rows get too long with single stems, you might try splitting them into two or more parts.

There is a variation of stem-and-leaf displays that is useful for comparing distributions. The two distributions are placed back to back along a common column of stems. The result is a back-to-back stem-and-leaf display, as shown in Figure 2.11. It compares the numbers of TD passes in the 1998 and 2000 seasons. The stems are in the middle, the leaves to the left are for the 1998 data, and the leaves to the right are for the 2000 data. For example, the second-to-last row shows that in 1998 there were teams with 11, 12, and 13 TD passes, and in 2000 there were two teams with 12 and three teams with 14 TD passes.

Figure 2.11. Back-to-back stem-and-leaf display. The left side shows the 1998 TD data and the right side shows the 2000 TD data. (“Touchdown Passes Back-to-Back Stem and Leaf” by Judy Schmitt is licensed under CC BY-NC-SA 4.0.

11	4	
	3	7
332	3	233
8865	2	889
44331110	2	001112223
987776665	1	56888899
321	1	22444
7	0	69

Figure 2.11 helps us see that the two seasons were similar, but that only in 1998 did any teams throw more than 40 TD passes.

There are two things about the football data that make them easy to graph with stems and leaves. First, the data are limited to whole numbers that can be represented with a one-digit stem and a one-digit leaf. Second, all the numbers are positive. If the data include numbers with three or more digits, or contain decimals, they can be rounded to two-digit accuracy. Negative values are also easily handled. Let us look at another example.

Figure 2.12 shows data from the Weapons and Aggression case study developed at Rice University. Each value is the mean difference over a series of trials between the times it took an experimental subject to name aggressive words (like punch) under two conditions. In one condition, the words were preceded by a non-weapon word such as bug. In the second condition, the same words were preceded by a weapon word such as gun or knife. The issue addressed by the experiment was whether a preceding weapon word would speed up (or prime) pronunciation of the aggressive word compared to a non-weapon priming word. A positive difference implies greater priming of the aggressive word by the weapon word. Negative differences imply that the priming by the weapon word was less than for a neutral word.

Figure 2.12. The effects of priming (in thousandths of a second). (“Priming Effects Raw Data” by Judy Schmitt is licensed under CC BY-NC-SA 4.0.)

43.2, 42.9, 35.6, 25.6,
25.4, 23.6, 20.5, 19.9,
14.4, 12.7, 11.3, 10.2,
10.0, 9.1, 7.5, 5.4, 4.7,
3.8, 2.1, 1.2, -0.2,
-6.3, -6.7, -8.8, -10.4,
-10.5, -14.9, -14.9,
-15.0, -18.5, -27.4

You see that the numbers range from 43.2 to -27.4 . The first value indicates that one subject was 43.2 milliseconds faster pronouncing aggressive words when they were preceded by weapon words than when preceded by neutral words. The value -27.4 indicates that another subject was 27.4 milliseconds slower pronouncing aggressive words when they were preceded by weapon words.

The data are displayed with stems and leaves in Figure 2.13. Since stem-and-leaf displays can only portray two whole digits (one for the stem and one for the leaf) the numbers are first rounded. Thus, the value 43.2 is rounded to 43 and represented with a stem of 4 and a leaf of 3. Similarly, 42.9 is rounded to 43. To represent negative numbers, we simply use negative stems. For example, the bottom row of the figure represents the number -27 . The second-to-last row represents the numbers -10 , -10 , -15 , etc. Once again, we have rounded the original values from Figure 2.12.

Figure 2.13. Stem-and-leaf display with negative numbers and rounding. (“Priming Effects Stem and Leaf” by Judy Schmitt is licensed under CC BY-NC-SA 4.0.)

4 | 33
3 | 6
2 | 00456
1 | 00134
0 | 1245589
-0 | 0679
-1 | 005559
-2 | 7

Observe that the figure contains a row headed by “0” and another headed by “-0.” The stem of 0 is for numbers between 0 and 9, whereas the stem of -0 is for numbers between 0 and -9 . For example, the fifth row of the table holds the numbers 1, 2, 4, 5, 5, 8, 9 and the sixth row holds 0, -6 , -7 , and -9 . Values that are exactly 0 before rounding should be split as evenly as possible between the “0” and “-0” rows. In Figure 2.12, none of the values are 0 before rounding. The “0” that appears in the “-0” row comes from the original value of -0.2 in the table.

Although stem-and-leaf displays are unwieldy for large datasets, they are often useful for datasets with up to 200 observations. Figure 2.14 portrays the distribution of populations of 185 U.S. cities in 1998. To be included, a city had to have between 100,000 and 500,000 residents.

Figure 2.14. Stem-and-leaf display of populations of 185 U.S. cities with populations between 100,000 and 500,000 in 1988. Stems represent units of 100,000, and leaves represent units of 10,000. (“US Populations Stem and Leaf” by Judy Schmitt is licensed under CC BY-NC-SA 4.0.)

```

4 | 899
4 | 6
4 | 4455
4 | 333
4 | 01
3 | 99
3 | 677777
3 | 55
3 | 223
3 | 111
2 | 8899
2 | 666667
2 | 444455
2 | 22333
2 | 000000
1 | 8888888888889999999999
1 | 666666777777
1 | 444444444444555555555555
1 | 2222222222222222222233333333
1 | 0000000000000000111111111111111111111111

```

Since a stem-and-leaf plot shows only two-place accuracy, we had to round the numbers to the nearest 10,000. For example the largest number (493,559) was rounded to 490,000 and then plotted with a stem of 4 and a leaf of 9. The fourth highest number (463,201) was rounded to 460,000 and plotted with a stem of 4 and a leaf of 6. Thus, the stems represent units of 100,000, and the leaves represent units of 10,000. Notice that each stem value is split into five parts: 0 to 1, 2 to 3, 4 to 5, 6 to 7, and 8 to 9.

Whether your data can be suitably represented by a stem-and-leaf display depends on whether they can be rounded without loss of important information. Also, their extreme values must fit into two successive digits, as the data in [Figure 2.14](#) fit into the 10,000 and 100,000 places (for leaves and stems, respectively). Deciding what kind of graph is best suited to displaying your data thus requires good judgment. Statistics is not just recipes!

Histograms

A histogram is a graphical method for displaying the shape of a distribution. It is particularly useful when there are a large number of observations. We begin with an example consisting of the scores of 642 students on a psychology test. The test consists of 197 items, each graded as “correct” or “incorrect.” The students’ scores ranged from 46 to 167.

The first step is to create a frequency table. Unfortunately, a simple frequency table would be too big, containing over 100 rows. To simplify the table, we group scores together as shown in [Table 2.2](#).

Table 2.2. Grouped frequency distribution of psychology test scores.

Interval's Lower Limit	Interval's Upper Limit	Class Frequency
39.5	49.5	3
49.5	59.5	10
59.5	69.5	53
69.5	79.5	107
79.5	89.5	147
89.5	99.5	130
99.5	109.5	78
109.5	119.5	59

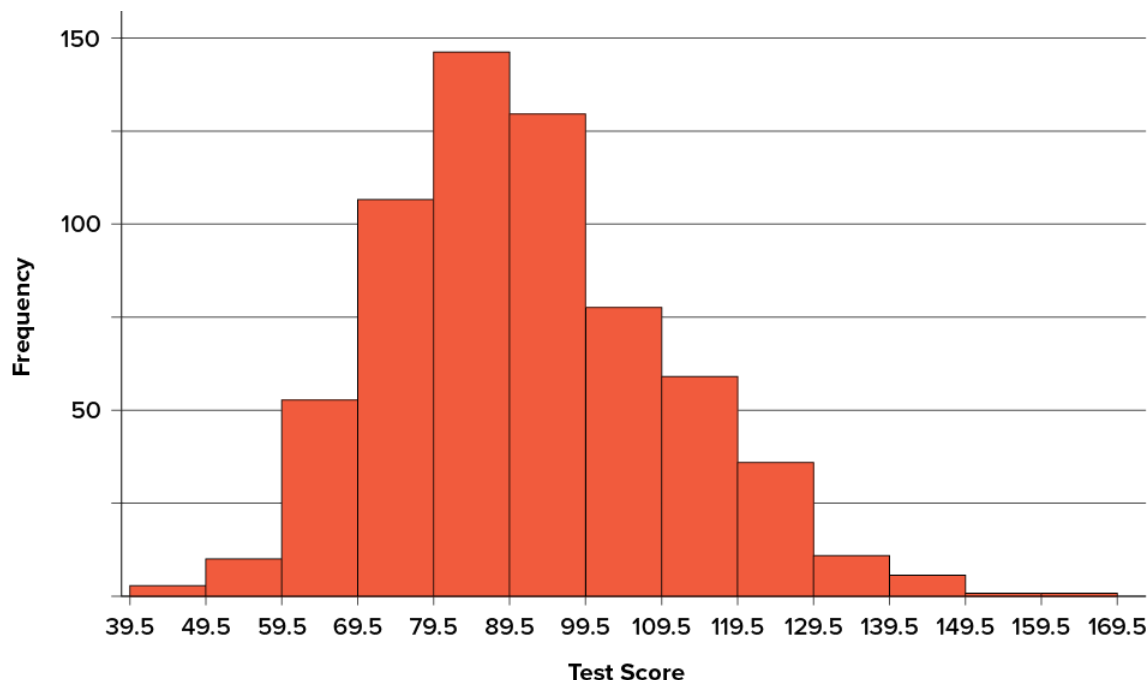
Interval's Lower Limit	Interval's Upper Limit	Class Frequency
119.5	129.5	36
129.5	139.5	11
139.5	149.5	6
149.5	159.5	1
159.5	169.5	1

To create this table, the range of scores was broken into intervals, called class intervals. The first interval is from 39.5 to 49.5, the second from 49.5 to 59.5, etc. Next, the number of scores falling into each interval was counted to obtain the class frequencies. There are 3 scores in the first interval, 10 in the second, etc.

Class intervals of width 10 provide enough detail about the distribution to be revealing without making the graph too “choppy.” More information on choosing the widths of class intervals is presented later in this section. Placing the limits of the class intervals midway between two numbers (e.g., 49.5) ensures that every score will fall in an interval rather than on the boundary between intervals.

In a histogram, the class frequencies are represented by bars. The height of each bar corresponds to its class frequency. A histogram of these data is shown in [Figure 2.15](#).

Figure 2.15. Histogram of scores on a psychology test. (“[Psychology Test Scores Histogram](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



The histogram makes it plain that most of the scores are in the middle of the distribution, with fewer scores in the extremes. You can also see that the distribution is not symmetric: the scores extend farther to the right than they do to the left. The distribution is therefore said to be skewed. (We’ll have more to say about shapes of distributions in [Chapter 3](#).)

In our example, the observations are whole numbers. Histograms can also be used when the scores are measured on a more continuous scale such as the length of time (in milliseconds) required to perform a task. In this case, there is no need to worry about fence sitters since they are improbable. (It would be quite a coincidence for a task to require exactly 7 seconds, measured to the nearest thousandth of a second.) We are therefore free to choose whole numbers as boundaries for our class intervals, for example, 4000, 5000, etc. The class frequency is then the number of observations that are greater than or equal to the lower bound, and strictly less than the upper bound. For example, one interval might hold times from 4000 to 4999 milliseconds. Using whole

numbers as boundaries avoids a cluttered appearance, and is the practice of many computer programs that create histograms. Note also that some computer programs label the middle of each interval rather than the end points.

Histograms can be based on relative frequencies instead of actual frequencies. Histograms based on relative frequencies show the proportion of scores in each interval rather than the number of scores. In this case, the y-axis runs from 0 to 1 (or somewhere in between if there are no extreme proportions). You can change a histogram based on frequencies to one based on relative frequencies by (a) dividing each class frequency by the total number of observations, and then (b) plotting the quotients on the y-axis (labeled as proportion).

There is more to be said about the widths of the class intervals, sometimes called bin widths. Your choice of bin width determines the number of class intervals. This decision, along with the choice of starting point for the first interval, affects the shape of the histogram. The best advice is to experiment with different choices of width, and to choose a histogram according to how well it communicates the shape of the distribution.

Frequency Polygons

Frequency polygons are a graphical device for understanding the shapes of distributions. They serve the same purpose as histograms, but are especially helpful for comparing sets of data. Frequency polygons are also a good choice for displaying cumulative frequency distributions.

To create a frequency polygon, start just as for histograms, by choosing a class interval. Then draw an x-axis representing the values of the scores in your data. Mark the middle of each class interval with a tick mark, and label it with the middle value represented by the class. Draw the y-axis to indicate the frequency of each class. Place a point in the middle of each class interval at the height corresponding to its frequency. Finally, connect the points. You should include one class interval below the lowest value in your data and one above the highest value. The graph will then touch the x-axis on both sides.

The frequency distribution of 642 psychology test scores, shown in [Table 2.3](#), was used to create the frequency polygon shown in [Figure 2.16](#).

Table 2.3. Frequency distribution of psychology test scores.

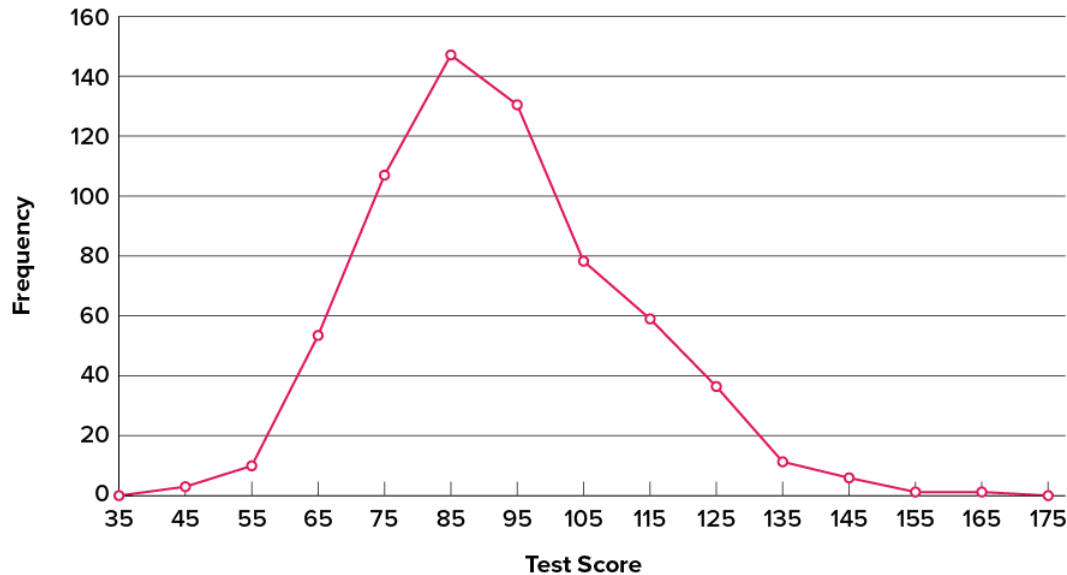
Lower Limit	Upper Limit	Count	Cumulative Count
29.5	39.5	0	0
39.5	49.5	3	3
49.5	59.5	10	13
59.5	69.5	53	66
69.5	79.5	107	173
79.5	89.5	147	320
89.5	99.5	130	450
99.5	109.5	78	528
109.5	119.5	59	587
119.5	129.5	36	623
129.5	139.5	11	634
139.5	149.5	6	640
149.5	159.5	1	641
159.5	169.5	1	642
169.5	170.5	0	642

The first label on the x-axis is 35. This represents an interval extending from 29.5 to 39.5. Since the lowest test score is 46, this interval has a frequency of 0. The point labeled 45 represents the interval from 39.5 to 49.5. There are three scores in this interval.

There are 147 scores in the interval that surrounds 85.

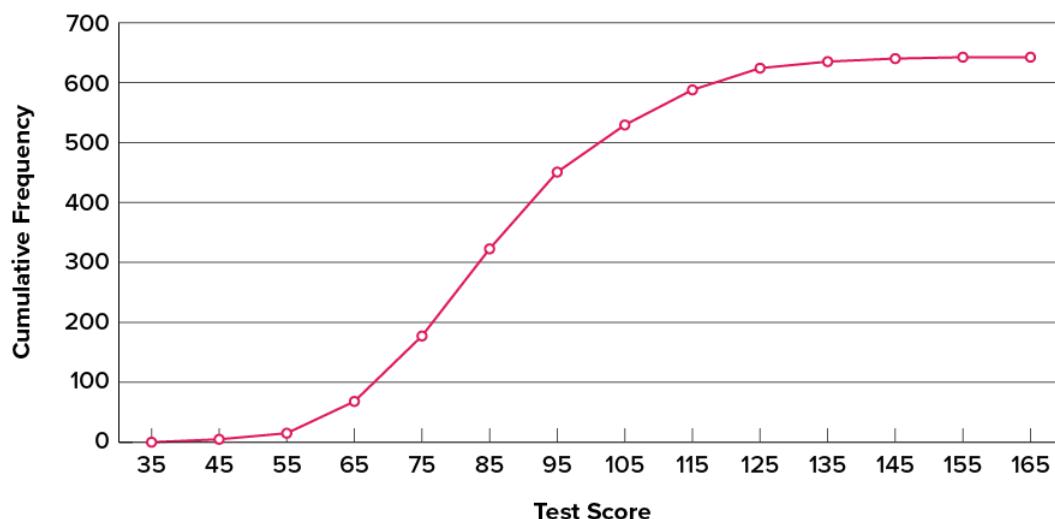
You can easily discern the shape of the distribution from [Figure 2.16](#). Most of the scores are between 65 and 115. It is clear that the distribution is not symmetric inasmuch as good scores (to the right) trail off more gradually than poor scores (to the left). In the terminology of [Chapter 3](#) (where we will study shapes of distributions more systematically), the distribution is skewed.

Figure 2.16. Frequency polygon for the psychology test scores. (“[Psychology Test Scores Frequency Polygon](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



A cumulative frequency polygon for the same test scores is shown in [Figure 2.17](#). The graph is the same as before except that the y value for each point is the number of students in the corresponding class interval plus all numbers in lower intervals. For example, there are no scores in the interval labeled “35,” three in the interval “45,” and 10 in the interval “55.” Therefore, the y value corresponding to “55” is 13. Since 642 students took the test, the cumulative frequency for the last interval is 642.

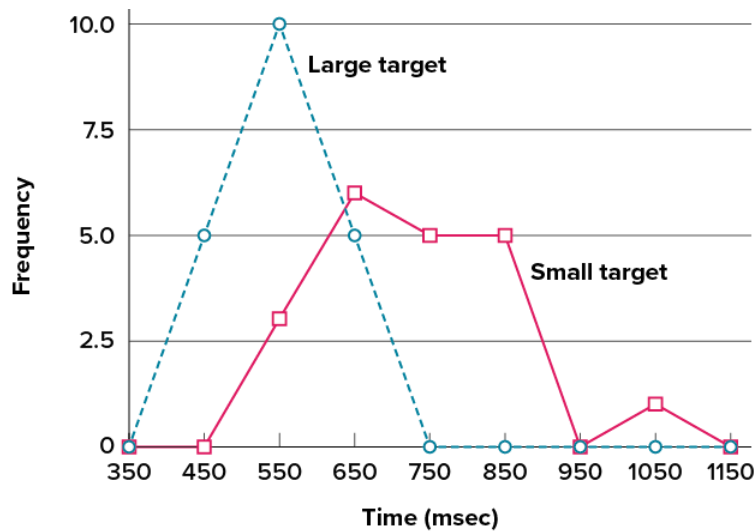
Figure 2.17. Cumulative frequency polygon for the psychology test scores. (“[Psychology Test Scores Cumulative Frequency Polygon](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Frequency polygons are useful for comparing distributions. This is achieved by overlaying the frequency polygons drawn for different datasets. [Figure 2.18](#) provides an example. The data come from a task in which the goal is to move a computer cursor to a target on the screen as fast as possible. On 20 of the trials, the target was a small rectangle; on the other 20, the target was a large rectangle. Time to reach the target was recorded on each trial. The two distributions (one for each target) are plotted together in

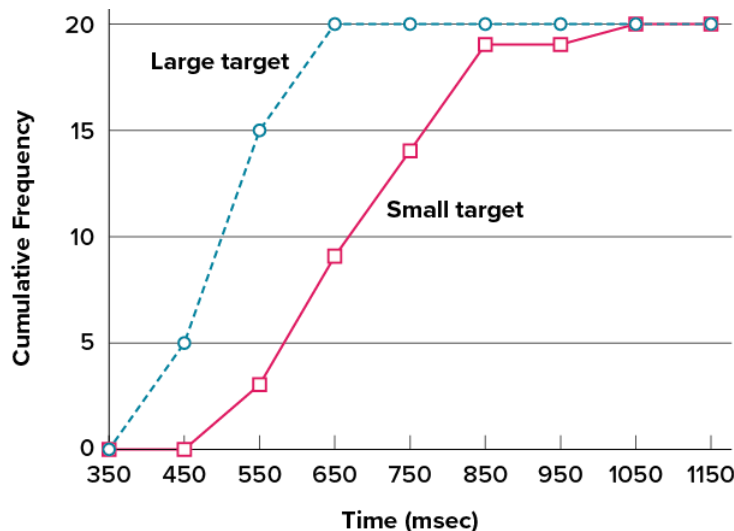
[Figure 2.18](#). The figure shows that, although there is some overlap in times, it generally took longer to move the cursor to the small target than to the large one.

Figure 2.18. Overlaid frequency polygons for the cursor task. (“[Cursor Task Frequency Polygons](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



It is also possible to plot two cumulative frequency distributions in the same graph. This is illustrated in [Figure 2.19](#) using the same data from the cursor task. The difference in distributions for the two targets is again evident.

Figure 2.19. Overlaid cumulative frequency polygons for the cursor task. (“[Cursor Task Cumulative Frequency Polygons](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Box Plots

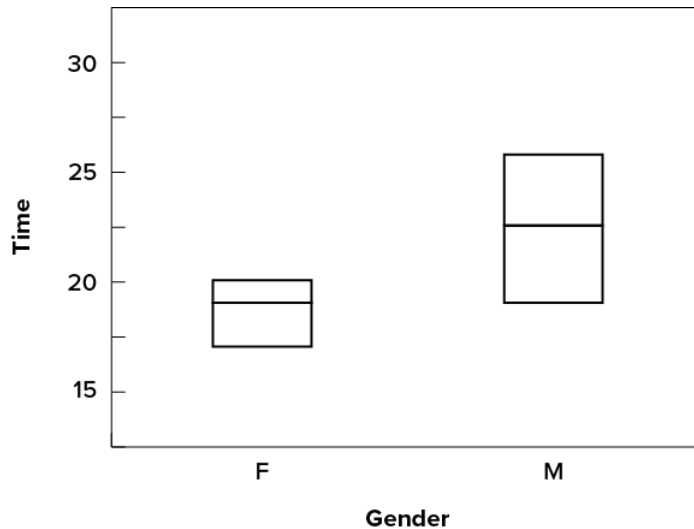
We have already discussed techniques for visually representing data (see histograms and frequency polygons). In this section we present another important graph, called a box plot. Box plots are useful for identifying outliers and for comparing distributions. We will explain box plots with the help of data from an in-class experiment. Students in Introductory Statistics were presented with a page containing 30 colored rectangles. Their task was to name the colors as quickly as possible. Their times (in seconds) were recorded. We’ll compare the scores for the 16 men and 31 women who participated in the experiment by making separate box plots for each gender. Such a display is said to involve parallel box plots. The data for the women in our sample are shown in [Figure 2.20](#).

Figure 2.20. Women’s times. (“[Women’s Times Raw Data](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

14, 15, 16, 16, 17, 17, 17, 17, 17, 18, 18, 18, 18, 18, 18, 19, 19, 19
20, 20, 20, 20, 20, 20, 21, 21, 22, 23, 24, 24, 29

There are several steps in constructing a box plot. The first relies on the 25th, 50th, and 75th percentiles in the distribution of scores. Figure 2.21 shows how these three statistics are used. For each gender we draw a box extending from the 25th percentile to the 75th percentile. The 50th percentile is drawn inside the box. Therefore, the bottom of each box is the 25th percentile, the top is the 75th percentile, and the line in the middle is the 50th percentile.

Figure 2.21. The first step in creating box plots. (“Box Plot First Step” by Judy Schmitt is licensed under CC BY-NC-SA 4.0.)



For the data reflecting the women’s times, the 25th percentile is 17, the 50th percentile is 19, and the 75th percentile is 20. For the men (whose data are not shown), the 25th percentile is 19, the 50th percentile is 22.5, and the 75th percentile is 25.5.

Before proceeding, the terminology in Table 2.4 is helpful.

Table 2.4. Box plot terms and values for women’s times.

Name	Formula	Value
Upper Hinge	75th percentile	20
Lower Hinge	25th percentile	17
H-Spread	Upper Hinge – Lower Hinge	

3
Step
 $1.5 \times \text{H-Spread}$
4.5
Upper Inner Fence
Upper Hinge + 1 Step
24.5
Lower Inner Fence
Lower Hinge – 1 Step
12.5

Upper Outer Fence

Upper Hinge + 2 Steps

29

Lower Outer Fence

Lower Hinge – 2 Steps

8

Upper Adjacent

Largest value below Upper Inner Fence

24

Lower Adjacent

Smallest value above Lower Inner Fence

14

Outside Value

A value beyond an Inner Fence but not beyond an Outer Fence

29

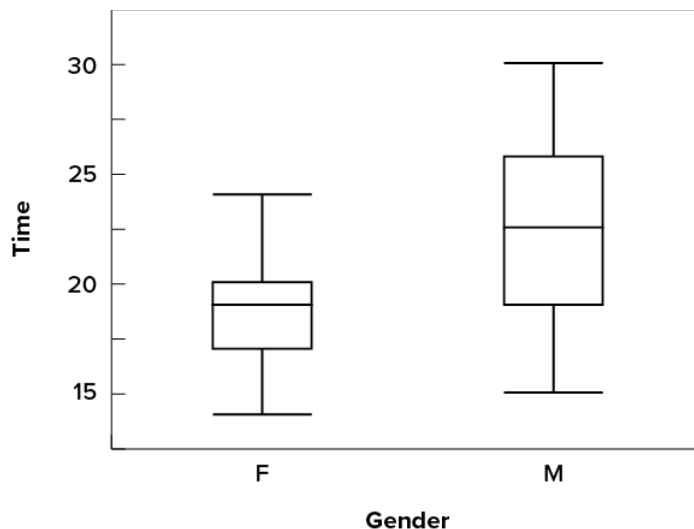
Far Out Value

A value beyond an Outer Fence

None

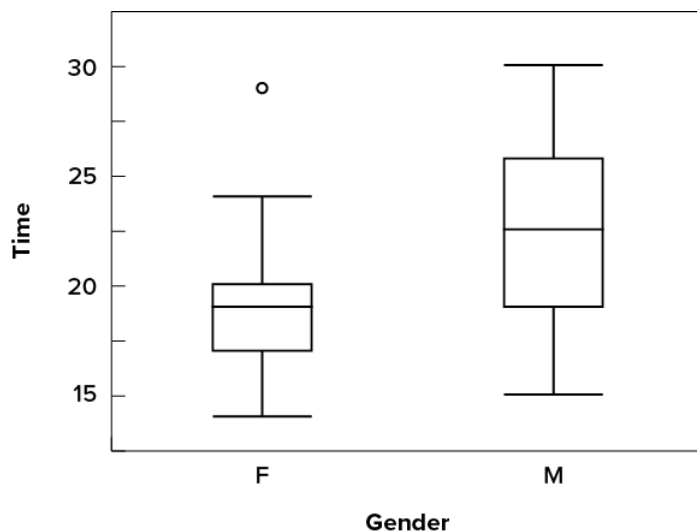
Continuing with the box plots, we put “whiskers” above and below each box to give additional information about the spread of data. Whiskers are vertical lines that end in a horizontal stroke. Whiskers are drawn from the upper and lower hinges to the upper and lower adjacent values (24 and 14 for the women’s data), as shown in [Figure 2.22](#).

Figure 2.22. The box plots with the whiskers drawn. (“[Box Plot Whiskers](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Although we don’t draw whiskers all the way to outside or far out values, we still wish to represent them in our box plots. This is achieved by adding additional marks beyond the whiskers. Specifically, outside values are indicated by small circles, and far out values are indicated by asterisks (*). In our data, there are no far-out values and just one outside value. This outside value of 29 is for the women and is shown in [Figure 2.23](#).

Figure 2.23. The box plots with the outside value shown. (“[Box Plot Outside Value](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



There is one more mark to include in box plots (although sometimes it is omitted). We indicate the mean score for a group by inserting a plus sign. Figure 2.24 shows the result of adding means to our box plots.

Figure 2.24. The completed box plots. (“Box Plot Mean Scores” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0.](https://creativecommons.org/licenses/by-nc-sa/4.0/))

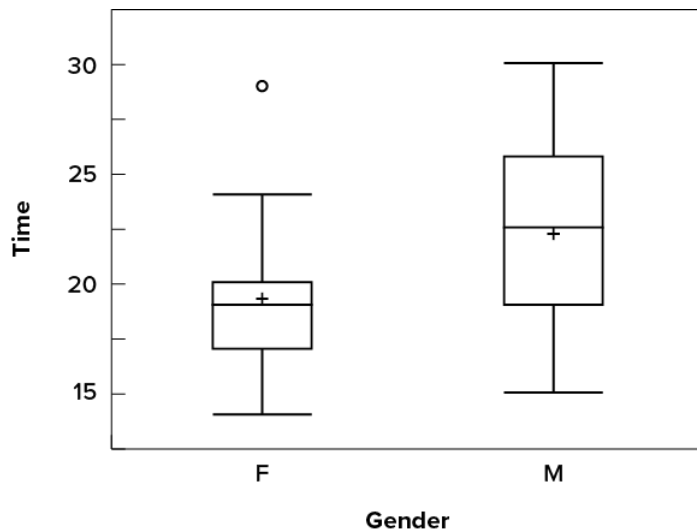
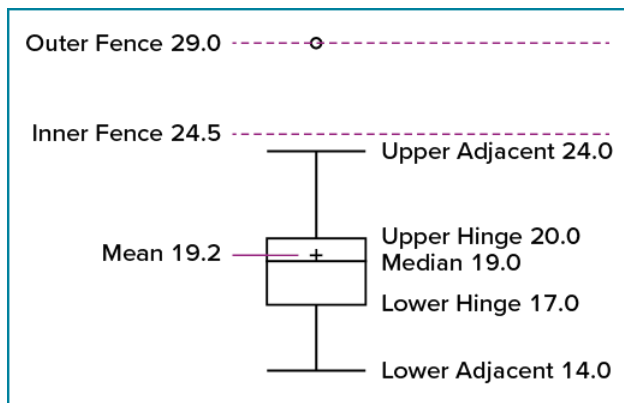


Figure 2.24 provides a revealing summary of the data. Since half the scores in a distribution are between the hinges (recall that the hinges are the 25th and 75th percentiles), we see that half the women’s times are between 17 and 20 seconds whereas half the men’s times are between 19 and 25.5 seconds. We also see that women generally named the colors faster than the men did, although one woman was slower than almost all of the men. Figure 2.25 shows the box plot for the women’s data with detailed labels.

Figure 2.25. The box plots for the women’s data with detailed labels. (“Women’s Data Labeled Box Plot” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0.](https://creativecommons.org/licenses/by-nc-sa/4.0/))



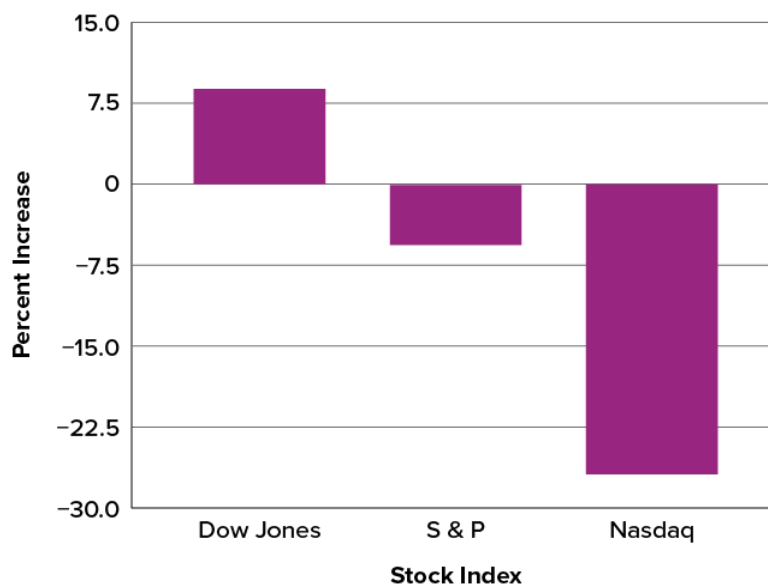
Box plots provide basic information about a distribution. For example, a distribution with a positive skew would have a longer whisker in the positive direction than in the negative direction. A larger mean than median would also indicate a positive skew. Box plots are good at portraying extreme values and are especially good at showing differences between distributions. However, many of the details of a distribution are not revealed in a box plot; to examine these details one should create a histogram and/or a stem-and-leaf display.

Bar Charts

In the [section on qualitative variables](#), we saw how bar charts could be used to illustrate the frequencies of different categories. For example, as we saw earlier in this chapter, the bar chart shown in [Figure 2.2](#) shows how many purchasers of iMac computers were previous Macintosh users, previous Windows users, and new computer purchasers.

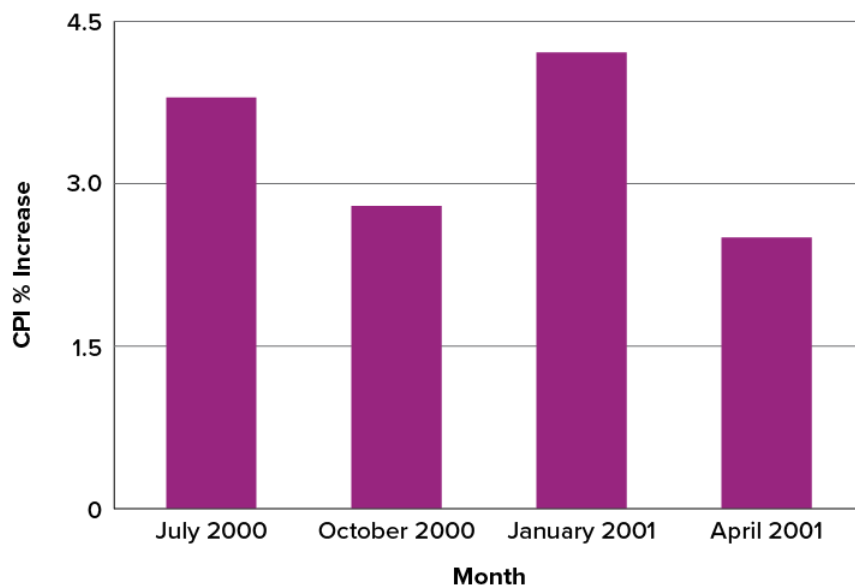
In this section we show how bar charts can be used to present other kinds of quantitative information, not just frequency counts. The bar chart in [Figure 2.26](#) shows the percent increases in the Dow Jones, Standard & Poor 500 (S&P), and Nasdaq stock indexes from May 24, 2000, to May 24, 2001. Notice that both the S&P and the Nasdaq had “negative increases” which means that they decreased in value. In this bar chart, the y-axis is not frequency but rather the signed quantity percentage increase.

Figure 2.26. Percent increase in three stock indexes from May 24, 2000, to May 24, 2001. (“[Percent Increase in Stock Indexes](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



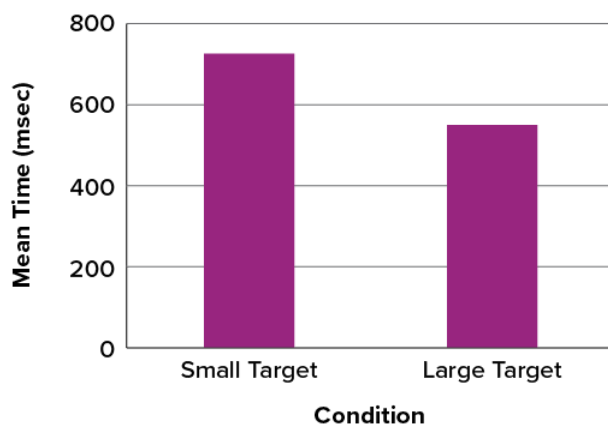
Bar charts are particularly effective for showing change over time. [Figure 2.27](#), for example, shows the percent increase in the Consumer Price Index (CPI) over four three-month periods. The fluctuation in inflation is apparent in the graph.

Figure 2.27. Percent change in the CPI over time. Each bar represents percent increase for the three months ending at the date indicated. (“[Percent Change in CPI](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



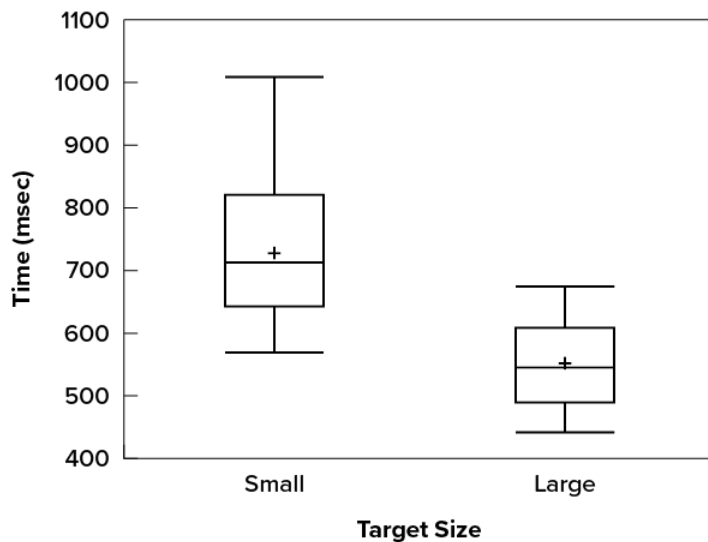
Bar charts are often used to compare the means of different experimental conditions. [Figure 2.28](#) shows the mean time it took one person to move the cursor to either a small target or a large target. On average, more time was required for small targets than for large ones.

Figure 2.28. Bar chart showing the means for the two conditions. (“[Means of Two Conditions](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Although bar charts can display means, we do not recommend them for this purpose. Box plots should be used instead since they provide more information than bar charts without taking up more space. For example, a box plot of the cursor-movement data is shown in [Figure 2.29](#). You can see that [Figure 2.29](#) reveals more about the distribution of movement times than does [Figure 2.28](#).

Figure 2.29. Box plots of times to move the cursor to the small and large targets. (“[Cursor Task Box Plot](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

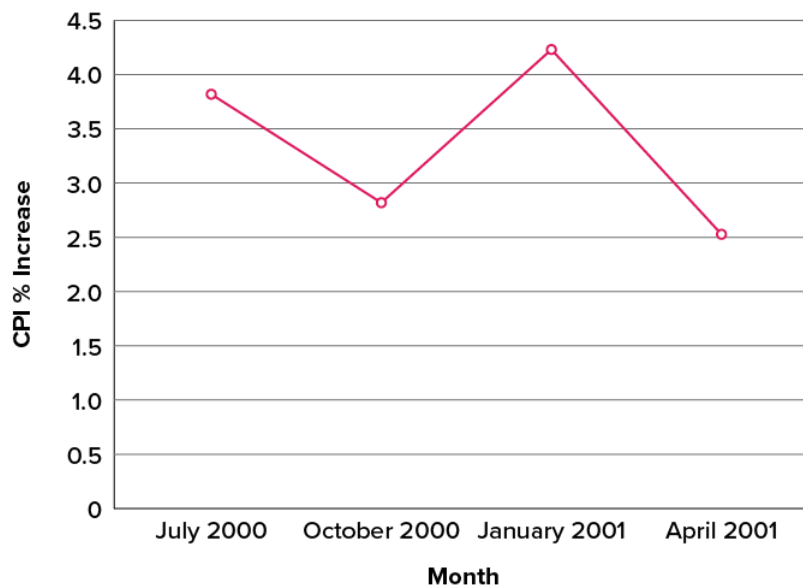


The [section on qualitative variables](#) presented earlier in this chapter discussed the use of bar charts for comparing distributions. Some common graphical mistakes were also noted. The earlier discussion applies equally well to the use of bar charts to display quantitative variables.

Line Graphs

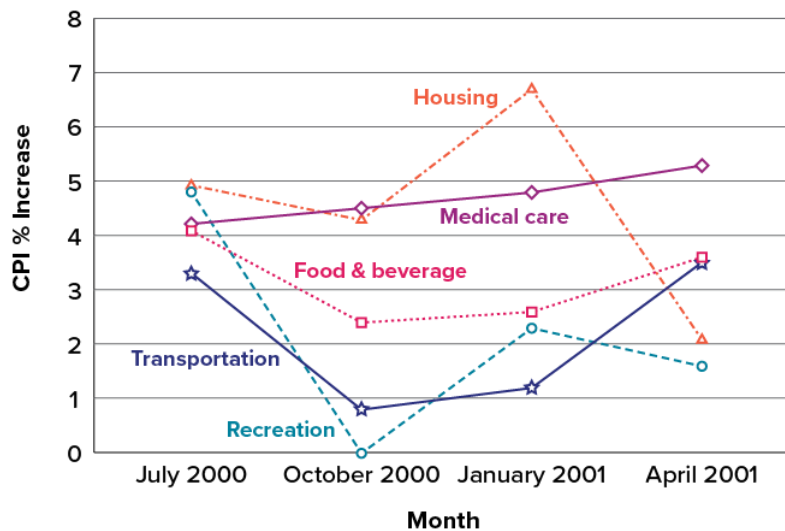
A line graph is a bar graph with the tops of the bars represented by points joined by lines (the rest of the bar is suppressed). For example, [Figure 2.27](#), which was presented in the section on bar charts, shows changes in the Consumer Price Index (CPI) over time. A line graph of these same data is shown in [Figure 2.30](#). Although the figures are similar, the line graph emphasizes the change from period to period.

Figure 2.30. A line graph of the percent change in the CPI over time. Each point represents percent increase for the three months ending at the date indicated. (“[Percent Change in CPI Line Graph](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)



Line graphs are appropriate only when both the x- and y-axes display ordered (rather than qualitative) variables. Although bar charts can also be used in this situation, line graphs are generally better at comparing changes over time. [Figure 2.31](#), for example, shows percent increases and decreases in five components of the CPI. The figure makes it easy to see that medical costs had a steadier progression than the other components. Although you could create an analogous bar chart, its interpretation would not be as easy.

Figure 2.31. A line graph of the percent change in five components of the CPI over time. (“Percent Change in CPI x5 Line Graph” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0.](https://creativecommons.org/licenses/by-nc-sa/4.0/))



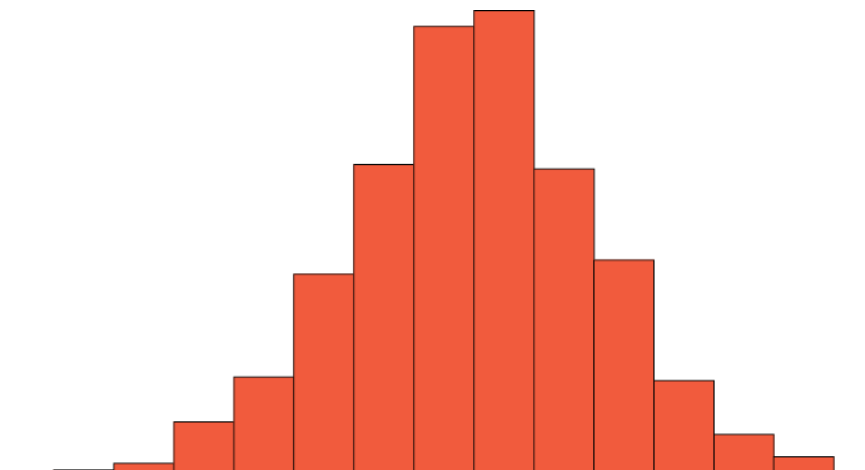
Let us stress that it is misleading to use a line graph when the x-axis contains merely qualitative variables. As we saw earlier in this chapter, [Figure 2.7](#) inappropriately shows a line graph of the card game data from Yahoo, discussed in the section on qualitative variables. The defect in [Figure 2.7](#) is that it gives the false impression that the games are naturally ordered in a numerical way.

The Shape of Distribution

Finally, it is useful to present discussion on how we describe the shapes of distributions, which we will revisit in [Chapter 3](#) to learn how different shapes affect our numerical descriptors of data and distributions.

The primary characteristic we are concerned about when assessing the shape of a distribution is whether the distribution is symmetrical or skewed. A symmetrical distribution, as the name suggests, can be cut down the center to form two mirror images. Although in practice we will never get a perfectly symmetrical distribution, we would like our data to be as close to symmetrical as possible for reasons we delve into in [Chapter 3](#). Many types of distributions are symmetrical, but by far the most common and pertinent distribution at this point is the normal distribution, shown in [Figure 2.32](#). Notice that although the symmetry is not perfect (for instance, the bar just to the right of the center is taller than the one just to the left), the two sides are roughly the same shape. The normal distribution has a single peak, known as the center, and two tails that extend out equally, forming what is known as a bell shape or bell curve.

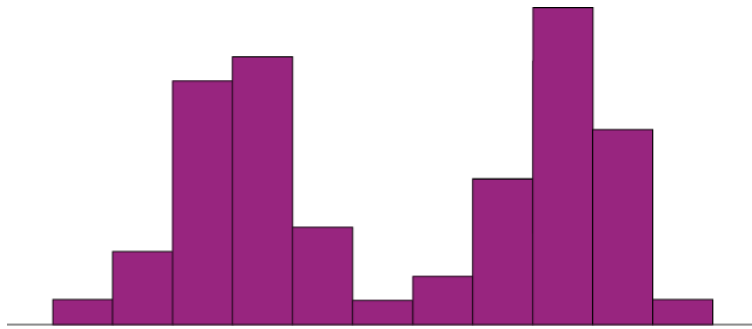
Figure 2.32. A symmetrical distribution. (“Symmetrical Distribution” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0.](https://creativecommons.org/licenses/by-nc-sa/4.0/))



Symmetrical distributions can also have multiple peaks. [Figure 2.33](#) shows a bimodal distribution, named for the two peaks that lie roughly symmetrically on either side of the center point. As we will see in [Chapter 3](#), this is not a particularly desirable

characteristic of our data, and, worse, this is a relatively difficult characteristic to detect numerically. Thus, it is important to visualize your data before moving ahead with any formal analyses.

Figure 2.33. A bimodal distribution. (“[Bimodal Distribution](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

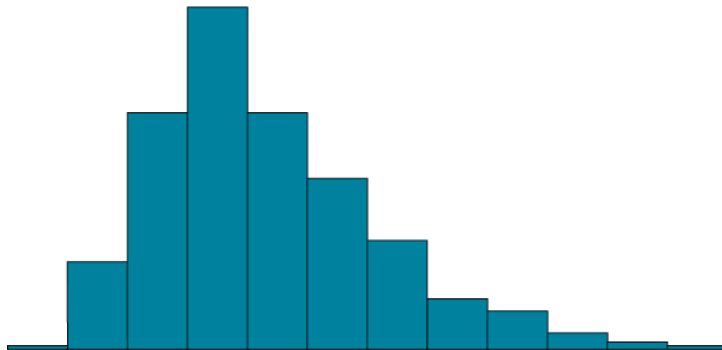


Distributions that are not symmetrical also come in many forms, more than can be described here. The most common asymmetry to be encountered is referred to as skew, in which one of the two tails of the distribution is disproportionately longer than the other. This property can affect the value of the averages we use in our analyses and make them an inaccurate representation of our data, which causes many problems.

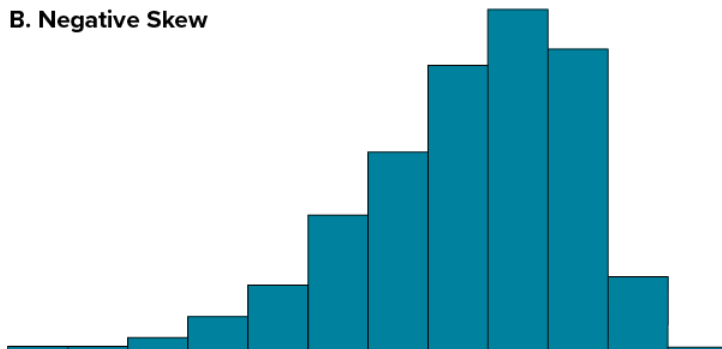
Skew can either be positive or negative (also known as right or left, respectively), based on which tail is longer. It is very easy to get the two confused at first; many students want to describe the skew by where the bulk of the data (larger portion of the histogram, known as the body) is placed, but the correct determination is based on which tail is longer. You can think of the tail as an arrow; whichever direction the arrow is pointing is the direction of the skew. [Figure 2.34](#) shows positive (right) and negative (left) skew, respectively.

Figure 2.34. Positively skewed (A) and negatively skewed (B) distributions. (“[Skewed Distributions](#)” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0](#).)

A. Positive Skew



B. Negative Skew



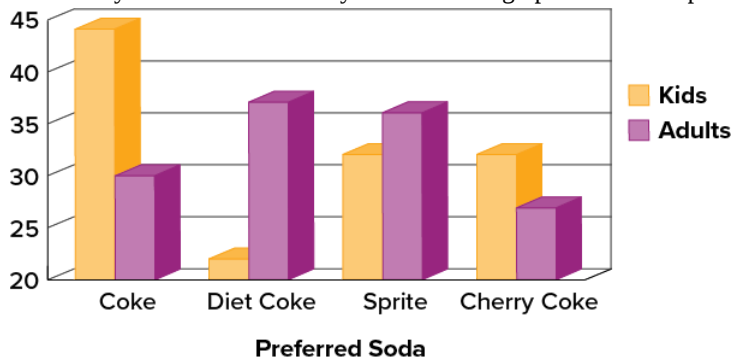
Exercises

1. Name some ways to graph quantitative variables and some ways to graph qualitative variables.
2. Given the following data, construct a pie chart and a bar chart. Which do you think is the more appropriate or useful way to display the data?

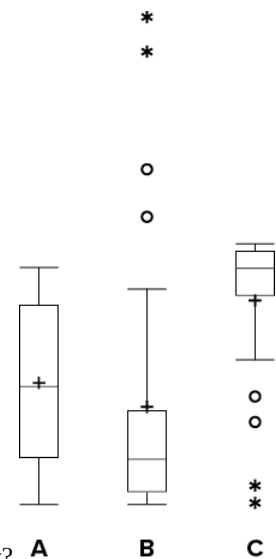
Favorite Movie Genre	Frequency
Comedy	14
Horror	9
Romance	8
Action	12

3. Pretend you are constructing a histogram for describing the distribution of salaries for individuals who are 40 years or older but not yet retired.
 1. What is on the y-axis? Explain.
 2. What is on the x-axis? Explain.
 3. What would be the probable shape of the salary distribution? Explain why.

4. A graph appears below showing the number of adults and children who prefer each type of soda. There were 130 adults and kids surveyed. Discuss some ways in which the graph could be improved.



("Improvable Bar Chart" by Judy Schmitt is licensed under [CC BY-NC-SA 4.0.](https://creativecommons.org/licenses/by-nc-sa/4.0/))



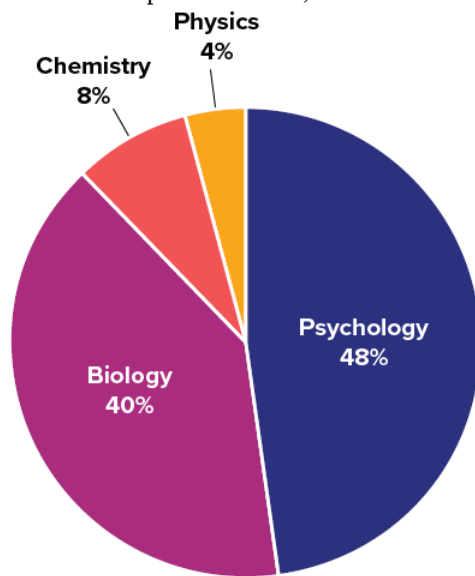
5. Which of the box plots on the graph has a large positive skew? Which has a large negative skew?

("Skewed Box Plots" by Judy Schmitt is licensed under [CC BY-NC-SA 4.0.](https://creativecommons.org/licenses/by-nc-sa/4.0/))

6. Create a histogram of the following data representing how many shows children said they watch each day:

Number of TV Shows	Frequency
0	2
1	18
2	36
3	7
4	3

7. Explain the differences between bar charts and histograms. When would each be used?
8. Draw a histogram of a distribution that is
1. Negatively skewed
 2. Symmetrical
 3. Positively skewed
9. Based on the pie chart below, which was made from a sample of 300 students, construct a frequency table of college majors.



(“College Majors Pie Chart” by Judy Schmitt is licensed under [CC BY-NC-SA 4.0.](https://creativecommons.org/licenses/by-nc-sa/4.0/))

10. Create a histogram of the following data. Label the tails and body and determine if it is skewed (and direction, if so) or symmetrical.

Hours Worked per Week	Proportion
0–10	4
10–20	8
20–30	11
30–40	51
40–50	12
50–60	9
60+	5

Answers to Odd-Numbered Exercises

1)

Qualitative variables are displayed using pie charts and bar charts. Quantitative variables are displayed as box plots, histograms, etc.

3)

[You do not need to draw the histogram, only describe it.]

1. The y-axis would show the frequency or proportion because this is always the case in histograms.
2. The x-axis would show income, because this is our quantitative variable of interest.
3. Because most income data are positively skewed, this histogram would likely be skewed positively too.

5)

Chart B has the positive skew because the outliers (dots and asterisks) are on the upper (higher) end; Chart C has the negative skew because the outliers are on the lower end.

7)

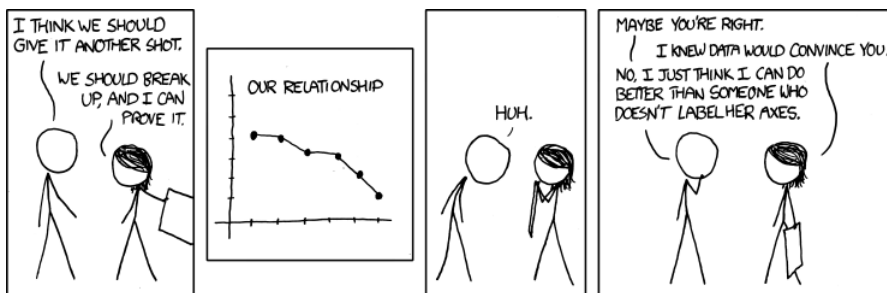
In bar charts, the bars do not touch; in histograms, the bars do touch. Bar charts are appropriate for qualitative variables, whereas histograms are better for quantitative variables.

9)

The frequency table appears below:

Major	Frequency
Psychology	144
Biology	120
Chemistry	24
Physics	12

“[Convincing](#)” by Randall Munroe/xkcd.com is licensed under [CC BY-NC 2.5](#).)



This page titled [1.2: Chapter 2- Describing Data Using Distributions and Graphs](#) is shared under a [not declared](#) license and was authored, remixed, and/or curated by [Linda R. Cote](#), [Rupa G. Gordon](#), [Chrislyn E. Randell](#), [Judy Schmitt](#), and [Helena Marvin](#).