

7.1: heoretical Models

Models, as discussed earlier, are an essential component in theory building. They simplify theoretical concepts, provide a precise way to evaluate relationships between variables, and serve as a vehicle for hypothesis testing. As discussed in Chapter 1, one of the central features of a theoretical model is the presumption of causality, and causality is based on three factors: time ordering (observational or theoretical), co-variation, and non-spuriousness. Of these three assumptions, co-variation is the one analyzed using OLS. The often-repeated adage, correlation is not causation” is key. Causation is driven by theory, but co-variation is a critical part of empirical hypothesis testing.

When describing relationships, it is important to distinguish between those that are *deterministic* versus *stochastic*. Deterministic relationships are “fully determined” such that, knowing the values of the independent variable, you can perfectly explain (or predict) the value of the dependent variable. Philosophers of Old (like Kant) imagined the universe to be like a massive and complex clock which, once wound up and set ticking, would permit perfect prediction of the future if you had all the information on the starting conditions. There is no “error” in the prediction. Stochastic relationships, on the other hand, include an irreducible random component, such that the independent variables permit only a partial prediction of the dependent variable. But that stochastic (or random) component of the variation in the dependent variable has a probability distribution that can be analyzed statistically.

7.1.1 Deterministic Linear Model

The deterministic linear model serves as the basis for evaluating theoretical models. It is expressed as:

$$Y_i = \alpha + \beta X_i \quad (7.1)$$

A deterministic model is **systematic** and contains no error, therefore *YY is perfectly predicted by XX*. This is illustrated in Figure 7.1.1. α and β are the model parameters and are constant terms. β is the slope or the change in YY over the change in XX. α is the intercept, or the value of YY when XX is zero.

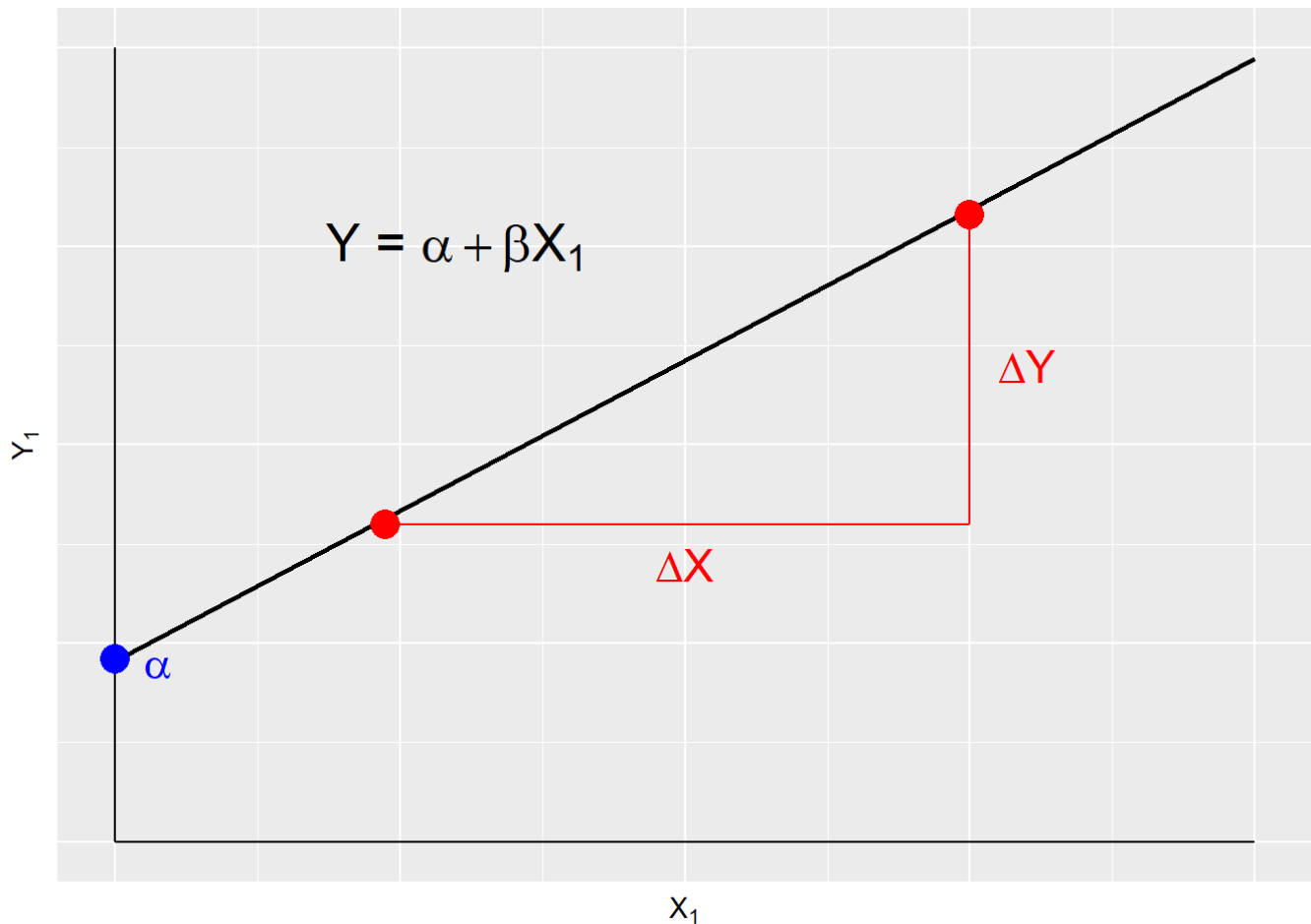


Figure 7.1.1: Deterministic Model

Given that in social science we rarely work with deterministic models, nearly all models contain a stochastic, or random, component.

7.1.2 Stochastic Linear Model

The stochastic, or statistical, the linear model contains a systematic component, $Y = \alpha + \beta X_1$, and a stochastic component called the **error term**. The error term is the difference between the expected value of Y_i and the observed value of Y_i ; $Y_i - \mu_{Y_i}$. This model is expressed as:

$$Y_i = \alpha + \beta X_i + \epsilon_i \quad (7.2)$$

where ϵ_i is the error term. In the deterministic model, each value of Y fits along the regression line, however in a stochastic model, the expected value of Y is conditioned by the values of X . This is illustrated in Figure 7.1.2.

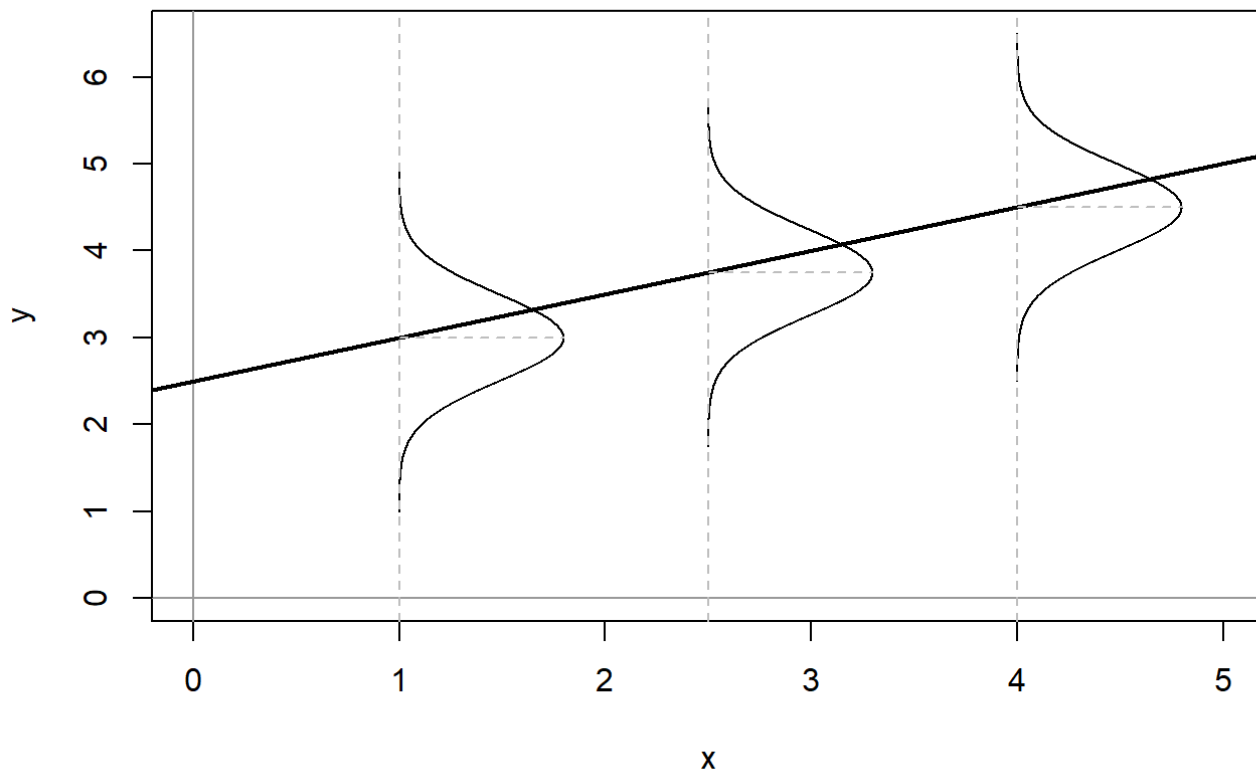


Figure 7.1.2: Stochastic Linear Model

Figure 7.1.2 shows the conditional population distributions of Y for several values of X , $p(Y|X)$. The conditional means of Y given X are denoted μ_i .

$$\mu_i = E(Y_i) = E(Y|X_i) = \alpha + \beta X_i \quad (7.3)$$

where $\alpha = E(Y) = \mu$ when $X=0$ - Each 1 unit increase in X increases $E(Y)$ by β

However, in the stochastic linear model variation in Y is caused by more than X , it is also caused by the error term ϵ . The error term is expressed as:

$$\begin{aligned} \epsilon_i &= Y_i - E(Y_i) = Y_i - (\alpha + \beta X_i) = Y_i - \alpha - \beta X_i \\ Y_i &= E(Y_i) + \epsilon_i = \alpha + \beta X_i + \epsilon_i \end{aligned}$$

We make several important assumptions about the error term that are discussed in the next section.

7.1.3 Assumptions about the Error Term

There are three key assumptions about the error term; a) errors have identical distributions, b) errors are independent, and c) errors are normally distributed.¹⁴

Error Assumptions

- Errors have identical distributions

$$E(\epsilon^2) = \sigma^2 \quad E(\epsilon) = 0$$

- Errors are independent of X and other ϵ_i

$$E(\epsilon_i) = E(\epsilon_j) = 0 \quad E(\epsilon_i \epsilon_j) = 0 \quad (i \neq j)$$

and

$$E(\epsilon_i) \neq E(\epsilon_j) \text{ for } i \neq j$$

- Errors are normally distributed

$$\epsilon_i \sim N(0, \sigma^2) \text{ for } i = 1, 2, \dots, n$$

Taken together these assumptions mean that the error term has a normal, independent, and identical distribution (normal i.i.d.). However, we don't know if, in any particular case, these assumptions are met. Therefore we must estimate a linear model.

This page titled [7.1: heoretical Models](#) is shared under a [CC BY 4.0](#) license and was authored, remixed, and/or curated by [Jenkins-Smith et al. \(University of Oklahoma Libraries\)](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.