

Review of Probability

This chapter reviews the core ideas of the theory of probability that are needed to understand regression analysis and econometrics. We assume that you have taken an introductory course in probability and statistics. If your knowledge of probability is stale, you should refresh it by reading this chapter. If you feel confident with the material, you still should skim the chapter and the terms and concepts at the end to make sure you are familiar with the ideas and notation.

Most aspects of the world around us have an element of randomness. The theory of probability provides mathematical tools for quantifying and describing this randomness. Section 2.1 reviews probability distributions for a single random variable, and Section 2.2 covers the mathematical expectation, mean, and variance of a single random variable. Most of the interesting problems in economics involve more than one variable, and Section 2.3 introduces the basic elements of probability theory for two random variables. Section 2.4 discusses three special probability distributions that play a central role in statistics and econometrics: the normal, chi-squared, and F distributions.

The final two sections of this chapter focus on a specific source of randomness of central importance in econometrics: the randomness that arises by randomly drawing a sample of data from a larger population. For example, suppose you survey ten recent college graduates selected at random, record (or “observe”) their earnings, and compute the average earnings using these ten data points (or “observations”). Because you chose the sample at random, you could have chosen ten different graduates by pure random chance; had you done so, you would have observed ten different earnings, and you would have computed a different sample average. Because the average earnings vary from one randomly chosen sample to the next, the sample average is itself a random variable. Therefore, the sample average has a probability distribution, which is referred to as its sampling distribution because this distribution describes the different possible values of the sample average that would have occurred had a different sample been drawn.

Section 2.5 discusses random sampling and the sampling distribution of the sample average. This sampling distribution is, in general, complicated. When the sample size is sufficiently large, however, the sampling distribution of the sample average is approximately normal, a result known as the central limit theorem, which is discussed in Section 2.6.

2.1 Random Variables and Probability Distributions

Probabilities, the Sample Space, and Random Variables

Probabilities and outcomes. The sex of the next new person you meet, your grade on an exam, and the number of times your wireless network connection fails while you are writing a term paper all have an element of chance or randomness. In each of these examples, there is something not yet known that is eventually revealed.

The mutually exclusive potential results of a random process are called the **outcomes**. For example, while writing your term paper, the wireless connection might never fail, it might fail once, it might fail twice, and so on. Only one of these outcomes will actually occur (the outcomes are mutually exclusive), and the outcomes need not be equally likely.

The **probability** of an outcome is the proportion of the time that the outcome occurs in the long run. If the probability of your wireless connection not failing while you are writing a term paper is 80%, then over the course of writing many term papers, you will complete 80% without a wireless connection failure.

The sample space and events. The set of all possible outcomes is called the **sample space**. An **event** is a subset of the sample space; that is, an event is a set of one or more outcomes. The event “my wireless connection will fail no more than once” is the set consisting of two outcomes: “no failures” and “one failure.”

Random variables. A random variable is a numerical summary of a random outcome. The number of times your wireless connection fails while you are writing a term paper is random and takes on a numerical value, so it is a random variable.

Some random variables are discrete and some are continuous. As their names suggest, a **discrete random variable** takes on only a discrete set of values, like 0, 1, 2, . . . , whereas a **continuous random variable** takes on a continuum of possible values.

Probability Distribution of a Discrete Random Variable

Probability distribution. The **probability distribution** of a discrete random variable is the list of all possible values of the variable and the probability that each value will occur. These probabilities sum to 1.

For example, let M be the number of times your wireless network connection fails while you are writing a term paper. The probability distribution of the random variable M is the list of probabilities of all possible outcomes: The probability that $M = 0$, denoted $\Pr(M = 0)$, is the probability of no wireless connection failures; $\Pr(M = 1)$ is the probability of a single connection failure; and so forth. An example of a probability distribution for M is given in the first row of Table 2.1. According to this distribution, the probability of no connection failures is 80%; the probability of one failure is 10%; and the probabilities of two, three, and four failures are,

TABLE 2.1 Probability of Your Wireless Network Connection Failing M Times

	Outcome (number of failures)				
	0	1	2	3	4
Probability distribution	0.80	0.10	0.06	0.03	0.01
Cumulative probability distribution	0.80	0.90	0.96	0.99	1.00

respectively, 6%, 3%, and 1%. These probabilities sum to 100%. This probability distribution is plotted in Figure 2.1.

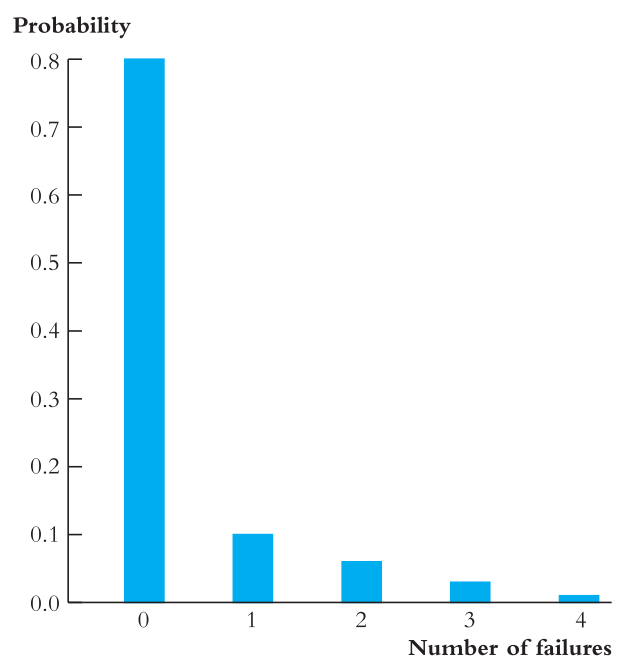
Probabilities of events. The probability of an event can be computed from the probability distribution. For example, the probability of the event of one or two failures is the sum of the probabilities of the constituent outcomes. That is, $\Pr(M = 1 \text{ or } M = 2) = \Pr(M = 1) + \Pr(M = 2) = 0.10 + 0.06 = 0.16$, or 16%.

Cumulative probability distribution. The **cumulative probability distribution** is the probability that the random variable is less than or equal to a particular value. The final row of Table 2.1 gives the cumulative probability distribution of the random variable M . For example, the probability of at most one connection failure, $\Pr(M \leq 1)$, is 90%, which is the sum of the probabilities of no failures (80%) and of one failure (10%).

A cumulative probability distribution is also referred to as a **cumulative distribution function**, a **c.d.f.**, or a **cumulative distribution**.

FIGURE 2.1 Probability Distribution of the Number of Wireless Network Connection Failures

The height of each bar is the probability that the wireless connection fails the indicated number of times. The height of the first bar is 0.8, so the probability of 0 connection failures is 80%. The height of the second bar is 0.1, so the probability of 1 failure is 10%, and so forth for the other bars.



The Bernoulli distribution. An important special case of a discrete random variable is when the random variable is binary; that is, the outcome is 0 or 1. A binary random variable is called a **Bernoulli random variable** (in honor of the 17th-century Swiss mathematician and scientist Jacob Bernoulli), and its probability distribution is called the **Bernoulli distribution**.

For example, let G be the sex of the next new person you meet, where $G = 0$ indicates that the person is male and $G = 1$ indicates that the person is female. The outcomes of G and their probabilities thus are

$$G = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p, \end{cases} \quad (2.1)$$

where p is the probability of the next new person you meet being a woman. The probability distribution in Equation (2.1) is the Bernoulli distribution.

Probability Distribution of a Continuous Random Variable

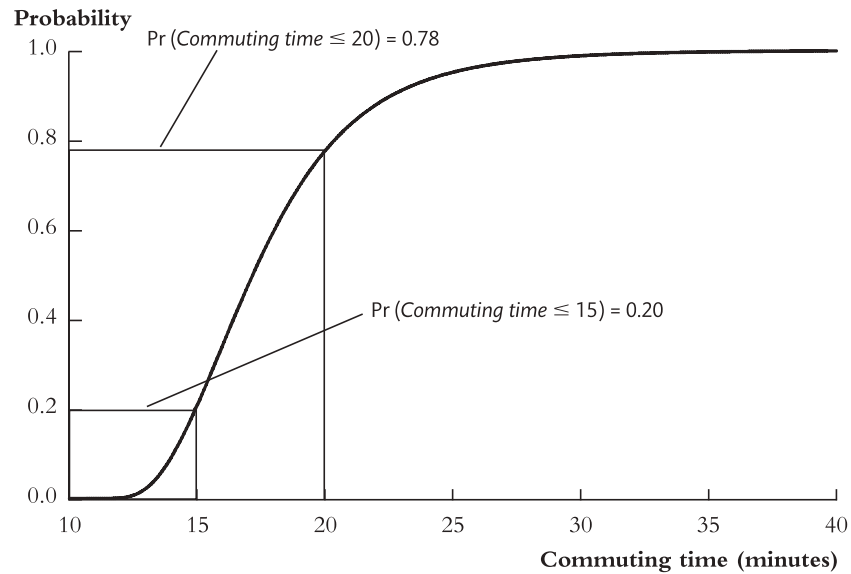
Cumulative probability distribution. The cumulative probability distribution for a continuous variable is defined just as it is for a discrete random variable. That is, the cumulative probability distribution of a continuous random variable is the probability that the random variable is less than or equal to a particular value.

For example, consider a student who drives from home to school. This student's commuting time can take on a continuum of values, and because it depends on random factors such as the weather and traffic conditions, it is natural to treat it as a continuous random variable. Figure 2.2a plots a hypothetical cumulative distribution of commuting times. For example, the probability that the commute takes less than 15 minutes is 20%, and the probability that it takes less than 20 minutes is 78%.

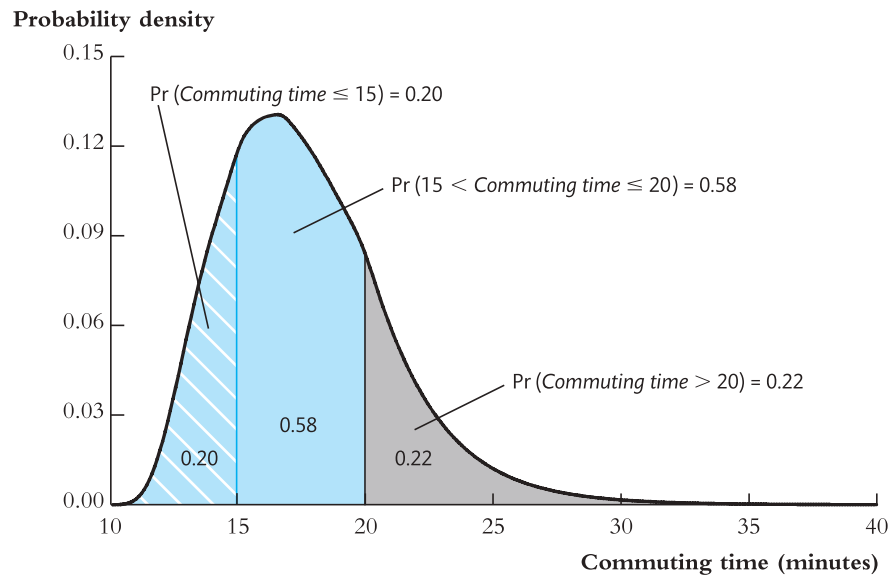
Probability density function. Because a continuous random variable can take on a continuum of possible values, the probability distribution used for discrete variables, which lists the probability of each possible value of the random variable, is not suitable for continuous variables. Instead, the probability is summarized by the **probability density function**. The area under the probability density function between any two points is the probability that the random variable falls between those two points. A probability density function is also called a **p.d.f.**, a **density function**, or simply a **density**.

Figure 2.2b plots the probability density function of commuting times corresponding to the cumulative distribution in Figure 2.2a. The probability that the commute takes between 15 and 20 minutes is given by the area under the p.d.f. between 15 minutes and 20 minutes, which is 0.58, or 58%. Equivalently, this probability can be seen on the cumulative distribution in Figure 2.2a as the difference between the probability that the commute is less than 20 minutes (78%) and the probability that it is less than 15 minutes (20%). Thus the probability density function and the cumulative probability distribution show the same information in different formats.

FIGURE 2.2 Cumulative Probability Distribution and Probability Density Functions of Commuting Time



(a) Cumulative probability distribution function of commuting times



(b) Probability density function of commuting times

Figure 2.2a shows the cumulative probability distribution function (c.d.f.) of commuting times. The probability that a commuting time is less than 15 minutes is 0.20 (or 20%), and the probability that it is less than 20 minutes is 0.78 (78%). Figure 2.2b shows the probability density function (or p.d.f.) of commuting times. Probabilities are given by areas under the p.d.f. The probability that a commuting time is between 15 and 20 minutes is 0.58 (58%) and is given by the area under the curve between 15 and 20 minutes.

2.2 Expected Values, Mean, and Variance

The Expected Value of a Random Variable

Expected value. The **expected value** of a random variable Y , denoted $E(Y)$, is the long-run average value of the random variable over many repeated trials or occurrences. The expected value of a discrete random variable is computed as a weighted average of the possible outcomes of that random variable, where the weights are the probabilities of that outcome. The expected value of Y is also called the **expectation** of Y or the **mean** of Y and is denoted μ_Y .

For example, suppose you loan a friend \$100 at 10% interest. If the loan is repaid, you get \$110 (the principal of \$100 plus interest of \$10), but there is a risk of 1% that your friend will default and you will get nothing at all. Thus the amount you are repaid is a random variable that equals \$110 with probability 0.99 and equals \$0 with probability 0.01. Over many such loans, 99% of the time you would be paid back \$110, but 1% of the time you would get nothing, so on average you would be repaid $\$110 \times 0.99 + \$0 \times 0.01 = \$108.90$. Thus the expected value of your repayment is \$108.90.

As a second example, consider the number of wireless network connection failures M with the probability distribution given in Table 2.1. The expected value of M —that is, the mean of M —is the average number of failures over many term papers, weighted by the frequency with which a given number of failures occurs. Accordingly,

$$E(M) = 0 \times 0.80 + 1 \times 0.10 + 2 \times 0.06 + 3 \times 0.03 + 4 \times 0.01 = 0.35. \quad (2.2)$$

That is, the expected number of connection failures while writing a term paper is 0.35. Of course, the actual number of failures must always be an integer; it makes no sense to say that the wireless connection failed 0.35 times while writing a particular term paper! Rather, the calculation in Equation (2.2) means that the average number of failures over many such term papers is 0.35.

The formula for the expected value of a discrete random variable Y that can take on k different values is given in Key Concept 2.1. (Key Concept 2.1 uses summation notation, which is reviewed in Exercise 2.25.)

KEY CONCEPT

Expected Value and the Mean

2.1

Suppose that the random variable Y takes on k possible values, y_1, \dots, y_k , where y_1 denotes the first value, y_2 denotes the second value, and so forth, and that the probability that Y takes on y_1 is p_1 , the probability that Y takes on y_2 is p_2 , and so forth. The expected value of Y , denoted $E(Y)$, is

$$E(Y) = y_1p_1 + y_2p_2 + \cdots + y_kp_k = \sum_{i=1}^k y_i p_i, \quad (2.3)$$

where the notation $\sum_{i=1}^k y_i p_i$ means “the sum of $y_i p_i$ for i running from 1 to k .” The expected value of Y is also called the mean of Y or the expectation of Y and is denoted μ_Y .

Expected value of a Bernoulli random variable. An important special case of the general formula in Key Concept 2.1 is the mean of a Bernoulli random variable. Let G be the Bernoulli random variable with the probability distribution in Equation (2.1). The expected value of G is

$$E(G) = 0 \times (1 - p) + 1 \times p = p. \quad (2.4)$$

Thus the expected value of a Bernoulli random variable is p , the probability that it takes on the value 1.

Expected value of a continuous random variable. The expected value of a continuous random variable is also the probability-weighted average of the possible outcomes of the random variable. Because a continuous random variable can take on a continuum of possible values, the formal mathematical definition of its expectation involves calculus and its definition is given in Appendix 18.1.

The Standard Deviation and Variance

The variance and standard deviation measure the dispersion or the “spread” of a probability distribution. The **variance** of a random variable Y , denoted $\text{var}(Y)$, is the expected value of the square of the deviation of Y from its mean: $\text{var}(Y) = E[(Y - \mu_Y)^2]$.

Because the variance involves the square of Y , the units of the variance are the units of the square of Y , which makes the variance awkward to interpret. It is therefore common to measure the spread by the **standard deviation**, which is the square root of the variance and is denoted σ_Y . The standard deviation has the same units as Y . These definitions are summarized in Key Concept 2.2.

For example, the variance of the number of connection failures M is the probability-weighted average of the squared difference between M and its mean, 0.35:

$$\begin{aligned} \text{var}(M) &= (0 - 0.35)^2 \times 0.80 + (1 - 0.35)^2 \times 0.10 + (2 - 0.35)^2 \times 0.06 \\ &\quad + (3 - 0.35)^2 \times 0.03 + (4 - 0.35)^2 \times 0.01 = 0.6475. \end{aligned} \quad (2.5)$$

The standard deviation of M is the square root of the variance, so $\sigma_M = \sqrt{0.6475} \cong 0.80$.

Variance and Standard Deviation

KEY CONCEPT

2.2

The variance of the discrete random variable Y , denoted σ_Y^2 , is

$$\sigma_Y^2 = \text{var}(Y) = E[(Y - \mu_Y)^2] = \sum_{i=1}^k (y_i - \mu_Y)^2 p_i. \quad (2.6)$$

The standard deviation of Y is σ_Y , the square root of the variance. The units of the standard deviation are the same as the units of Y .

Variance of a Bernoulli random variable. The mean of the Bernoulli random variable G with the probability distribution in Equation (2.1) is $\mu_G = p$ [Equation (2.4)], so its variance is

$$\text{var}(G) = \sigma_G^2 = (0 - p)^2 \times (1 - p) + (1 - p)^2 \times p = p(1 - p). \quad (2.7)$$

Thus the standard deviation of a Bernoulli random variable is $\sigma_G = \sqrt{p(1 - p)}$.

Mean and Variance of a Linear Function of a Random Variable

This section discusses random variables (say, X and Y) that are related by a linear function. For example, consider an income tax scheme under which a worker is taxed at a rate of 20% on his or her earnings and then given a (tax-free) grant of \$2000. Under this tax scheme, after-tax earnings Y are related to pre-tax earnings X by the equation

$$Y = 2000 + 0.8X. \quad (2.8)$$

That is, after-tax earnings Y is 80% of pre-tax earnings X , plus \$2000.

Suppose an individual's pre-tax earnings next year are a random variable with mean μ_X and variance σ_X^2 . Because pre-tax earnings are random, so are after-tax earnings. What are the mean and standard deviations of her after-tax earnings under this tax? After taxes, her earnings are 80% of the original pre-tax earnings, plus \$2000. Thus the expected value of her after-tax earnings is

$$E(Y) = \mu_Y = 2000 + 0.8\mu_X. \quad (2.9)$$

The variance of after-tax earnings is the expected value of $(Y - \mu_Y)^2$. Because $Y = 2000 + 0.8X$, $Y - \mu_Y = 2000 + 0.8X - (2000 + 0.8\mu_X) = 0.8(X - \mu_X)$. Thus $E[(Y - \mu_Y)^2] = E\{[0.8(X - \mu_X)]^2\} = 0.64E[(X - \mu_X)^2]$. It follows that $\text{var}(Y) = 0.64\text{var}(X)$, so, taking the square root of the variance, the standard deviation of Y is

$$\sigma_Y = 0.8\sigma_X. \quad (2.10)$$

That is, the standard deviation of the distribution of her after-tax earnings is 80% of the standard deviation of the distribution of her pre-tax earnings.

This analysis can be generalized so that Y depends on X with an intercept a (instead of \$2000) and a slope b (instead of 0.8) so that

$$Y = a + bX. \quad (2.11)$$

Then the mean and variance of Y are

$$\mu_Y = a + b\mu_X \quad \text{and} \quad (2.12)$$

$$\sigma_Y^2 = b^2\sigma_X^2, \quad (2.13)$$

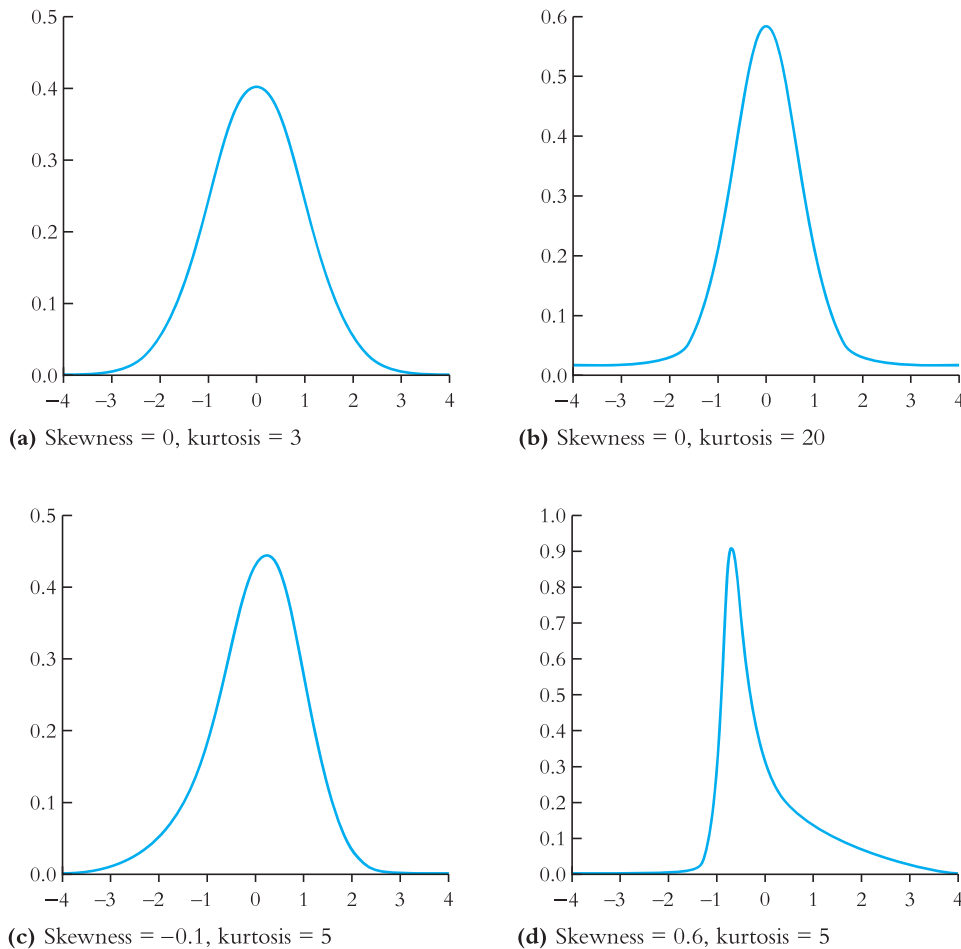
and the standard deviation of Y is $\sigma_Y = b\sigma_X$. The expressions in Equations (2.9) and (2.10) are applications of the more general formulas in Equations (2.12) and (2.13) with $a = 2000$ and $b = 0.8$.

Other Measures of the Shape of a Distribution

The mean and standard deviation measure two important features of a distribution: its center (the mean) and its spread (the standard deviation). This section discusses measures of two other features of a distribution: the skewness, which measures the lack of symmetry of a distribution, and the kurtosis, which measures how thick, or “heavy,” are its tails. The mean, variance, skewness, and kurtosis are all based on what are called the **moments of a distribution**.

Skewness. Figure 2.3 plots four distributions, two that are symmetric (Figures 2.3a and 2.3b) and two that are not (Figures 2.3c and 2.3d). Visually, the distribution in Figure 2.3d appears to deviate more from symmetry than does the distribution in

FIGURE 2.3 Four Distributions with Different Skewness and Kurtosis



All of these distributions have a mean of 0 and a variance of 1. The distributions with skewness of 0 (a and b) are symmetric; the distributions with nonzero skewness (c and d) are not symmetric. The distributions with kurtosis exceeding 3 (b, c, and d) have heavy tails.

Figure 2.3c. The skewness of a distribution provides a mathematical way to describe how much a distribution deviates from symmetry.

The **skewness** of the distribution of a random variable Y is

$$\text{Skewness} = \frac{E[(Y - \mu_Y)^3]}{\sigma_Y^3}, \quad (2.14)$$

where σ_Y is the standard deviation of Y . For a symmetric distribution, a value of Y a given amount above its mean is just as likely as a value of Y the same amount below its mean. If so, then positive values of $(Y - \mu_Y)^3$ will be offset on average (in expectation) by equally likely negative values. Thus, for a symmetric distribution, $E(Y - \mu_Y)^3 = 0$: The skewness of a symmetric distribution is 0. If a distribution is not symmetric, then a positive value of $(Y - \mu_Y)^3$ generally is not offset on average by an equally likely negative value, so the skewness is nonzero for a distribution that is not symmetric. Dividing by σ_Y^3 in the denominator of Equation (2.14) cancels the units of Y^3 in the numerator, so the skewness is unit free; in other words, changing the units of Y does not change its skewness.

Below each of the four distributions in Figure 2.3 is its skewness. If a distribution has a long right tail, positive values of $(Y - \mu_Y)^3$ are not fully offset by negative values, and the skewness is positive. If a distribution has a long left tail, its skewness is negative.

Kurtosis. The **kurtosis** of a distribution is a measure of how much mass is in its tails and therefore is a measure of how much of the variance of Y arises from extreme values. An extreme value of Y is called an **outlier**. The greater the kurtosis of a distribution, the more likely are outliers.

The kurtosis of the distribution of Y is

$$\text{Kurtosis} = \frac{E[(Y - \mu_Y)^4]}{\sigma_Y^4}. \quad (2.15)$$

If a distribution has a large amount of mass in its tails, then some extreme departures of Y from its mean are likely, and these departures will lead to large values, on average (in expectation), of $(Y - \mu_Y)^4$. Thus, for a distribution with a large amount of mass in its tails, the kurtosis will be large. Because $(Y - \mu_Y)^4$ cannot be negative, the kurtosis cannot be negative.

The kurtosis of a normally distributed random variable is 3, so a random variable with kurtosis exceeding 3 has more mass in its tails than a normal random variable. A distribution with kurtosis exceeding 3 is called **leptokurtic** or, more simply, heavy-tailed. Like skewness, the kurtosis is unit free, so changing the units of Y does not change its kurtosis.

Below each of the four distributions in Figure 2.3 is its kurtosis. The distributions in Figures 2.3b–d are heavy-tailed.

Moments. The mean of Y , $E(Y)$, is also called the first moment of Y , and the expected value of the square of Y , $E(Y^2)$, is called the second moment of Y . In general, the

expected value of Y^r is called the r^{th} **moment** of the random variable Y . That is, the r^{th} moment of Y is $E(Y^r)$. The skewness is a function of the first, second, and third moments of Y , and the kurtosis is a function of the first through fourth moments of Y .

Standardized Random Variables

A random variable can be transformed into a random variable with mean 0 and variance 1 by subtracting its mean and then dividing by its standard deviation, a process called standardization. Specifically, let Y have mean μ_Y and variance σ_Y^2 . Then the **standardized random variable** computed from Y is $(Y - \mu_Y)/\sigma_Y$. The mean of the standardized random variable is $E(Y - \mu_Y)/\sigma_Y = (EY - \mu_Y)/\sigma_Y = 0$, and its variance is $\text{var}[(Y - \mu_Y)/\sigma_Y] = \text{var}(Y)/\sigma_Y^2 = 1$. Standardized random variables do not have any units, such as dollars or meters, because the units of Y are canceled by dividing through by σ_Y , which also has the units of Y .

2.3 Two Random Variables

Most of the interesting questions in economics involve two or more variables. Are college graduates more likely to have a job than nongraduates? How does the distribution of income for women compare to that for men? These questions concern the distribution of two random variables, considered together (education and employment status in the first example, income and sex in the second). Answering such questions requires an understanding of the concepts of joint, marginal, and conditional probability distributions.

Joint and Marginal Distributions

Joint distribution. The **joint probability distribution** of two discrete random variables, say X and Y , is the probability that the random variables simultaneously take on certain values, say x and y . The probabilities of all possible (x, y) combinations sum to 1. The joint probability distribution can be written as the function $\text{Pr}(X = x, Y = y)$.

For example, weather conditions—whether or not it is raining—affect the commuting time of the student commuter in Section 2.1. Let Y be a binary random variable that equals 1 if the commute is short (less than 20 minutes) and that equals 0 otherwise, and let X be a binary random variable that equals 0 if it is raining and 1 if not. Between these two random variables, there are four possible outcomes: it rains and the commute is long ($X = 0, Y = 0$); rain and short commute ($X = 0, Y = 1$); no rain and long commute ($X = 1, Y = 0$); and no rain and short commute ($X = 1, Y = 1$). The joint probability distribution is the frequency with which each of these four outcomes occurs over many repeated commutes.

An example of a joint distribution of these two variables is given in Table 2.2. According to this distribution, over many commutes, 15% of the days have rain and a long commute ($X = 0, Y = 0$); that is, the probability of a long rainy commute is

TABLE 2.2 Joint Distribution of Weather Conditions and Commuting Times

	Rain ($X = 0$)	No Rain ($X = 1$)	Total
Long commute ($Y = 0$)	0.15	0.07	0.22
Short commute ($Y = 1$)	0.15	0.63	0.78
Total	0.30	0.70	1.00

15%, or $\Pr(X = 0, Y = 0) = 0.15$. Also, $\Pr(X = 0, Y = 1) = 0.15$, $\Pr(X = 1, Y = 0) = 0.07$, and $\Pr(X = 1, Y = 1) = 0.63$. These four possible outcomes are mutually exclusive and constitute the sample space, so the four probabilities sum to 1.

Marginal probability distribution. The **marginal probability distribution** of a random variable Y is just another name for its probability distribution. This term is used to distinguish the distribution of Y alone (the marginal distribution) from the joint distribution of Y and another random variable.

The marginal distribution of Y can be computed from the joint distribution of X and Y by adding up the probabilities of all possible outcomes for which Y takes on a specified value. If X can take on l different values x_1, \dots, x_l , then the marginal probability that Y takes on the value y is

$$\Pr(Y = y) = \sum_{i=1}^l \Pr(X = x_i, Y = y). \quad (2.16)$$

For example, in Table 2.2, the probability of a long rainy commute is 15%, and the probability of a long commute with no rain is 7%, so the probability of a long commute (rainy or not) is 22%. The marginal distribution of commuting times is given in the final column of Table 2.2. Similarly, the marginal probability that it will rain is 30%, as shown in the final row of Table 2.2.

Conditional Distributions

Conditional distribution. The distribution of a random variable Y conditional on another random variable X taking on a specific value is called the **conditional distribution** of Y given X . The conditional probability that Y takes on the value y when X takes on the value x is written $\Pr(Y = y | X = x)$.

For example, what is the probability of a long commute ($Y = 0$) if you know it is raining ($X = 0$)? From Table 2.2, the joint probability of a rainy short commute is 15%, and the joint probability of a rainy long commute is 15%, so if it is raining, a long commute and a short commute are equally likely. Thus the probability of a long commute ($Y = 0$) conditional on it being rainy ($X = 0$) is 50%, or $\Pr(Y = 0 | X = 0) = 0.50$. Equivalently, the marginal probability of rain is 30%; that is, over many commutes, it rains 30% of the time. Of this 30% of commutes, 50% of the time the commute is long ($0.15/0.30$).

TABLE 2.3 Joint and Conditional Distributions of Number of Wireless Connection Failures (M) and Network Age (A)

A. Joint Distribution						
	$M = 0$	$M = 1$	$M = 2$	$M = 3$	$M = 4$	Total
Old network ($A = 0$)	0.35	0.065	0.05	0.025	0.01	0.50
New network ($A = 1$)	0.45	0.035	0.01	0.005	0.00	0.50
Total	0.80	0.10	0.06	0.03	0.01	1.00
B. Conditional Distributions of M given A						
	$M = 0$	$M = 1$	$M = 2$	$M = 3$	$M = 4$	Total
$\Pr(M A = 0)$	0.70	0.13	0.10	0.05	0.02	1.00
$\Pr(M A = 1)$	0.90	0.07	0.02	0.01	0.00	1.00

In general, the conditional distribution of Y given $X = x$ is

$$\Pr(Y = y|X = x) = \frac{\Pr(X = x, Y = y)}{\Pr(X = x)}. \quad (2.17)$$

For example, the conditional probability of a long commute given that it is rainy is $\Pr(Y = 0|X = 0) = \Pr(X = 0, Y = 0)/\Pr(X = 0) = 0.15/0.30 = 0.50$.

As a second example, consider a modification of the network connection failure example. Suppose that half the time you write your term paper in the school library, which has a new wireless network; otherwise, you write it in your room, which has an old wireless network. If we treat the location where you write the term paper as random, then the network age A ($= 1$ if the network is new, $= 0$ if it is old) is a random variable. Suppose the joint distribution of the random variables M and A is given in Part A of Table 2.3. Then the conditional distributions of connection failures given the age of the network are shown in Part B of the table. For example, the joint probability of $M = 0$ and $A = 0$ is 0.35; because half the time you use the old network, the conditional probability of no failures given that you use the old network is $\Pr(M = 0|A = 0) = \Pr(M = 0, A = 0)/\Pr(A = 0) = 0.35/0.50 = 0.70$, or 70%. In contrast, the conditional probability of no failures given that you use the new network is 90%. According to the conditional distributions in Part B of Table 2.3, the new network is less likely to fail than the old one; for example, the probability of three failures is 5% using the old network but 1% using the new network.

Conditional expectation. The **conditional expectation** of Y given X , also called the **conditional mean** of Y given X , is the mean of the conditional distribution of Y given X . That is, the conditional expectation is the expected value of Y , computed using the conditional distribution of Y given X . If Y takes on k values y_1, \dots, y_k , then the conditional mean of Y given $X = x$ is

$$E(Y|X = x) = \sum_{i=1}^k y_i \Pr(Y = y_i|X = x). \quad (2.18)$$

For example, based on the conditional distributions in Table 2.3, the expected number of connection failures, given that the network is old, is $E(M|A = 0) = 0 \times 0.70 + 1 \times 0.13 + 2 \times 0.10 + 3 \times 0.05 + 4 \times 0.02 = 0.56$. The expected number of failures, given that the network is new, is $E(M|A = 1) = 0.14$, less than for the old network.

The conditional expectation of Y given $X = x$ is just the mean value of Y when $X = x$. In the example of Table 2.3, the mean number of failures is 0.56 for the old network, so the conditional expectation of Y given that the network is old is 0.56. Similarly, for the new network, the mean number of failures is 0.14; that is, the conditional expectation of Y given that the network is new is 0.14.

The law of iterated expectations. The mean of Y is the weighted average of the conditional expectation of Y given X , weighted by the probability distribution of X . For example, the mean height of adults is the weighted average of the mean height of men and the mean height of women, weighted by the proportions of men and women. Stated mathematically, if X takes on the l values x_1, \dots, x_l , then

$$E(Y) = \sum_{i=1}^l E(Y|X = x_i) \Pr(X = x_i). \quad (2.19)$$

Equation (2.19) follows from Equations (2.18) and (2.17) (see Exercise 2.19).

Stated differently, the expectation of Y is the expectation of the conditional expectation of Y given X ,

$$E(Y) = E[E(Y|X)], \quad (2.20)$$

where the inner expectation on the right-hand side of Equation (2.20) is computed using the conditional distribution of Y given X and the outer expectation is computed using the marginal distribution of X . Equation (2.20) is known as the **law of iterated expectations**.

For example, the mean number of connection failures M is the weighted average of the conditional expectation of M given that it is old and the conditional expectation of M given that it is new, so $E(M) = E(M|A = 0) \times \Pr(A = 0) + E(M|A = 1) \times \Pr(A = 1) = 0.56 \times 0.50 + 0.14 \times 0.50 = 0.35$. This is the mean of the marginal distribution of M , as calculated in Equation (2.2).

The law of iterated expectations implies that if the conditional mean of Y given X is 0, then the mean of Y is 0. This is an immediate consequence of Equation (2.20): if $E(Y|X) = 0$, then $E(Y) = E[E(Y|X)] = E[0] = 0$. Said differently, if the mean of Y given X is 0, then it must be that the probability-weighted average of these conditional means is 0; that is, the mean of Y must be 0.

The law of iterated expectations also applies to expectations that are conditional on multiple random variables. For example, let X , Y , and Z be random variables that are jointly distributed. Then the law of iterated expectations says that $E(Y) = E[E(Y|X, Z)]$, where $E(Y|X, Z)$ is the conditional expectation of Y

given both X and Z . For example, in the network connection illustration of Table 2.3, let P denote the number of people using the network; then $E(M|A, P)$ is the expected number of failures for a network with age A that has P users. The expected number of failures overall, $E(M)$, is the weighted average of the expected number of failures for a network with age A and number of users P , weighted by the proportion of occurrences of both A and P .

Exercise 2.20 provides some additional properties of conditional expectations with multiple variables.

Conditional variance. The variance of Y conditional on X is the variance of the conditional distribution of Y given X . Stated mathematically, the **conditional variance** of Y given X is

$$\text{var}(Y|X = x) = \sum_{i=1}^k [y_i - E(Y|X = x)]^2 \Pr(Y = y_i|X = x). \quad (2.21)$$

For example, the conditional variance of the number of failures given that the network is old is $\text{var}(M|A = 0) = (0 - 0.56)^2 \times 0.70 + (1 - 0.56)^2 \times 0.13 + (2 - 0.56)^2 \times 0.10 + (3 - 0.56)^2 \times 0.05 + (4 - 0.56)^2 \times 0.02 \cong 0.99$. The standard deviation of the conditional distribution of M given that $A = 0$ is thus $\sqrt{0.99} = 0.99$. The conditional variance of M given that $A = 1$ is the variance of the distribution in the second row of Part B of Table 2.3, which is 0.22, so the standard deviation of M for the new network is $\sqrt{0.22} = 0.47$. For the conditional distributions in Table 2.3, the expected number of failures for the new network (0.14) is less than that for the old network (0.56), and the spread of the distribution of the number of failures, as measured by the conditional standard deviation, is smaller for the new network (0.47) than for the old (0.99).

Bayes' rule. Bayes' rule says that the conditional probability of Y given X is the conditional probability of X given Y times the relative marginal probabilities of Y and X :

$$\Pr(Y = y|X = x) = \frac{\Pr(X = x|Y = y)\Pr(Y = y)}{\Pr(X = x)} \quad (\text{Bayes' rule}). \quad (2.22)$$

Equation (2.22) obtains from the definition of the conditional distribution in Equation (2.17), which implies that $\Pr(X = x, Y = y) = \Pr(Y = y|X = x) \Pr(X = x)$ and that $\Pr(X = x, Y = y) = \Pr(X = x|Y = y)\Pr(Y = y)$; equating the second parts of these two equalities and rearranging gives Bayes' rule.

Bayes' rule can be used to deduce conditional probabilities from the reverse conditional probability, with the help of marginal probabilities. For example, suppose you told your friend that you were dropped by the network three times last night while working on your term paper and your friend knows that half the time you work in the library and half the time you work in your room. Then your friend could deduce from Table 2.3 that the probability you worked in your room last night given three network failures is 83% (Exercise 2.28).

The conditional mean is the minimum mean squared error prediction. The conditional mean plays a central role in prediction; in fact it is, in a precise sense, the optimal prediction of Y given $X = x$.

A common formulation of the statistical prediction problem is to posit that the cost of making a prediction error increases with the square of that error. The motivation for this squared-error prediction loss is that small errors in prediction might not matter much, but large errors can be very costly in real-world applications. Stated mathematically, the prediction problem thus is: what is the function $g(X)$ that minimizes the mean squared prediction error, $E\{[Y - g(X)]^2\}$? The answer is the conditional mean $E(Y|X)$: Of all possible ways to use the information X , the conditional mean minimizes the mean squared prediction error. This result is proven in Appendix 2.2.

Independence

Two random variables X and Y are **independently distributed**, or **independent**, if knowing the value of one of the variables provides no information about the other. Specifically, X and Y are independent if the conditional distribution of Y given X equals the marginal distribution of Y . That is, X and Y are independently distributed if, for all values of x and y ,

$$\Pr(Y = y|X = x) = \Pr(Y = y) \text{ (independence of } X \text{ and } Y\text{)}. \quad (2.23)$$

Substituting Equation (2.23) into Equation (2.17) gives an alternative expression for independent random variables in terms of their joint distribution. If X and Y are independent, then

$$\Pr(X = x, Y = y) = \Pr(X = x)\Pr(Y = y). \quad (2.24)$$

That is, the joint distribution of two independent random variables is the product of their marginal distributions.

Covariance and Correlation

Covariance. One measure of the extent to which two random variables move together is their covariance. The **covariance** between X and Y is the expected value $E[(X - \mu_X)(Y - \mu_Y)]$, where μ_X is the mean of X and μ_Y is the mean of Y . The covariance is denoted $\text{cov}(X, Y)$ or σ_{XY} . If X can take on l values and Y can take on k values, then the covariance is given by the formula

$$\begin{aligned} \text{cov}(X, Y) &= \sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)] \\ &= \sum_{i=1}^k \sum_{j=1}^l (x_j - \mu_X)(y_i - \mu_Y)\Pr(X = x_j, Y = y_i). \end{aligned} \quad (2.25)$$

To interpret this formula, suppose that when X is greater than its mean (so that $X - \mu_X$ is positive), then Y tends to be greater than its mean (so that $Y - \mu_Y$ is

positive) and that when X is less than its mean (so that $X - \mu_X < 0$), then Y tends to be less than its mean (so that $Y - \mu_Y < 0$). In both cases, the product $(X - \mu_X) \times (Y - \mu_Y)$ tends to be positive, so the covariance is positive. In contrast, if X and Y tend to move in opposite directions (so that X is large when Y is small, and vice versa), then the covariance is negative. Finally, if X and Y are independent, then the covariance is 0 (see Exercise 2.19).

Correlation. Because the covariance is the product of X and Y , deviated from their means, its units are, awkwardly, the units of X multiplied by the units of Y . This “units” problem can make numerical values of the covariance difficult to interpret.

The correlation is an alternative measure of dependence between X and Y that solves the “units” problem of the covariance. Specifically, the **correlation** between X and Y is the covariance between X and Y divided by their standard deviations:

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{var}(X) \text{var}(Y)}} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}. \quad (2.26)$$

Because the units of the numerator in Equation (2.26) are the same as those of the denominator, the units cancel, and the correlation is unit free. The random variables X and Y are said to be **uncorrelated** if $\text{corr}(X, Y) = 0$.

The correlation always is between -1 and 1 ; that is, as proven in Appendix 2.1,

$$-1 \leq \text{corr}(X, Y) \leq 1 \quad (\text{correlation inequality}). \quad (2.27)$$

Correlation and conditional mean. If the conditional mean of Y does not depend on X , then Y and X are uncorrelated. That is,

$$\text{if } E(Y|X) = \mu_Y, \text{ then } \text{cov}(Y, X) = 0 \text{ and } \text{corr}(Y, X) = 0. \quad (2.28)$$

We now show this result. First, suppose Y and X have mean 0, so that $\text{cov}(Y, X) = E[(Y - \mu_Y)(X - \mu_X)] = E(YX)$. By the law of iterated expectations [Equation (2.20)], $E(YX) = E[E(YX|X)] = E[E(Y|X)X] = 0$ because $E(Y|X) = 0$, so $\text{cov}(Y, X) = 0$. Equation (2.28) follows by substituting $\text{cov}(Y, X) = 0$ into the definition of correlation in Equation (2.26). If Y and X do not have mean 0, subtract off their means, and then the preceding proof applies.

It is *not* necessarily true, however, that if X and Y are uncorrelated, then the conditional mean of Y given X does not depend on X . Said differently, it is possible for the conditional mean of Y to be a function of X but for Y and X nonetheless to be uncorrelated. An example is given in Exercise 2.23.

The Mean and Variance of Sums of Random Variables

The mean of the sum of two random variables, X and Y , is the sum of their means:

$$E(X + Y) = E(X) + E(Y) = \mu_X + \mu_Y. \quad (2.29)$$

The Distribution of Adulthood Earnings in the United Kingdom by Childhood Socioeconomic Circumstances

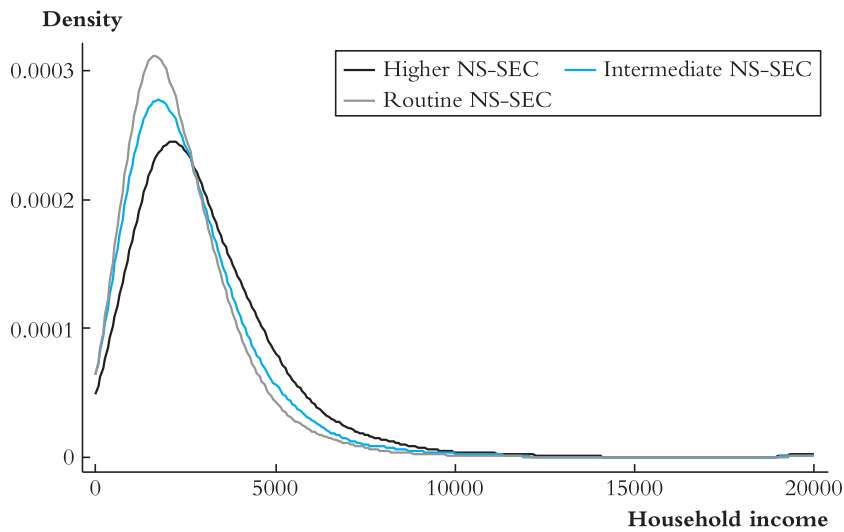
Politicians sometimes talk about how inequality in income arises as a result of differences in individual ability and effort. Are these politicians right? Or, in contrast, do childhood circumstances affect an individual's income during adulthood? For example, do children who grow up with fewer advantages go on to be part of households with lower average income?

One way to answer these questions is by considering how an individual's household income as

an adult varies according to their father's occupational type. While no two occupations are identical, researchers often group similar jobs into a given number of meaningful classes. One method of doing this, as seen in the United Kingdom's National Statistics Socio-economic Classification (NS-SEC),¹ is grouping jobs into a hierarchy of three classes: higher, intermediate, and routine.

Figure 2.4 illustrates these three conditional distributions of household income for individuals in

FIGURE 2.4 Conditional Distributions of Household Income of U.K. individuals in 2009–2010, by Occupational Type of Father



The three distributions of household incomes are for individuals in the United Kingdom, based on the National Statistics Socio-economic Classification (NS-SEC) of their father—higher, intermediate, and routine jobs.

¹For further details refer to “The National Statistics Socio-economic classification (NS-SEC),” The Office for National Statistics, <https://www.ons.gov.uk/>, 2010.

TABLE 2.4 Summaries of the Conditional Distribution of Monthly Household Income for Individuals in the United Kingdom Given NS-SEC of Father's Occupation

NS-SEC of Father's Job	Mean	Standard Deviation	Percentile			
			25%	50% (median)	75%	90%
(a) Higher	£3,149.27	£2,434.33	£1,663.33	£2,626.92	£3,973.74	£5,629.00
(b) Intermediate	2,692.01	2,187.53	1,362.44	2,237.56	3,382.00	4,881.99
(c) Routine	2,440.94	1,878.58	1,291.00	2,049.74	3,067.76	4,339.84

the United Kingdom in 2009 and 2010 according to the NS-SEC of their father's occupation in that individual's childhood.² The lower the classification of paternal occupation, the more concentrated in the lower end of the distribution is household income in adulthood.

The statistics for monthly household income for these individuals by NS-SEC classification are summarized in Table 2.4. For example, the mean income of individuals whose father's occupation is classified as routine, that is, $E(\text{Income}|\text{Father's social class} = \text{routine})$, was £2,440.94. This is over £700 less than that for individuals whose father's occupation is classified as higher, that is, $E(\text{Income}|\text{Father's social class} = \text{higher})$, which is £3,149.27. Furthermore, these differences are much greater at higher ends of the

distribution, with the difference in income between these groups being over £900 at the 75th percentile and almost £1,300 at the 90th percentile. The standard deviation of household income also increases with occupation classification, meaning that the spread of household income is also greater according to this measure.

This information is critical when examining the sort of claim discussed earlier. It appears that childhood circumstances may play some part in determining an individual's socioeconomic circumstances later in life. Can we say this for certain? Is there anything more to consider? These circumstances and others like a "gender gap" in earnings are an important aspect of the distribution of income. We revisit this topic in later chapters.

²Conditional distributions were estimated from data from the first wave of the United Kingdom's Understanding Society dataset (gathered during 2009 and 2010). More details are available at <https://www.understandingsociety.ac.uk/>. Individuals with missing observations are excluded.

KEY CONCEPT

2.3

Means, Variances, and Covariances of Sums of Random Variables

Let X , Y , and V be random variables; let μ_X and σ_X^2 be the mean and variance of X and let σ_{XY} be the covariance between X and Y (and so forth for the other variables); and let a , b , and c be constants. Equations (2.30) through (2.36) follow from the definitions of the mean, variance, and covariance:

$$E(a + bX + cY) = a + b\mu_X + c\mu_Y, \quad (2.30)$$

$$\text{var}(a + bY) = b^2\sigma_Y^2, \quad (2.31)$$

$$\text{var}(aX + bY) = a^2\sigma_X^2 + 2ab\sigma_{XY} + b^2\sigma_Y^2, \quad (2.32)$$

$$E(Y^2) = \sigma_Y^2 + \mu_Y^2, \quad (2.33)$$

$$\text{cov}(a + bX + cV, Y) = b\sigma_{XY} + c\sigma_{VY}, \quad (2.34)$$

$$E(XY) = \sigma_{XY} + \mu_X\mu_Y, \quad (2.35)$$

$$|\text{corr}(X, Y)| \leq 1 \text{ and } |\sigma_{XY}| \leq \sqrt{\sigma_X^2\sigma_Y^2} \text{ (correlation inequality)}. \quad (2.36)$$

The variance of the sum of X and Y is the sum of their variances plus two times their covariance:

$$\text{var}(X + Y) = \text{var}(X) + \text{var}(Y) + 2\text{cov}(X, Y) = \sigma_X^2 + \sigma_Y^2 + 2\sigma_{XY}. \quad (2.37)$$

If X and Y are independent, then the covariance is 0, and the variance of their sum is the sum of their variances:

$$\begin{aligned} \text{var}(X + Y) &= \text{var}(X) + \text{var}(Y) = \sigma_X^2 + \sigma_Y^2 \\ &\text{(if } X \text{ and } Y \text{ are independent)}. \end{aligned} \quad (2.38)$$

Useful expressions for means, variances, and covariances involving weighted sums of random variables are collected in Key Concept 2.3. The results in Key Concept 2.3 are derived in Appendix 2.1.

2.4 The Normal, Chi-Squared, Student t , and F Distributions

The probability distributions most often encountered in econometrics are the normal, chi-squared, Student t , and F distributions.

The Normal Distribution

A continuous random variable with a **normal distribution** has the familiar bell-shaped probability density shown in Figure 2.5. The function defining the normal probability density is given in Appendix 18.1. As Figure 2.5 shows, the normal density with mean μ and variance σ^2 is symmetric around its mean and has 95% of its probability between $\mu - 1.96\sigma$ and $\mu + 1.96\sigma$.

Some special notation and terminology have been developed for the normal distribution. The normal distribution with mean μ and variance σ^2 is expressed concisely as $N(\mu, \sigma^2)$. The **standard normal distribution** is the normal distribution with mean $\mu = 0$ and variance $\sigma^2 = 1$ and is denoted $N(0, 1)$. Random variables that have a $N(0, 1)$ distribution are often denoted Z , and the standard normal cumulative distribution function is denoted by the Greek letter Φ ; accordingly, $\Pr(Z \leq c) = \Phi(c)$, where c is a constant. Values of the standard normal cumulative distribution function are tabulated in Appendix Table 1.

To look up probabilities for a normal variable with a general mean and variance, we must first standardize the variable. For example, suppose Y is distributed $N(1, 4)$ —that is, Y is normally distributed with a mean of 1 and a variance of 4. What is the probability that $Y \leq 2$ —that is, what is the shaded area in Figure 2.6a? The standardized version of Y is Y minus its mean, divided by its standard deviation; that is, $(Y - 1)/\sqrt{4} = \frac{1}{2}(Y - 1)$. Accordingly, the random variable $\frac{1}{2}(Y - 1)$ is normally distributed with mean 0 and variance 1 (see Exercise 2.8); it has the standard normal

FIGURE 2.5 The Normal Probability Density

The normal probability density function with mean μ and variance σ^2 is a bell-shaped curve, centered at μ . The area under the normal p.d.f. between $\mu - 1.96\sigma$ and $\mu + 1.96\sigma$ is 0.95. The normal distribution is denoted $N(\mu, \sigma^2)$.

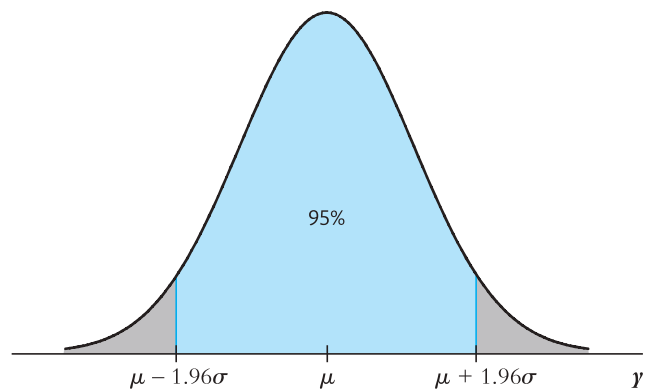
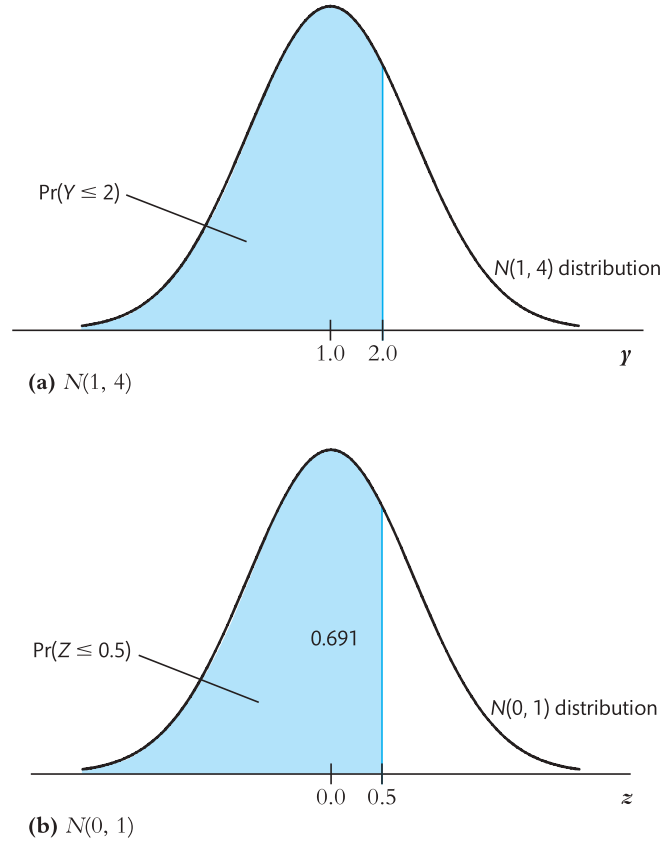


FIGURE 2.6 Calculating the Probability That $Y \leq 2$ When Y Is Distributed $N(1, 4)$

To calculate $\Pr(Y \leq 2)$, standardize Y , then use the standard normal distribution table. Y is standardized by subtracting its mean ($\mu = 1$) and dividing by its standard deviation ($\sigma = 2$). The probability that $Y \leq 2$ is shown in Figure 2.6a, and the corresponding probability after standardizing Y is shown in Figure 2.6b. Because the standardized random variable, $(Y - 1) / 2$, is a standard normal (Z) random variable, $\Pr(Y \leq 2) = \Pr\left(\frac{Y-1}{2} \leq \frac{2-1}{2}\right) = \Pr(Z \leq 0.5)$. From Appendix Table 1, $\Pr(Z \leq 0.5) = \Phi(0.5) = 0.691$.



KEY CONCEPT

2.4

Computing Probabilities and Involving Normal Random Variables

Suppose Y is normally distributed with mean μ and variance σ^2 ; in other words, Y is distributed $N(\mu, \sigma^2)$. Then Y is standardized by subtracting its mean and dividing by its standard deviation, that is, by computing $Z = (Y - \mu) / \sigma$.

Let c_1 and c_2 denote two numbers with $c_1 < c_2$, and let $d_1 = (c_1 - \mu) / \sigma$ and $d_2 = (c_2 - \mu) / \sigma$. Then

$$\Pr(Y \leq c_2) = \Pr(Z \leq d_2) = \Phi(d_2), \tag{2.39}$$

$$\Pr(Y \geq c_1) = \Pr(Z \geq d_1) = 1 - \Phi(d_1), \tag{2.40}$$

$$\Pr(c_1 \leq Y \leq c_2) = \Pr(d_1 \leq Z \leq d_2) = \Phi(d_2) - \Phi(d_1). \tag{2.41}$$

The normal cumulative distribution function Φ is tabulated in Appendix Table 1.

distribution shown in Figure 2.6b. Now $Y \leq 2$ is equivalent to $\frac{1}{2}(Y - 1) \leq \frac{1}{2}(2 - 1)$; that is, $\frac{1}{2}(Y - 1) \leq \frac{1}{2}$. Thus

$$\Pr(Y \leq 2) = \Pr\left[\frac{1}{2}(Y - 1) \leq \frac{1}{2}\right] = \Pr\left(Z \leq \frac{1}{2}\right) = \Phi(0.5) = 0.691, \quad (2.42)$$

where the value 0.691 is taken from Appendix Table 1.

The same approach can be used to compute the probability that a normally distributed random variable exceeds (or is less than) some value or that it falls in a certain range. These steps are discussed in Key Concept 2.4. The box “The Unpegging of the Swiss Franc” presents an unusual application of the cumulative normal distribution.

The normal distribution is symmetric, so its skewness is 0. The kurtosis of the normal distribution is 3.

The multivariate normal distribution. The normal distribution can be generalized to describe the joint distribution of a set of random variables. In this case, the distribution is called the **multivariate normal distribution** or, if only two variables are being considered, the **bivariate normal distribution**. The formula for the bivariate normal p.d.f. is given in Appendix 18.1, and the formula for the general multivariate normal p.d.f. is given in Appendix 19.2.

The multivariate normal distribution has four important properties. If X and Y have a bivariate normal distribution with covariance σ_{XY} and if a and b are two constants, then $aX + bY$ has the normal distribution:

$$aX + bY \text{ is distributed } N(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\sigma_{XY}) \\ (X, Y \text{ bivariate normal}). \quad (2.43)$$

The Unpegging of the Swiss Franc

On Thursday, January 15, 2015, the value of the euro fell by 17.472% from 1.201 to 0.991 against the Swiss franc. This was a huge shift, illustrated in the downward spike in Figure 2.7, given that the previous year had not seen a day’s movement greater than 0.544%. If you had woken up as a statistical analyst for a financial company on that Thursday morning, how might you have estimated the probability of this happening that day?

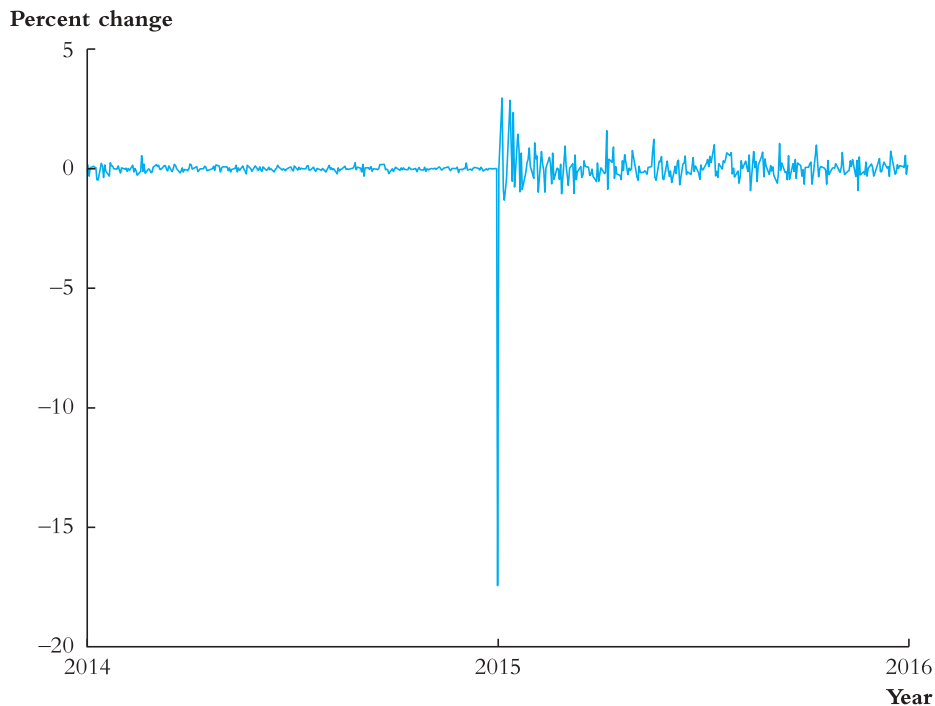
If you had assumed the data was normally distributed, you would have required an estimate of the standard deviation of daily percentage change in the euro/Swiss franc exchange rates. Using Datastream

data¹ for the year to January 14, 2015, you can estimate this as 0.112%.

What was the probability of a drop of 17.472%? We can first calculate the number of standard deviations that describes a change of this magnitude as $\frac{17.472}{0.112} = 156$. If the daily percentage changes are normally distributed, then the estimate of the probability of a fall at least as big as 156 standard deviations corresponds to an inconceivably small number— 8.175×10^{-5288} , which is derived using Equation (2.39).

¹Datastream, maintained by Thomson Reuters, is a global financial and macroeconomic data platform that acts as a repository of financial and economic data.

continued on next page

FIGURE 2.7 Daily Percentage Change in the Euro/Swiss Franc Exchange Rate

The day-on-day percentage change in the value of the euro in Swiss francs for a year before and a year after the unpegging of the Swiss franc on January 15, 2015.

So was the probability of a fall at least this large really so small? Well, no. The error here is to not investigate the nature of our data further, and to fail to understand the actual process that determined the value of the currency. The Swiss franc had in fact been kept within very small bounds due to the actions of the country's central bank in setting a so-called "peg" for the currency. In the previous twelve months, this had been within the range of 1.2008 and 1.236 Swiss francs per euro. In fact, the introduction of this peg over three years earlier had caused an appreciation of the euro against the Swiss franc of over 20 standard deviations (again, assuming a normal distribution derived from previous daily changes!).²

It was the introduction of the peg that had caused such little volatility in—or such a low standard deviation of—the value of the currency. Once this peg was removed, as happened on that particular Thursday, the value of the currency was able to float and vary according to market factors. Investors responded to the removal of the peg by bidding down the value of the euro against the franc substantially.

It is not only the removal of a currency peg in this way that can cause extreme fluctuations. The result

²See the article published in Reuters, "Charts of the Day, Swiss Franc Edition," by Felix Salmon, September 6, 2011.

of the 2016 “Brexit” referendum in the United Kingdom—an event that, while seen as unlikely, was at least partly foreseeable—led to an appreciation in the value of the euro against British pound sterling on June 24, 2016, of 6.17%. This is equivalent to 9.80 standard deviations (based on data from the previous year), or an event with an apparent probability of 5.629×10^{-23} . While it may seem substantially more likely to occur, the probability of such an event actually taking place is less than once every 1,000,000,000,000,000,000 years (a total of 18 zeros)!³ Again, it seems unlikely that this

is an accurate characterization of the probability of such an event occurring.

Clearly, it is dangerous to assume that data is normally distributed or that recent observations of a variable will provide a useful prediction of the range of future values. Indeed, it is partly for this reason that advertisements for financial products in the United Kingdom must carry a disclaimer that “past performance is not a guide to future performance.”

³This is based on the assumption of 260 trading days per year.

More generally, if n random variables have a multivariate normal distribution, then any linear combination of these variables (such as their sum) is normally distributed.

Second, if a set of variables has a multivariate normal distribution, then the marginal distribution of each of the variables is normal [this follows from Equation (2.43) by setting $a = 1$ and $b = 0$].

Third, if variables with a multivariate normal distribution have covariances that equal 0, then the variables are independent. Thus, if X and Y have a bivariate normal distribution and $\sigma_{XY} = 0$, then X and Y are independent (this is shown in Appendix 18.1). In Section 2.3, it was shown that if X and Y are independent, then, regardless of their joint distribution, $\sigma_{XY} = 0$. If X and Y are jointly normally distributed, then the converse is also true. This result—that 0 covariance implies independence—is a special property of the multivariate normal distribution that is not true in general.

Fourth, if X and Y have a bivariate normal distribution, then the conditional expectation of Y given X is linear in X ; that is, $E(Y|X = x) = a + bx$, where a and b are constants (Exercise 18.11). Joint normality implies linearity of conditional expectations, but linearity of conditional expectations does not imply joint normality.

The Chi-Squared Distribution

The chi-squared distribution is used when testing certain types of hypotheses in statistics and econometrics.

The **chi-squared distribution** is the distribution of the sum of m squared independent standard normal random variables. This distribution depends on m , which is called the degrees of freedom of the chi-squared distribution. For example, let Z_1 , Z_2 , and Z_3 be independent standard normal random variables. Then $Z_1^2 + Z_2^2 + Z_3^2$ has a chi-squared distribution with 3 degrees of freedom. The name for this distribution derives from the Greek letter used to denote it: A chi-squared distribution with m degrees of freedom is denoted χ_m^2 .

Selected percentiles of the χ_m^2 distribution are given in Appendix Table 3. For example, Appendix Table 3 shows that the 95th percentile of the χ_3^2 distribution is 7.81, so $\Pr(Z_1^2 + Z_2^2 + Z_3^2 \leq 7.81) = 0.95$.

The Student t Distribution

The **Student t distribution** with m degrees of freedom is defined to be the distribution of the ratio of a standard normal random variable to the square root of an independently distributed chi-squared random variable with m degrees of freedom divided by m . That is, let Z be a standard normal random variable, let W be a random variable with a chi-squared distribution with m degrees of freedom, and let Z and W be independently distributed. Then the random variable $Z / \sqrt{W/m}$ has a Student t distribution (also called the **t distribution**) with m degrees of freedom. This distribution is denoted t_m . Selected percentiles of the Student t distribution are given in Appendix Table 2.

The Student t distribution depends on the degrees of freedom m . Thus the 95th percentile of the t_m distribution depends on the degrees of freedom m . The Student t distribution has a bell shape similar to that of the normal distribution, but it has more mass in the tails; that is, it is a “fatter” bell shape than the normal. When m is 30 or more, the Student t distribution is well approximated by the standard normal distribution, and the t_∞ distribution equals the standard normal distribution.

The F Distribution

The **F distribution** with m and n degrees of freedom, denoted $F_{m,n}$, is defined to be the distribution of the ratio of a chi-squared random variable with degrees of freedom m , divided by m , to an independently distributed chi-squared random variable with degrees of freedom n , divided by n . To state this mathematically, let W be a chi-squared random variable with m degrees of freedom and let V be a chi-squared random variable with n degrees of freedom, where W and V are independently distributed. Then $\frac{W/m}{V/n}$ has an $F_{m,n}$ distribution—that is, an F distribution with numerator degrees of freedom m and denominator degrees of freedom n .

In statistics and econometrics, an important special case of the F distribution arises when the denominator degrees of freedom is large enough that the $F_{m,n}$

distribution can be approximated by the $F_{m,\infty}$ distribution. In this limiting case, the denominator random variable V/n is the mean of infinitely many squared standard normal random variables, and that mean is 1 because the mean of a squared standard normal random variable is 1 (see Exercise 2.24). Thus the $F_{m,\infty}$ distribution is the distribution of a chi-squared random variable with m degrees of freedom divided by m : W/m is distributed $F_{m,\infty}$. For example, from Appendix Table 4, the 95th percentile of the $F_{3,\infty}$ distribution is 2.60, which is the same as the 95th percentile of the χ^2_3 distribution, 7.81 (from Appendix Table 2), divided by the degrees of freedom, which is $3(7.81/3 = 2.60)$.

The 90th, 95th, and 99th percentiles of the $F_{m,n}$ distribution are given in Appendix Table 5 for selected values of m and n . For example, the 95th percentile of the $F_{3,30}$ distribution is 2.92, and the 95th percentile of the $F_{3,90}$ distribution is 2.71. As the denominator degrees of freedom n increases, the 95th percentile of the $F_{3,n}$ distribution tends to the $F_{3,\infty}$ limit of 2.60.

2.5 Random Sampling and the Distribution of the Sample Average

Almost all the statistical and econometric procedures used in this text involve averages or weighted averages of a sample of data. Characterizing the distributions of sample averages therefore is an essential step toward understanding the performance of econometric procedures.

This section introduces some basic concepts about random sampling and the distributions of averages that are used throughout the book. We begin by discussing random sampling. The act of random sampling—that is, randomly drawing a sample from a larger population—has the effect of making the sample average itself a random variable. Because the sample average is a random variable, it has a probability distribution, which is called its sampling distribution. This section concludes with some properties of the sampling distribution of the sample average.

Random Sampling

Simple random sampling. Suppose our commuting student from Section 2.1 aspires to be a statistician and decides to record her commuting times on various days. She selects these days at random from the school year, and her daily commuting time has the cumulative distribution function in Figure 2.2a. Because these days were selected at random, knowing the value of the commuting time on one of these randomly selected days provides no information about the commuting time on another of the days; that is, because the days were selected at random, the values of the commuting time on the different days are independently distributed random variables.

The situation described in the previous paragraph is an example of the simplest sampling scheme used in statistics, called **simple random sampling**, in which n objects are

KEY CONCEPT

Simple Random Sampling and i.i.d. Random Variables

2.5

In a simple random sample, n objects are drawn at random from a population, and each object is equally likely to be drawn. The value of the random variable Y for the i^{th} randomly drawn object is denoted Y_i . Because each object is equally likely to be drawn and the distribution of Y_i is the same for all i , the random variables Y_1, \dots, Y_n are independently and identically distributed (i.i.d.); that is, the distribution of Y_i is the same for all $i = 1, \dots, n$, and Y_1 is distributed independently of Y_2, \dots, Y_n and so forth.

selected at random from a **population** (the population of commuting days) and each member of the population (each day) is equally likely to be included in the sample.

The n observations in the sample are denoted Y_1, \dots, Y_n , where Y_1 is the first observation, Y_2 is the second observation, and so forth. In the commuting example, Y_1 is the commuting time on the first of the n randomly selected days, and Y_i is the commuting time on the i^{th} of the randomly selected days.

Because the members of the population included in the sample are selected at random, the values of the observations Y_1, \dots, Y_n are themselves random. If different members of the population are chosen, their values of Y will differ. Thus the act of random sampling means that Y_1, \dots, Y_n can be treated as random variables. Before they are sampled, Y_1, \dots, Y_n can take on many possible values; after they are sampled, a specific value is recorded for each observation.

i.i.d. draws. Because Y_1, \dots, Y_n are randomly drawn from the same population, the marginal distribution of Y_i is the same for each $i = 1, \dots, n$; this marginal distribution is the distribution of Y in the population being sampled. When Y_i has the same marginal distribution for $i = 1, \dots, n$, then Y_1, \dots, Y_n are said to be **identically distributed**.

Under simple random sampling, knowing the value of Y_1 provides no information about Y_2 , so the conditional distribution of Y_2 given Y_1 is the same as the marginal distribution of Y_2 . In other words, under simple random sampling, Y_1 is distributed independently of Y_2, \dots, Y_n .

When Y_1, \dots, Y_n are drawn from the same distribution and are independently distributed, they are said to be **independently and identically distributed (i.i.d.)**.

Simple random sampling and i.i.d. draws are summarized in Key Concept 2.5.

The Sampling Distribution of the Sample Average

The **sample average** or **sample mean**, \bar{Y} , of the n observations Y_1, \dots, Y_n is

$$\bar{Y} = \frac{1}{n}(Y_1 + Y_2 + \cdots + Y_n) = \frac{1}{n} \sum_{i=1}^n Y_i. \quad (2.44)$$

An essential concept is that the act of drawing a random sample has the effect of making the sample average \bar{Y} a random variable. Because the sample was drawn at random, the value of each Y_i is random. Because Y_1, \dots, Y_n are random, their average is random. Had a different sample been drawn, then the observations and their sample average would have been different: The value of \bar{Y} differs from one randomly drawn sample to the next.

For example, suppose our student commuter selected five days at random to record her commute times, then computed the average of those five times. Had she chosen five different days, she would have recorded five different times—and thus would have computed a different value of the sample average.

Because \bar{Y} is random, it has a probability distribution. The distribution of \bar{Y} is called the **sampling distribution** of \bar{Y} because it is the probability distribution associated with possible values of \bar{Y} that could be computed for different possible samples Y_1, \dots, Y_n .

The sampling distribution of averages and weighted averages plays a central role in statistics and econometrics. We start our discussion of the sampling distribution of \bar{Y} by computing its mean and variance under general conditions on the population distribution of Y .

Mean and variance of \bar{Y} . Suppose that the observations Y_1, \dots, Y_n are i.i.d., and let μ_Y and σ_Y^2 denote the mean and variance of Y_i (because the observations are i.i.d., the mean is the same for all $i = 1, \dots, n$, and so is the variance). When $n = 2$, the mean of the sum $Y_1 + Y_2$ is given by applying Equation (2.29): $E(Y_1 + Y_2) = \mu_Y + \mu_Y = 2\mu_Y$. Thus the mean of the sample average is $E[\frac{1}{2}(Y_1 + Y_2)] = \frac{1}{2} \times 2\mu_Y = \mu_Y$. In general,

$$E(\bar{Y}) = \frac{1}{n} \sum_{i=1}^n E(Y_i) = \mu_Y. \tag{2.45}$$

The variance of \bar{Y} is found by applying Equation (2.38). For example, for $n = 2$, $\text{var}(Y_1 + Y_2) = 2\sigma_Y^2$, so [by applying Equation (2.32) with $a = b = \frac{1}{2}$ and $\text{cov}(Y_1, Y_2) = 0$], $\text{var}(\bar{Y}) = \frac{1}{2}\sigma_Y^2$. For general n , because Y_1, \dots, Y_n are i.i.d., Y_i and Y_j are independently distributed for $i \neq j$, so $\text{cov}(Y_i, Y_j) = 0$. Thus

$$\begin{aligned} \text{var}(\bar{Y}) &= \text{var}\left(\frac{1}{n} \sum_{i=1}^n Y_i\right) \\ &= \frac{1}{n^2} \sum_{i=1}^n \text{var}(Y_i) + \frac{1}{n^2} \sum_{i=1}^n \sum_{j=1, j \neq i}^n \text{cov}(Y_i, Y_j) \\ &= \frac{\sigma_Y^2}{n}. \end{aligned} \tag{2.46}$$

The standard deviation of \bar{Y} is the square root of the variance, σ_Y/\sqrt{n} .

Financial Diversification and Portfolios

The principle of diversification says that you can reduce your risk by holding small investments in multiple assets, compared to putting all your money into one asset. That is, you shouldn't put all your eggs in one basket.

The math of diversification follows from Equation (2.46). Suppose you divide \$1 equally among n assets. Let Y_i represent the payout in one year of \$1 invested in the i^{th} asset. Because you invested $1/n$ dollars in each asset, the actual payoff of your portfolio after one year is $(Y_1 + Y_2 + \cdots + Y_n)/n = \bar{Y}$. To keep things simple, suppose that each asset has the same expected payout, μ_Y , the same variance, σ^2 , and the same positive correlation, ρ , across assets [so that $\text{cov}(Y_i, Y_j) = \rho\sigma^2$]. Then the expected payout is

$E(\bar{Y}) = \mu_Y$, and for large n , the variance of the portfolio payout is $\text{var}(\bar{Y}) = \rho\sigma^2$ (Exercise 2.26). Putting all your money into one asset or spreading it equally across all n assets has the same expected payout, but diversifying reduces the variance from σ^2 to $\rho\sigma^2$.

The math of diversification has led to financial products such as stock mutual funds, in which the fund holds many stocks and an individual owns a share of the fund, thereby owning a small amount of many stocks. But diversification has its limits: For many assets, payouts are positively correlated, so $\text{var}(\bar{Y})$ remains positive even if n is large. In the case of stocks, risk is reduced by holding a portfolio, but that portfolio remains subject to the unpredictable fluctuations of the overall stock market.

In summary, if Y_1, \dots, Y_n are i.i.d., the mean, the variance, and the standard deviation of \bar{Y} are

$$E(\bar{Y}) = \mu_Y, \quad (2.47)$$

$$\text{var}(\bar{Y}) = \sigma_{\bar{Y}}^2 = \frac{\sigma_Y^2}{n}, \text{ and} \quad (2.48)$$

$$\text{std.dev}(\bar{Y}) = \sigma_{\bar{Y}} = \frac{\sigma_Y}{\sqrt{n}}. \quad (2.49)$$

These results hold whatever the distribution of Y is; that is, the distribution of Y does not need to take on a specific form, such as the normal distribution, for Equations (2.47) through (2.49) to hold.

The notation $\sigma_{\bar{Y}}^2$ denotes the variance of the sampling distribution of the sample average \bar{Y} . In contrast, σ_Y^2 is the variance of each individual Y_i , that is, the variance of the population distribution from which the observation is drawn. Similarly, $\sigma_{\bar{Y}}$ denotes the standard deviation of the sampling distribution of \bar{Y} .

Sampling distribution of \bar{Y} when Y is normally distributed. Suppose that Y_1, \dots, Y_n are i.i.d. draws from the $N(\mu_Y, \sigma_Y^2)$ distribution. As stated following Equation (2.43), the sum of n normally distributed random variables is itself normally distributed. Because the mean of \bar{Y} is μ_Y and the variance of \bar{Y} is σ_Y^2/n , this means that, if Y_1, \dots, Y_n are i.i.d. draws from the $N(\mu_Y, \sigma_Y^2)$ distribution, then \bar{Y} is distributed $N(\mu_Y, \sigma_Y^2/n)$.

2.6 Large-Sample Approximations to Sampling Distributions

Sampling distributions play a central role in the development of statistical and econometric procedures, so it is important to know, in a mathematical sense, what the sampling distribution of \bar{Y} is. There are two approaches to characterizing sampling distributions: an “exact” approach and an “approximate” approach.

The exact approach entails deriving a formula for the sampling distribution that holds exactly for any value of n . The sampling distribution that exactly describes the distribution of \bar{Y} for any n is called the **exact distribution** or **finite-sample distribution** of \bar{Y} . For example, if Y is normally distributed and Y_1, \dots, Y_n are i.i.d., then (as discussed in Section 2.5) the exact distribution of \bar{Y} is normal with mean μ_Y and variance σ_Y^2/n . Unfortunately, if the distribution of Y is not normal, then in general the exact sampling distribution of \bar{Y} is very complicated and depends on the distribution of Y .

The approximate approach uses approximations to the sampling distribution that rely on the sample size being large. The large-sample approximation to the sampling distribution is often called the **asymptotic distribution**—“asymptotic” because the approximations become exact in the limit that $n \rightarrow \infty$. As we see in this section, these approximations can be very accurate even if the sample size is only $n = 30$ observations. Because sample sizes used in practice in econometrics typically number in the hundreds or thousands, these asymptotic distributions can be counted on to provide very good approximations to the exact sampling distribution.

This section presents the two key tools used to approximate sampling distributions when the sample size is large: the law of large numbers and the central limit theorem. The law of large numbers says that when the sample size is large, \bar{Y} will be close to μ_Y with very high probability. The central limit theorem says that when the sample size is large, the sampling distribution of the standardized sample average, $(\bar{Y} - \mu_Y) / \sigma_{\bar{Y}}$, is approximately normal.

Although exact sampling distributions are complicated and depend on the distribution of Y , the asymptotic distributions are simple. Moreover—remarkably—the asymptotic normal distribution of $(\bar{Y} - \mu_Y) / \sigma_{\bar{Y}}$ does *not* depend on the distribution of Y . This normal approximate distribution provides enormous simplifications and underlies the theory of regression used throughout this text.

The Law of Large Numbers and Consistency

The **law of large numbers** states that, under general conditions, \bar{Y} will be near μ_Y with very high probability when n is large. This is sometimes called the “law of averages.” When a large number of random variables with the same mean are averaged together, the large values tend to balance the small values, and their sample average is close to their common mean.

For example, consider a simplified version of our student commuter’s experiment in which she simply records whether her commute was short (less than

KEY CONCEPT

2.6

Convergence in Probability, Consistency, and the Law of Large Numbers

The sample average \bar{Y} converges in probability to μ_Y (or, equivalently, \bar{Y} is consistent for μ_Y) if the probability that \bar{Y} is in the range $(\mu_Y - c)$ to $(\mu_Y + c)$ becomes arbitrarily close to 1 as n increases for any constant $c > 0$. The convergence of \bar{Y} to μ_Y in probability is written $\bar{Y} \xrightarrow{p} \mu_Y$.

The law of large numbers says that if Y_1, \dots, Y_n are independently and identically distributed with $E(Y_i) = \mu_Y$ and if large outliers are unlikely (technically if $\text{var}(Y_i) = \sigma_Y^2 < \infty$), then $\bar{Y} \xrightarrow{p} \mu_Y$.

20 minutes) or long. Let $Y_i = 1$ if her commute was short on the i^{th} randomly selected day and $Y_i = 0$ if it was long. Because she used simple random sampling, Y_1, \dots, Y_n are i.i.d. Thus Y_1, \dots, Y_n are i.i.d. draws of a Bernoulli random variable, where (from Table 2.2) the probability that $Y_i = 1$ is 0.78. Because the expectation of a Bernoulli random variable is its success probability, $E(Y_i) = \mu_Y = 0.78$. The sample average \bar{Y} is the fraction of days in her sample in which her commute was short.

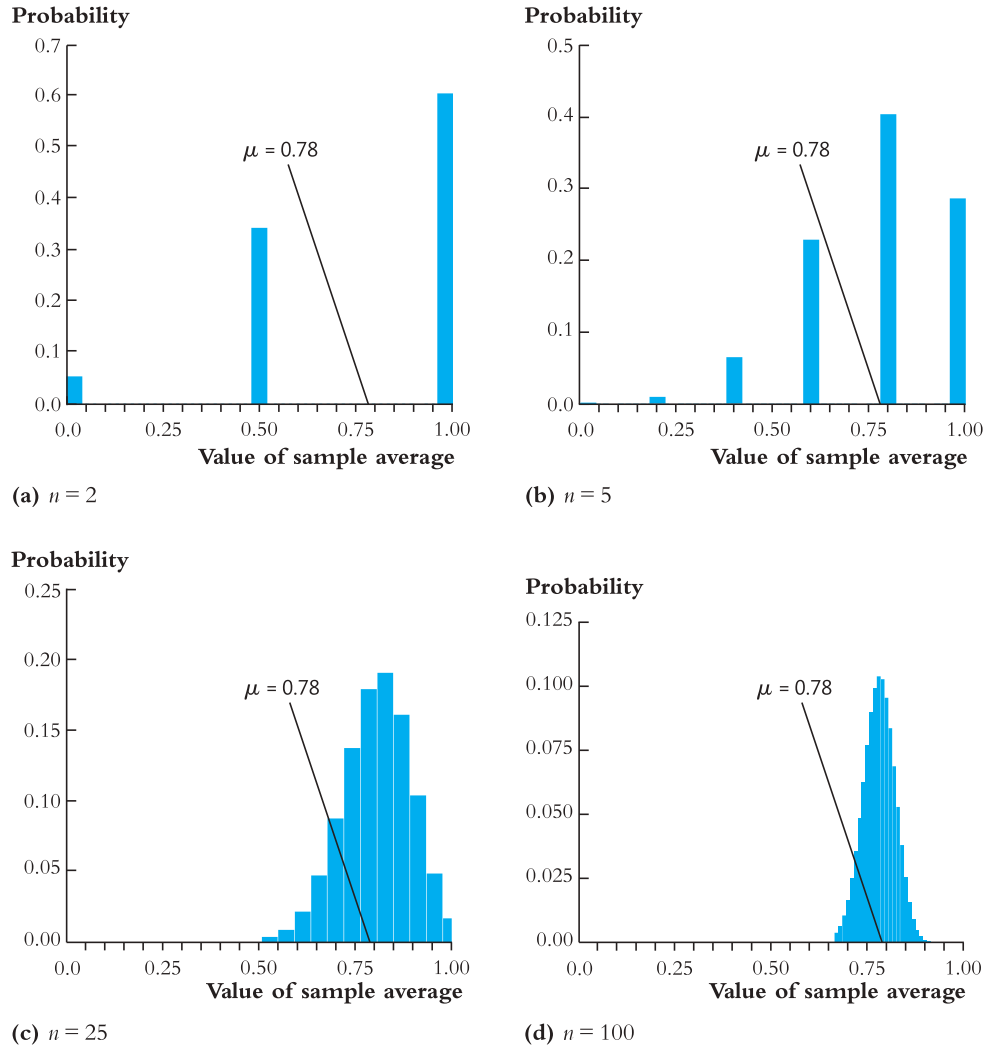
Figure 2.8 shows the sampling distribution of \bar{Y} for various sample sizes n . When $n = 2$ (Figure 2.8a), \bar{Y} can take on only three values: 0, $\frac{1}{2}$, and 1 (neither commute was short, one was short, and both were short), none of which is particularly close to the true proportion in the population, 0.78. As n increases, however (Figures 2.8b–d), \bar{Y} takes on more values, and the sampling distribution becomes tightly centered on μ_Y .

The property that \bar{Y} is near μ_Y with probability increasing to 1 as n increases is called **convergence in probability** or, more concisely, **consistency** (see Key Concept 2.6). The law of large numbers states that under certain conditions \bar{Y} converges in probability to μ_Y or, equivalently, that \bar{Y} is consistent for μ_Y .

The conditions for the law of large numbers that we will use in this text are that Y_1, \dots, Y_n are i.i.d. and that the variance of Y_i , σ_Y^2 , is finite. The mathematical role of these conditions is made clear in Section 18.2, where the law of large numbers is proven. If the data are collected by simple random sampling, then the i.i.d. assumption holds. The assumption that the variance is finite says that extremely large values of Y_i —that is, outliers—are unlikely and are observed infrequently; otherwise, these large values could dominate \bar{Y} , and the sample average would be unreliable. This assumption is plausible for the applications in this text. For example, because there is an upper limit to our student's commuting time (she could park and walk if the traffic is dreadful), the variance of the distribution of commuting times is finite.

The Central Limit Theorem

The **central limit theorem** says that, under general conditions, the distribution of \bar{Y} is well approximated by a normal distribution when n is large. Recall that the mean of Y is μ_Y and its variance is $\sigma_{\bar{Y}}^2 = \sigma_Y^2/n$. According to the central limit theorem, when

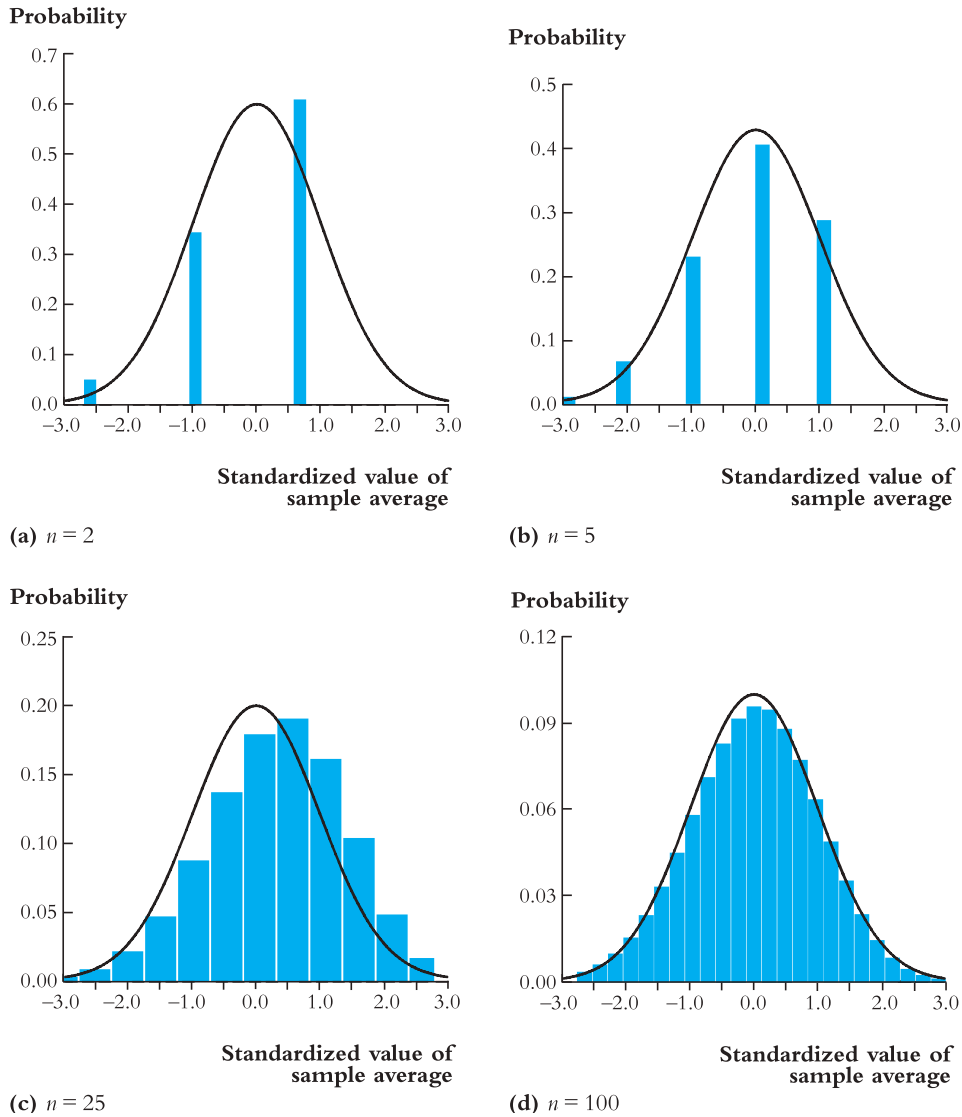
FIGURE 2.8 Sampling Distribution of the Sample Average of n Bernoulli Random Variables

The distributions are the sampling distributions of \bar{Y} , the sample average of n independent Bernoulli random variables with $p = \Pr(Y_i = 1) = 0.78$ (the probability of a short commute is 78%). The variance of the sampling distribution of \bar{Y} decreases as n gets larger, so the sampling distribution becomes more tightly concentrated around its mean, $\mu = 0.78$, as the sample size n increases.

n is large, the distribution of \bar{Y} is approximately $N(\mu_Y, \sigma_{\bar{Y}}^2)$. As discussed at the end of Section 2.5, the distribution of \bar{Y} is *exactly* $N(\mu_Y, \sigma_{\bar{Y}}^2)$ when the sample is drawn from a population with the normal distribution $N(\mu_Y, \sigma_Y^2)$. The central limit theorem says that this same result is *approximately* true when n is large even if Y_1, \dots, Y_n are not themselves normally distributed.

The convergence of the distribution of \bar{Y} to the bell-shaped, normal approximation can be seen (a bit) in Figure 2.8. However, because the distribution gets quite tight for large n , this requires some squinting. It would be easier to see the shape of

FIGURE 2.9 Distribution of the Standardized Sample Average of n Bernoulli Random Variables with $p = 0.78$



The sampling distributions of \bar{Y} in Figure 2.8 are plotted here after standardizing \bar{Y} . Standardization centers the distributions in Figure 2.8 and magnifies the scale on the horizontal axis by a factor of \sqrt{n} . When the sample size is large, the sampling distributions are increasingly well approximated by the normal distribution (the solid line), as predicted by the central limit theorem. The normal distribution is scaled so that the height of the distribution is approximately the same in all figures.

the distribution of \bar{Y} if you used a magnifying glass or had some other way to zoom in or to expand the horizontal axis of the figure.

One way to do this is to standardize \bar{Y} so that it has a mean of 0 and a variance of 1. This process leads to examining the distribution of the standardized version of \bar{Y} , $(\bar{Y} - \mu_Y) / \sigma_{\bar{Y}}$. According to the central limit theorem, this distribution should be well approximated by a $N(0, 1)$ distribution when n is large.

The distribution of the standardized average $(\bar{Y} - \mu_Y) / \sigma_{\bar{Y}}$ is plotted in Figure 2.9 for the distributions in Figure 2.8; the distributions in Figure 2.9 are exactly the same as in Figure 2.8, except that the scale of the horizontal axis is changed so that the standardized variable has a mean of 0 and a variance of 1. After this change of scale, it is easy to see that, if n is large enough, the distribution of \bar{Y} is well approximated by a normal distribution.

One might ask, how large is “large enough”? That is, how large must n be for the distribution of \bar{Y} to be approximately normal? The answer is, “It depends.” The quality of the normal approximation depends on the distribution of the underlying Y_i that make up the average. At one extreme, if the Y_i are themselves normally distributed, then \bar{Y} is exactly normally distributed for all n . In contrast, when the underlying Y_i themselves have a distribution that is far from normal, then this approximation can require $n = 30$ or even more.

This point is illustrated in Figure 2.10 for a population distribution, shown in Figure 2.10a, that is quite different from the Bernoulli distribution. This distribution has a long right tail (it is skewed to the right). The sampling distribution of \bar{Y} , after centering and scaling, is shown in Figures 2.10b–d for $n = 5, 25,$ and $100,$ respectively. Although the sampling distribution is approaching the bell shape for $n = 25,$ the normal approximation still has noticeable imperfections. By $n = 100,$ however, the normal approximation is quite good. In fact, for $n \geq 100,$ the normal approximation to the distribution of \bar{Y} typically is very good for a wide variety of population distributions.

The central limit theorem is a remarkable result. While the “small n ” distributions of \bar{Y} in parts b and c of Figures 2.9 and 2.10 are complicated and quite different from each other, the “large n ” distributions in Figures 2.9d and 2.10d are simple and, amazingly, have a similar shape. Because the distribution of \bar{Y} approaches the normal as n grows large, \bar{Y} is said to have an **asymptotic normal distribution**.

The convenience of the normal approximation, combined with its wide applicability because of the central limit theorem, makes it a key underpinning of applied econometrics. The central limit theorem is summarized in Key Concept 2.7.

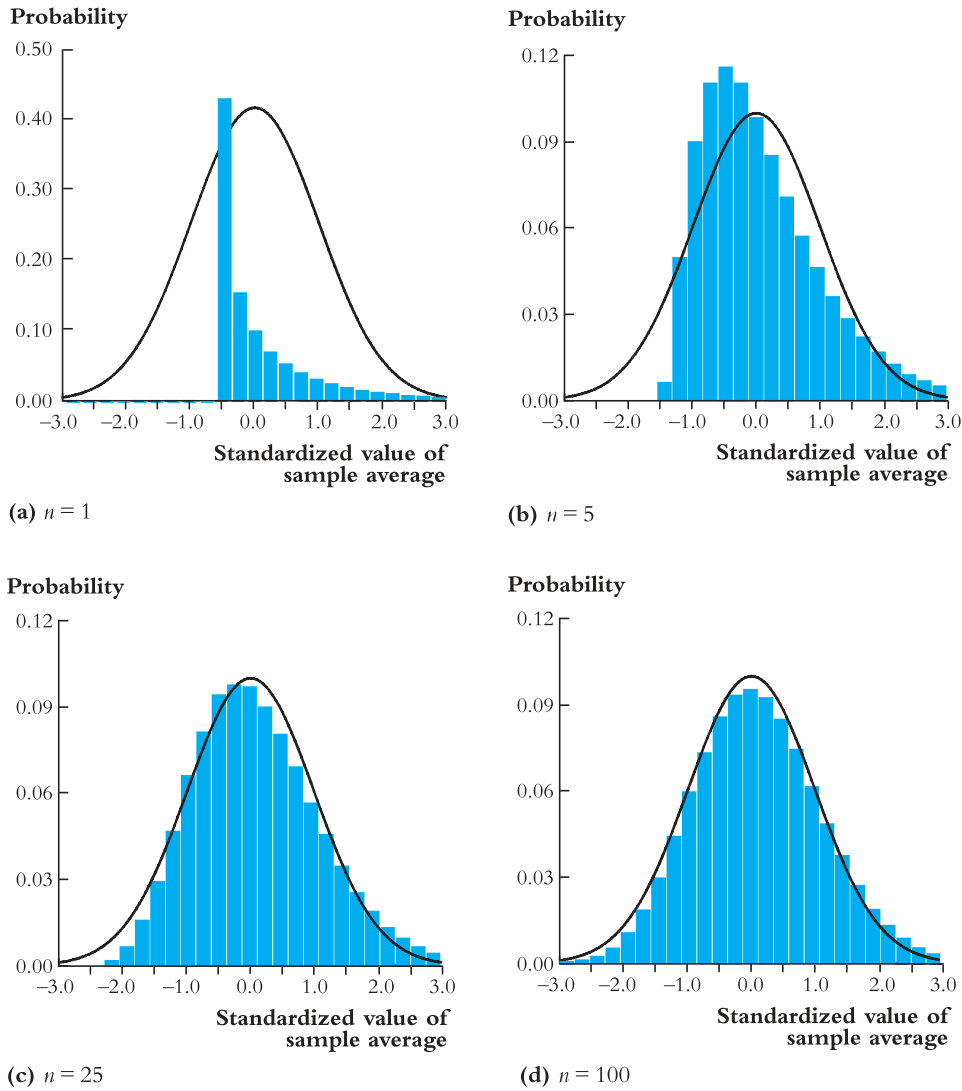
The Central Limit Theorem

KEY CONCEPT

2.7

Suppose that Y_1, \dots, Y_n are i.i.d. with $E(Y_i) = \mu_Y$ and $\text{var}(Y_i) = \sigma_Y^2$, where $0 < \sigma_Y^2 < \infty$. As $n \rightarrow \infty$, the distribution of $(\bar{Y} - \mu_Y) / \sigma_{\bar{Y}}$ (where $\sigma_{\bar{Y}}^2 = \sigma_Y^2 / n$) becomes arbitrarily well approximated by the standard normal distribution.

FIGURE 2.10 Distribution of the Standardized Sample Average of n Draws from a Skewed Population Distribution



The figures show sampling distributions of the standardized sample average of n draws from the skewed (asymmetric) population distribution shown in Figure 2.10a. When n is small ($n = 5$), the sampling distribution, like the population distribution, is skewed. But when n is large ($n = 100$), the sampling distribution is well approximated by a standard normal distribution (solid line), as predicted by the central limit theorem. The normal distribution is scaled so that the height of the distribution is approximately the same in all figures.

Summary

1. The probabilities with which a random variable takes on different values are summarized by the cumulative distribution function, the probability distribution function (for discrete random variables), and the probability density function (for continuous random variables).
2. The expected value of a random variable Y (also called its mean, μ_Y), denoted $E(Y)$, is its probability-weighted average value. The variance of Y is $\sigma_Y^2 = E[(Y - \mu_Y)^2]$, and the standard deviation of Y is the square root of its variance.
3. The joint probabilities for two random variables, X and Y , are summarized by their joint probability distribution. The conditional probability distribution of Y given $X = x$ is the probability distribution of Y , conditional on X taking on the value x .
4. A normally distributed random variable has the bell-shaped probability density in Figure 2.5. To calculate a probability associated with a normal random variable, first standardize the variable, and then use the standard normal cumulative distribution tabulated in Appendix Table 1.
5. Simple random sampling produces n random observations, Y_1, \dots, Y_n , that are independently and identically distributed (i.i.d.).
6. The sample average, \bar{Y} , varies from one randomly chosen sample to the next and thus is a random variable with a sampling distribution. If Y_1, \dots, Y_n are i.i.d., then
 - a. the sampling distribution of \bar{Y} has mean μ_Y and variance $\sigma_{\bar{Y}}^2 = \sigma_Y^2/n$;
 - b. the law of large numbers says that \bar{Y} converges in probability to μ_Y ; and
 - c. the central limit theorem says that the standardized version of \bar{Y} , $(\bar{Y} - \mu_Y) / \sigma_{\bar{Y}}$, has a standard normal distribution [$N(0, 1)$ distribution] when n is large.

Key Terms

outcomes (56)	probability density function (p.d.f.) (58)
probability (56)	density function (58)
sample space (56)	density (58)
event (56)	expected value (60)
discrete random variable (56)	expectation (60)
continuous random variable (56)	mean (60)
probability distribution (56)	variance (61)
cumulative probability distribution (57)	standard deviation (61)
cumulative distribution function (c.d.f.) (57)	moments of a distribution (63)
cumulative distribution (57)	skewness (64)
Bernoulli random variable (58)	kurtosis (64)
Bernoulli distribution (58)	outlier (64)
	leptokurtic (64)

- r^{th} moment (65)
- standardized random variable (65)
- joint probability distribution (65)
- marginal probability distribution (66)
- conditional distribution (66)
- conditional expectation (67)
- conditional mean (67)
- law of iterated expectations (68)
- conditional variance (69)
- Bayes' rule (69)
- independently distributed (70)
- independent (70)
- covariance (70)
- correlation (71)
- uncorrelated (71)
- normal distribution (75)
- standard normal distribution (75)
- multivariate normal distribution (77)
- bivariate normal distribution (77)
- chi-squared distribution (80)
- Student t distribution (80)
- t distribution (80)
- F distribution (80)
- simple random sampling (81)
- population (82)
- identically distributed (82)
- independently and identically distributed (i.i.d.) (82)
- sample average (82)
- sample mean (82)
- sampling distribution (83)
- exact (finite-sample) distribution (85)
- asymptotic distribution (85)
- law of large numbers (85)
- convergence in probability (86)
- consistency (86)
- central limit theorem (86)
- asymptotic normal distribution (89)

MyLab Economics Can Help You Get a Better Grade

MyLab Economics If your exam were tomorrow, would you be ready? For each chapter, **MyLab Economics** Practice Tests and Study Plan help you prepare for your exams. You can also find the Exercises and all Review the Concepts Questions available now in **MyLab Economics**. To see how it works, turn to the **MyLab Economics** spread on the inside front cover of this text and then go to www.pearson.com/mylab/economics.

For additional Empirical Exercises and Data Sets, log on to the Companion Website at www.pearsonglobaleditions.com.

Review the Concepts

- 2.1** Examples of random variables used in this chapter included (a) the sex of the next person you meet, (b) the number of times a wireless network fails, (c) the time it takes to commute to school, and (d) whether it is raining or not. Explain why each can be thought of as random.
- 2.2** Suppose that the random variables X and Y are independent and you know their distributions. Explain why knowing the value of X tells you nothing about the value of Y .
- 2.3** Suppose that X denotes the amount of rainfall in your hometown during a randomly selected month and Y denotes the number of children born in Los Angeles during the same month. Are X and Y independent? Explain.

- 2.4** A math class has 100 students, and the mean student weight is 65 kg. A random sample of five students is selected from the class, and their average weight is calculated. Will the average weight of the students in the sample equal 65 kg? Why or why not? Use this example to explain why the sample average, \bar{Y} , is a random variable.
- 2.5** Suppose that Y_1, \dots, Y_n are i.i.d. random variables with a $N(2, 6)$ distribution. Sketch the probability density of \bar{Y} when $n = 2$. Repeat this for $n = 15$ and $n = 200$. Describe how the densities differ. What is the relationship between your answers and the law of large numbers?
- 2.6** Suppose that Y_1, \dots, Y_n are i.i.d. random variables with probability distribution given in Figure 2.10a. You want to calculate $\Pr(\bar{Y} \leq 0.2)$. Would it be reasonable to use normal approximation if $n = 8$? How about when $n = 30$ and $n = 150$? Explain.
- 2.7** Y is a random variable with $\mu_Y = 0$; $\sigma_Y = 1$, skewness = 0, and kurtosis = 90. Sketch a hypothetical probability distribution of Y . Explain why n random variables drawn from this distribution might have some large outliers.

Exercises

- 2.1** Let Y denote the number of “heads” that occur when two coins are tossed. Assume the probability of a heads is 0.4 on either coin.
- Derive the probability distribution of Y .
 - Derive the mean and variance of Y .
- 2.2** Use the probability distribution given in Table 2.2 to compute (a) $E(Y)$ and $E(X)$; (b) σ_X^2 and σ_Y^2 ; and (c) σ_{XY} and $\text{corr}(X, Y)$.
- 2.3** Using the random variables X and Y from Table 2.2, consider two new random variables, $W = 4 + 8X$ and $V = 11 - 2Y$. Compute (a) $E(W)$ and $E(V)$; (b) σ_W^2 and σ_V^2 ; and (c) σ_{WV} and $\text{corr}(W, V)$.
- 2.4** Suppose X is a Bernoulli random variable with $\Pr(X = 1) = p$.
- Show $E(X^4) = p$.
 - Show $E(X^k) = p$ for $k > 0$.
 - Suppose that $p = 0.53$. Compute the mean, variance, skewness, and kurtosis of X . (*Hint:* You might find it helpful to use the formulas given in Exercise 2.21.)

- 2.5** In July, Lugano's, a city in Switzerland, daily high temperature has a mean of 65°F and a standard deviation of 5°F. What are the mean, standard deviation, and variance in degrees Celsius?
- 2.6** The following table gives the joint probability distribution between employment status and college graduation among those either employed or looking for work (unemployed) in the working-age population of South Africa.

	Unemployed ($Y = 0$)	Employed ($Y = 1$)	Total
Non-college grads ($X = 0$)	0.078	0.673	0.751
College grads ($X = 1$)	0.042	0.207	0.249
Total	0.12	0.88	1.000

- a.** Compute $E(Y)$.
- b.** The unemployment rate is the fraction of the labor force that is unemployed. Show that the unemployment rate is given by $1 - E(Y)$.
- c.** Calculate $E(Y|X = 1)$ and $E(Y|X = 0)$.
- d.** Calculate the unemployment rate for (i) college graduates and (ii) non-college graduates.
- e.** A randomly selected member of this population reports being unemployed. What is the probability that this worker is a college graduate?
A non-college graduate?
- f.** Are educational achievement and employment status independent? Explain.
- 2.7** In a given population of two-earner male-female couples, male earnings have a mean of \$50,000 per year and a standard deviation of \$15,000. Female earnings have a mean of \$48,000 per year and a standard deviation of \$13,000. The correlation between male and female earnings for a couple is 0.90. Let C denote the combined earnings for a randomly selected couple.
- a.** What is the mean of C ?
- b.** What is the covariance between male and female earnings?
- c.** What is the standard deviation of C ?
- d.** Convert the answers to (a) through (c) from U.S. dollars (\$) to euros (€).
- 2.8** The random variable Y has a mean of 4 and a variance of $\frac{1}{9}$. Let $Z = 3(Y - 4)$. Find the mean and the variance of Z .

2.9 X and Y are discrete random variables with the following joint distribution:

		Value of Y				
		2	4	6	8	10
Value of X	3	0.04	0.09	0.03	0.12	0.01
	6	0.10	0.06	0.15	0.03	0.02
	9	0.13	0.11	0.04	0.06	0.01

That is, $\Pr(X = 3, Y = 2) = 0.04$, and so forth.

- Calculate the probability distribution, mean, and variance of Y .
- Calculate the probability distribution, mean, and variance of Y given $X = 6$.
- Calculate the covariance and correlation between X and Y .

2.10 Compute the following probabilities:

- If Y is distributed $N(4, 9)$, find $\Pr(Y \leq 5)$.
- If Y is distributed $N(5, 16)$, find $\Pr(Y > 2)$.
- If Y is distributed $N(1, 4)$, find $\Pr(2 \leq Y \leq 5)$.
- If Y is distributed $N(2, 1)$, find $\Pr(1 \leq Y \leq 4)$.

2.11 Compute the following probabilities:

- If Y is distributed χ_3^2 , find $\Pr(Y \leq 6.25)$.
- If Y is distributed χ_8^2 , find $\Pr(Y \leq 15.51)$.
- If Y is distributed $F_{8, \infty}$, find $\Pr(Y \leq 1.94)$.
- Why are the answers to (b) and (c) the same?
- If Y is distributed χ_1^2 , find $\Pr(Y \leq 0.5)$. (*Hint:* Use the definition of the χ_1^2 distribution.)

2.12 Compute the following probabilities:

- If Y is distributed t_{12} , find $\Pr(Y \leq 1.36)$.
- If Y is distributed t_{30} , find $\Pr(-1.70 \leq Y \leq 1.70)$.
- If Y is distributed $N(0, 1)$, find $\Pr(-1.70 \leq Y \leq 1.70)$.
- When do the critical values of the normal and the t distribution coincide?
- If Y is distributed $F_{4, 11}$, find $\Pr(Y > 3.36)$.
- If Y is distributed $F_{3, 21}$, find $\Pr(Y > 4.87)$.

2.13 X is a Bernoulli random variable with $\Pr(X = 1) = 0.90$; Y is distributed $N(0, 4)$; W is distributed $N(0, 16)$; and X , Y , and W are independent. Let $S = XY + (1 - X)W$. (That is, $S = Y$ when $X = 1$, and $S = W$ when $X = 0$.)

- a. Show that $E(Y^2) = 4$ and $E(W^2) = 16$.
- b. Show that $E(Y^3) = 0$ and $E(W^3) = 0$. (*Hint: What is the skewness for a symmetric distribution?*)
- c. Show that $E(Y^4) = 3 \times 4^2$ and $E(W^4) = 3 \times 16^2$. (*Hint: Use the fact that the kurtosis is 3 for a normal distribution.*)
- d. Derive $E(S)$, $E(S^2)$, $E(S^3)$, and $E(S^4)$. (*Hint: Use the law of iterated expectations conditioning on $X = 0$ and $X = 1$.*)
- e. Derive the skewness and kurtosis for S .
- 2.14** In a population, $\mu_Y = 50$ and $\sigma_Y^2 = 21$. Use the central limit theorem to answer the following questions:
- a. In a random sample of size $n = 50$, find $\Pr(\bar{Y} \leq 51)$.
- b. In a random sample of size $n = 150$, find $\Pr(\bar{Y} > 49)$.
- c. In a random sample of size $n = 45$, find $\Pr(50.5 \leq \bar{Y} \leq 51)$.
- 2.15** Suppose $Y_i, i = 1, 2, \dots, n$ are i.i.d. random variables, each distributed $N(20, 4)$.
- a. Compute $\Pr(19.6 \leq \bar{Y} \leq 20.4)$ when (i) $n = 25$, (ii) $n = 100$, and (iii) $n = 800$.
- b. Suppose c is a positive number. Show that $\Pr(20 - c \leq \bar{Y} \leq 20 + c)$ becomes close to 1.0 as n grows large.
- c. Use your answer in (b) to argue that Y converges in probability to 20.
- 2.16** Y is distributed $N(10, 100)$ and you want to calculate $\Pr(Y \leq 5.8)$. Unfortunately, you do not have your textbook, and do not have access to a normal probability table like Appendix Table 1. However, you do have your computer and a computer program that can generate i.i.d. draws from the $N(10, 100)$ distribution. Explain how you can use your computer to compute an accurate approximation for $\Pr(Y \leq 5.8)$.
- 2.17** $Y_i, i = 1, \dots, n$, are i.i.d. Bernoulli random variables with $p = 0.6$. Let \bar{Y} denote the sample mean.
- a. Use the central limit theorem to compute approximations for
- i. $\Pr(Y \geq 0.64)$ when $n = 50$.
- ii. $\Pr(Y \leq 0.56)$ when $n = 200$.
- b. How large would n need to be to ensure that $\Pr(0.65 > Y > 0.55) = 0.95$? (Use the central limit theorem to compute an approximate answer.)
- 2.18** In any year, the weather can inflict storm damage to a home. From year to year, the damage is random. Let Y denote the dollar value of damage in any given year. Suppose that in 95% of the years $Y = \$0$, but in 5% of the years $Y = \$30,000$.

- a. What are the mean and standard deviation of the damage in any year?
- b. Consider an “insurance pool” of 120 people whose homes are sufficiently dispersed so that, in any year, the damage to different homes can be viewed as independently distributed random variables. Let \bar{Y} denote the average damage to these 120 homes in a year. (i) What is the expected value of the average damage \bar{Y} ? (ii) What is the probability that \bar{Y} exceeds \$3,000?
- 2.19** Consider two random variables, X and Y . Suppose that Y takes on k values y_1, \dots, y_k and that X takes on l values x_1, \dots, x_l .
- a. Show that $\Pr(Y = y_j) = \sum_{i=1}^l \Pr(Y = y_j | X = x_i) \Pr(X = x_i)$. [Hint: Use the definition of $\Pr(Y = y_j | X = x_i)$.]
- b. Use your answer to (a) to verify Equation (2.19).
- c. Suppose that X and Y are independent. Show that $\sigma_{XY} = 0$ and $\text{corr}(X, Y) = 0$.
- 2.20** Consider three random variables, X , Y , and Z . Suppose that Y takes on k values y_1, \dots, y_k ; that X takes on l values x_1, \dots, x_l ; and that Z takes on m values z_1, \dots, z_m . The joint probability distribution of X , Y , Z is $\Pr(X = x, Y = y, Z = z)$, and the conditional probability distribution of Y given X and Z is $\Pr(Y = y | X = x, Z = z) = \frac{\Pr(Y = y, X = x, Z = z)}{\Pr(X = x, Z = z)}$.
- a. Explain how the marginal probability that $Y = y$ can be calculated from the joint probability distribution. [Hint: This is a generalization of Equation (2.16).]
- b. Show that $E(Y) = E[E(Y | X, Z)]$. [Hint: This is a generalization of Equations (2.19) and (2.20).]
- 2.21** X is a random variable with moments $E(X)$, $E(X^2)$, $E(X^3)$, and so forth.
- a. Show $E(X - \mu)^3 = E(X^3) - 3[E(X^2)][E(X)] + 2[E(X)]^3$.
- b. Show $E(X - \mu)^4 = E(X^4) - 4[E(X)][E(X^3)] + 6[E(X)]^2[E(X^2)] - 3[E(X)]^4$.
- 2.22** Suppose you have some money to invest, for simplicity \$1, and you are planning to put a fraction w into a stock market mutual fund and the rest, $1 - w$, into a mutual fund. Suppose that \$1 invested in a stock fund yields R_s after one year and that \$1 invested in mutual fund yields R_b . Suppose that R_s is random with mean 0.06 and standard deviation 0.09, and suppose that R_b is random with mean 0.04 and standard deviation 0.05. The correlation between R_s and R_b is 0.3. If you place a fraction w of your money in the stock fund and the rest, $1 - w$, in the mutual fund, then the return on your investment is $R = wR_s + (1 - w)R_b$.
- a. Suppose that $w = 0.2$. Compute the mean and standard deviation of R .

- b.** Suppose that $w = 0.8$. Compute the mean and standard deviation of R .
- c.** What value of w makes the mean of R as large as possible? What is the standard deviation of R for this value of w ?
- d.** (Harder) What is the value of w that minimizes the standard deviation of R ? (Show using a graph, algebra, or calculus.)
- 2.23** This exercise provides an example of a pair of random variables, X and Y , for which the conditional mean of Y given X depends on X but $\text{corr}(X, Y) = 0$. Let X and Z be two independently distributed standard normal random variables, and let $Y = X^2 + Z$.
- a.** Show that $E(Y|X) = X^2$.
- b.** Show that $\mu_Y = 1$.
- c.** Show that $E(XY) = 0$. (*Hint:* Use the fact that the odd moments of a standard normal random variable are all 0.)
- d.** Show that $\text{cov}(X, Y) = 0$ and thus $\text{corr}(X, Y) = 0$.
- 2.24** Suppose Y_i is distributed i.i.d. $N(0, \sigma^2)$ for $i = 1, 2, \dots, n$.
- a.** Show that $E(Y_i^2/\sigma^2) = 1$.
- b.** Show that $W = (1/\sigma^2)\sum_{i=1}^n Y_i^2$ is distributed χ_n^2 .
- c.** Show that $E(W) = n$. [*Hint:* Use your answer to (a).]
- d.** Show that $V = Y_1/\sqrt{\frac{\sum_{i=2}^n Y_i^2}{n-1}}$ is distributed t_{n-1} .
- 2.25** (Review of summation notation) Let x_1, \dots, x_n denote a sequence of numbers; y_1, \dots, y_n denote another sequence of numbers; and a, b , and c denote three constants. Show that
- a.** $\sum_{i=1}^n ax_i = a \sum_{i=1}^n x_i$,
- b.** $\sum_{i=1}^n (x_i + y_i) = \sum_{i=1}^n x_i + \sum_{i=1}^n y_i$,
- c.** $\sum_{i=1}^n a = n \times a$, and
- d.** $\sum_{i=1}^n (a + bx_i + cy_i)^2 = na^2 + b^2 \sum_{i=1}^n x_i^2 + c^2 \sum_{i=1}^n y_i^2 + 2ab \sum_{i=1}^n x_i + 2ac \sum_{i=1}^n y_i + 2bc \sum_{i=1}^n x_i y_i$.
- 2.26** Suppose that Y_1, Y_2, \dots, Y_n are random variables with a common mean μ_Y ; a common variance σ_Y^2 ; and the same correlation ρ (so that the correlation between Y_i and Y_j is equal to ρ for all pairs i and j , where $i \neq j$).
- a.** Show that $\text{cov}(Y_i, Y_j) = \rho\sigma_Y^2$ for $i \neq j$.
- b.** Suppose that $n = 2$. Show that $E(\bar{Y}) = \mu_Y$ and $\text{var}(\bar{Y}) = \frac{1}{2}\sigma_Y^2 + \frac{1}{2}\rho\sigma_Y^2$.

- c. For $n \geq 2$, show that $E(\bar{Y}) = \mu_Y$ and $\text{var}(\bar{Y}) = \sigma_Y^2/n + [(n-1)/n]\rho\sigma_Y^2$.
- d. When n is very large, show that $\text{var}(\bar{Y}) \approx \rho\sigma_Y^2$.
- 2.27** Consider the problem of predicting Y using another variable, X , so that the prediction of Y is some function of X , say $g(X)$. Suppose that the quality of the prediction is measured by the squared prediction error made on average over all predictions, that is, by $E\{[Y - g(X)]^2\}$. This exercise provides a non-calculus proof that of all possible prediction functions g , the best prediction is made by the conditional expectation, $E(Y|X)$.
- a. Let $\hat{Y} = E(Y|X)$, and let $u = Y - \hat{Y}$ denote its prediction error. Show that $E(u) = 0$. (*Hint:* Use the law of iterated expectations.)
- b. Show that $E(uX) = 0$.
- c. Let $\tilde{Y} = g(X)$ denote a different prediction of Y using X , and let $v = Y - \tilde{Y}$ denote its error. Show that $E[(Y - \tilde{Y})^2] > E[(Y - \hat{Y})^2]$. [*Hint:* Let $h(X) = g(X) - E(Y|X)$, so that $v = [Y - E(Y|X)] - h(X)$. Derive $E(v^2)$.]
- 2.28** Refer to Part B of Table 2.3 for the conditional distribution of the number of network failures M given network age A . Let $\Pr(A = 0) = 0.5$; that is, you work in your room 50% of the time.
- a. Compute the probability of three network failures, $\Pr(M = 3)$.
- b. Use Bayes' rule to compute $\Pr(A = 0 | M = 3)$.
- c. Now suppose you work in your room one-fourth of the time, so $\Pr(A = 0) = 0.25$. Use Bayes' rule to compute $\Pr(A = 0 | M = 3)$.

Empirical Exercise

- E2.1** On the text website, <http://www.pearsonglobaleditions.com>, you will find the spreadsheet **Age_HourlyEarnings**, which contains the joint distribution of age (Age) and average hourly earnings (AHE) for 25- to 34-year-old full-time workers in 2015 with an education level that exceeds a high school diploma. Use this joint distribution to carry out the following exercises. (*Note:* For these exercises, you need to be able to carry out calculations and construct charts using a spreadsheet.)
- a. Compute the marginal distribution of Age .
- b. Compute the mean of AHE for each value of Age ; that is, compute, $E(AHE|Age = 25)$, and so forth.
- c. Compute and plot the mean of AHE versus Age . Are average hourly earnings and age related? Explain.

- d. Use the law of iterated expectations to compute the mean of AHE ; that is, compute $E(AHE)$.
- e. Compute the variance of AHE .
- f. Compute the covariance between AHE and Age .
- g. Compute the correlation between AHE and Age .
- h. Relate your answers in (f) and (g) to the plot you constructed in (c).

APPENDIX

2.1 Derivation of Results in Key Concept 2.3

This appendix derives the equations in Key Concept 2.3.

Equation (2.30) follows from the definition of the expectation.

To derive Equation (2.31), use the definition of the variance to write $\text{var}(a + bY) = E\{[a + bY - E(a + bY)]^2\} = E\{[b(Y - \mu_Y)]^2\} = b^2E[(Y - \mu_Y)^2] = b^2\sigma_Y^2$.

To derive Equation (2.32), use the definition of the variance to write

$$\begin{aligned}
 \text{var}(aX + bY) &= E\{[(aX + bY) - (a\mu_X + b\mu_Y)]^2\} \\
 &= E\{[a(X - \mu_X) + b(Y - \mu_Y)]^2\} \\
 &= E[a^2(X - \mu_X)^2] + 2E[ab(X - \mu_X)(Y - \mu_Y)] \\
 &\quad + E[b^2(Y - \mu_Y)^2] \\
 &= a^2\text{var}(X) + 2ab\text{cov}(X, Y) + b^2\text{var}(Y) \\
 &= a^2\sigma_X^2 + 2ab\sigma_{XY} + b^2\sigma_Y^2,
 \end{aligned} \tag{2.50}$$

where the second equality follows by collecting terms, the third equality follows by expanding the quadratic, and the fourth equality follows by the definition of the variance and covariance.

To derive Equation (2.33), write

$$E(Y^2) = E\{(Y - \mu_Y) + \mu_Y\}^2 = E[(Y - \mu_Y)^2] + 2\mu_Y E(Y - \mu_Y) + \mu_Y^2 = \sigma_Y^2 + \mu_Y^2$$

because $E(Y - \mu_Y) = 0$.

To derive Equation (2.34), use the definition of the covariance to write

$$\begin{aligned}
 \text{cov}(a + bX + cV, Y) &= E\{[a + bX + cV - E(a + bX + cV)][Y - \mu_Y]\} \\
 &= E\{[b(X - \mu_X) + c(V - \mu_V)][Y - \mu_Y]\} \\
 &= E\{[b(X - \mu_X)][Y - \mu_Y]\} + E\{[c(V - \mu_V)][Y - \mu_Y]\} \\
 &= b\sigma_{XY} + c\sigma_{VY},
 \end{aligned} \tag{2.51}$$

which is Equation (2.34).

To derive Equation (2.35), write

$$\begin{aligned} E(XY) &= E\{[(X - \mu_X) + \mu_X][(Y - \mu_Y) + \mu_Y]\} \\ &= E[(X - \mu_X)(Y - \mu_Y)] + \mu_X E(Y - \mu_Y) + \mu_Y E(X - \mu_X) + \mu_X \mu_Y \\ &= \sigma_{XY} + \mu_X \mu_Y. \end{aligned}$$

We now prove the correlation inequality in Equation (2.36); that is, $|\text{corr}(X, Y)| \leq 1$. Let $a = -\sigma_{XY}/\sigma_X^2$ and $b = 1$. Applying Equation (2.32), we have,

$$\begin{aligned} \text{var}(aX + Y) &= a^2\sigma_X^2 + \sigma_Y^2 + 2a\sigma_{XY} \\ &= (-\sigma_{XY}/\sigma_X^2)^2\sigma_X^2 + \sigma_Y^2 + 2(-\sigma_{XY}/\sigma_X^2)\sigma_{XY} \\ &= \sigma_Y^2 - \sigma_{XY}^2/\sigma_X^2. \end{aligned} \tag{2.52}$$

Because $\text{var}(aX + Y)$ is a variance, it cannot be negative, so from the final line of Equation (2.52), it must be that $\sigma_Y^2 - \sigma_{XY}^2/\sigma_X^2 \geq 0$. Rearranging this inequality yields

$$\sigma_{XY}^2 \leq \sigma_X^2\sigma_Y^2 \text{ (covariance inequality)}. \tag{2.53}$$

The covariance inequality implies that $\sigma_{XY}^2/(\sigma_X^2\sigma_Y^2) \leq 1$ or, equivalently, $|\sigma_{XY}/(\sigma_X\sigma_Y)| \leq 1$, which (using the definition of the correlation) proves the correlation inequality, $|\text{corr}(X, Y)| \leq 1$.

APPENDIX

2.2 The Conditional Mean as the Minimum Mean Squared Error Predictor

At a general level, the statistical prediction problem is, how does one best use the information in a random variable X to predict the value of another random variable Y ?

To answer to this question, we must first make precise mathematically what it means for one prediction to be better than another. A common way to do so is to consider the cost of making a prediction error. This cost, which is called the prediction loss, depends on the magnitude of the prediction error. For example, if your job is to predict sales so that a production supervisor can develop a production schedule, being off by a small amount is unlikely to inconvenience customers or to disrupt the production process. But if you are off by a large amount and production is set far too low, your company might lose customers who need to wait a long time to receive a product they order, or if production is far too high, the company will have costly excess inventory on its hands. In either case, a large prediction error can be disproportionately more costly than a small one.

One way to make this logic precise is to let the cost of a prediction error depend on the square of that error, so an error twice as large is four times as costly. Specifically, suppose that your prediction of Y , given the random variable X , is $g(X)$. The prediction error is $Y - g(X)$, and the quadratic loss associated with this prediction is,

$$\text{Loss} = E\{[Y - g(X)]^2\}. \quad (2.54)$$

We now show that, of all possible functions $g(X)$, the loss in Equation (2.54) is minimized by $g(X) = E(Y|X)$. We show this result using discrete random variables, however this result extends to continuous random variables. The proof here uses calculus; Exercise 2.27 works through a non-calculus proof of this result.

First consider the simpler problem of finding a number, m , that minimizes $E[(Y - m)^2]$. From the definition of the expectation, $E[(Y - m)^2] = \sum_{i=1}^k (Y_i - m)^2 p_i$. To find the value of m that minimizes $E[(Y - m)^2]$, take the derivative of $\sum_{i=1}^k (Y_i - m)^2 p_i$ with respect to m and set it to zero:

$$\begin{aligned} \frac{d}{dm} \sum_{i=1}^k (Y_i - m)^2 p_i &= -2 \sum_{i=1}^k (Y_i - m) p_i = -2 \left(\sum_{i=1}^k Y_i p_i - m \sum_{i=1}^k p_i \right) \\ &= -2 \left(\sum_{i=1}^k Y_i p_i - m \right) = 0, \end{aligned} \quad (2.55)$$

where the final equality uses the fact that probabilities sum to 1. It follows from the final equality in Equation (2.55) that the squared error prediction loss is minimized by $m = \sum_{i=1}^k Y_i p_i = E(Y)$, that is, by setting m equal to the mean of Y .

To find the predictor $g(X)$ that minimizes the loss in Equation (2.54), use the law of iterated expectations to write that loss as, $\text{Loss} = E\{[Y - g(X)]^2\} = E(E\{[Y - g(X)]^2|X\})$. Thus, if the function $g(X)$ minimizes $E\{[Y - g(X)]^2|X = x\}$ for each value of x , it minimizes the loss in Equation (2.54). But for a fixed value $X = x$, $g(X) = g(x)$ is a fixed number, so this problem is the same as the one just solved, and the loss is minimized by choosing $g(x)$ to be the mean of Y , given $X = x$. This is true for every value of x . Thus the squared error loss in Equation (2.54) is minimized by $g(X) = E(Y|X)$.

Review of Statistics

Statistics is the science of using data to learn about the world around us. Statistical tools help us answer questions about unknown characteristics of distributions in populations of interest. For example, what is the mean of the distribution of earnings of recent college graduates? Do mean earnings differ for men and women and, if so, by how much?

These questions relate to the distribution of earnings in the population of workers. One way to answer these questions would be to perform an exhaustive survey of the population of workers, measuring the earnings of each worker and thus finding the population distribution of earnings. In practice, however, such a comprehensive survey would be extremely expensive. Comprehensive surveys that do exist, also known as censuses, are often undertaken periodically (for example, every ten years in India, the United States of America and the United Kingdom). This is because the process of conducting a census is an extraordinary commitment, consisting of designing census forms, managing and conducting surveys, and compiling and analyzing data. Censuses across the world have a long history, with accounts of censuses recorded by Babylonians in 4000 BC. According to historians, censuses have been conducted as far back as Ancient Rome; the Romans would track the population by making people return to their birthplace every year in order to be counted.¹ In England and other parts of Wales, a notable census was the Domesday Book, which was compiled in 1086 by William the Conqueror. The U.K. census in its current form dates back to 1801 after essays by economist Thomas Malthus (1798) inspired parliament to want to accurately know the size of the population. Over time the census has evolved from amounting to a mere headcount to the much more ambitious survey of the 2011 U.K. census costing an estimated £482 million. In India, there are accounts of censuses recorded around 300 BC, but the census in its current form has been undertaken since 1872 and every ten years since 1881. In comparison to the U.K. census of 2011, the most recent census of India, also conducted in 2011, approximately cost a mere ₹2200 crore (US\$320 million)! Despite the considerable efforts made to ensure that the census records all individuals, many people slip through the cracks and are not surveyed. Thus a different, more practical approach is needed.

The key insight of statistics is that one can learn about a population distribution by selecting a random sample from that population. Rather than survey the entire population of China (1.4 billion in 2018), we might survey, say, 1000 members of the population, selected at random by simple random sampling. Using statistical methods, we

¹Source: Office for National Statistics, <https://www.ons.gov.uk>, accessed on August 23, 2018.

can use this sample to reach tentative conclusions—to draw statistical inferences—about characteristics of the full population.²

Three types of statistical methods are used throughout econometrics: estimation, hypothesis testing, and confidence intervals. Estimation entails computing a “best guess” numerical value for an unknown characteristic of a population distribution, such as its mean, from a sample of data. Hypothesis testing entails formulating a specific hypothesis about the population and then using sample evidence to decide whether it is true. Confidence intervals use a set of data to estimate an interval or range for an unknown population characteristic. Sections 3.1, 3.2, and 3.3 review estimation, hypothesis testing, and confidence intervals in the context of statistical inference about an unknown population mean.

Most of the interesting questions in economics involve relationships between two or more variables or comparisons between different populations. For example, is there a gap between the mean earnings for male and female recent college graduates? In Section 3.4, the methods for learning about the mean of a single population in Sections 3.1 through 3.3 are extended to compare means in two different populations. Section 3.5 discusses how the methods for comparing the means of two populations can be used to estimate causal effects in experiments. Sections 3.2 through 3.5 focus on the use of the normal distribution for performing hypothesis tests and for constructing confidence intervals when the sample size is large. In some special circumstances, hypothesis tests and confidence intervals can be based on the Student t distribution instead of the normal distribution; these special circumstances are discussed in Section 3.6. The chapter concludes with a discussion of the sample correlation and scatterplots in Section 3.7.

3.1 Estimation of the Population Mean

Suppose you want to know the mean value of Y (that is, μ_Y) in a population, such as the mean earnings of women recently graduated from college. A natural way to estimate this mean is to compute the sample average \bar{Y} from a sample of n independently and identically distributed (i.i.d.) observations, Y_1, \dots, Y_n (recall that Y_1, \dots, Y_n are i.i.d. if they are collected by simple random sampling). This section discusses estimation of μ_Y and the properties of \bar{Y} as an estimator of μ_Y .

Estimators and Their Properties

Estimators. The sample average \bar{Y} is a natural way to estimate μ_Y , but it is not the only way. For example, another way to estimate μ_Y is simply to use the first observation, Y_1 . Both \bar{Y} and Y_1 are functions of the data that are designed to estimate μ_Y ; using the terminology in Key Concept 3.1, both are estimators of μ_Y . When evaluated in repeated samples, \bar{Y} and Y_1 take on different values (they produce

²Estimates of the ‘live’ population of China can be found here using the ‘official’ China Population Clock: <http://data.stats.gov.cn/english/>

Estimators and Estimates

KEY CONCEPT

3.1

An **estimator** is a function of a sample of data to be drawn randomly from a population. An **estimate** is the numerical value of the estimator when it is actually computed using data from a specific sample. An estimator is a random variable because of randomness in selecting the sample, while an estimate is a nonrandom number.

different estimates) from one sample to the next. Thus the estimators \bar{Y} and Y_1 both have sampling distributions. There are, in fact, many estimators of μ_Y , of which \bar{Y} and Y_1 are two examples.

There are many possible estimators, so what makes one estimator “better” than another? Because estimators are random variables, this question can be phrased more precisely: What are desirable characteristics of the sampling distribution of an estimator? In general, we would like an estimator that gets as close as possible to the unknown true value, at least in some average sense; in other words, we would like the sampling distribution of an estimator to be as tightly centered on the unknown value as possible. This observation leads to three specific desirable characteristics of an estimator: unbiasedness (a lack of bias), consistency, and efficiency.

Unbiasedness. Suppose you evaluate an estimator many times over repeated randomly drawn samples. It is reasonable to hope that, on average, you would get the right answer. Thus a desirable property of an estimator is that the mean of its sampling distribution equals μ_Y ; if so, the estimator is said to be unbiased.

To state this concept mathematically, let $\hat{\mu}_Y$ denote some estimator of μ_Y , such as \bar{Y} or Y_1 . [The caret (^) notation will be used throughout this text to denote an estimator, so $\hat{\mu}_Y$ is an estimator of μ_Y .] The estimator $\hat{\mu}_Y$ is unbiased if $E(\hat{\mu}_Y) = \mu_Y$, where $E(\hat{\mu}_Y)$ is the mean of the sampling distribution of $\hat{\mu}_Y$; otherwise, $\hat{\mu}_Y$ is biased.

Bias, Consistency, and Efficiency

KEY CONCEPT

3.2

Let $\hat{\mu}_Y$ be an estimator of μ_Y . Then:

- The *bias* of $\hat{\mu}_Y$ is $E(\hat{\mu}_Y) - \mu_Y$.
- $\hat{\mu}_Y$ is an *unbiased estimator* of μ_Y if $E(\hat{\mu}_Y) = \mu_Y$.
- $\hat{\mu}_Y$ is a *consistent estimator* of μ_Y if $\hat{\mu}_Y \xrightarrow{p} \mu_Y$.
- Let $\tilde{\mu}_Y$ be another estimator of μ_Y , and suppose that both $\hat{\mu}_Y$ and $\tilde{\mu}_Y$ are unbiased. Then $\hat{\mu}_Y$ is said to be more *efficient* than $\tilde{\mu}_Y$ if $\text{var}(\hat{\mu}_Y) < \text{var}(\tilde{\mu}_Y)$.

Consistency. Another desirable property of an estimator μ_Y is that when the sample size is large, the uncertainty about the value of μ_Y arising from random variations in the sample is very small. Stated more precisely, a desirable property of $\hat{\mu}_Y$ is that the probability that it is within a small interval of the true value μ_Y approaches 1 as the sample size increases; that is, $\hat{\mu}_Y$ is consistent for μ_Y (Key Concept 2.6).

Variance and efficiency. Suppose you have two candidate estimators, $\hat{\mu}_Y$ and $\tilde{\mu}_Y$, both of which are unbiased. How might you choose between them? One way to do so is to choose the estimator with the tightest sampling distribution. This suggests choosing between $\hat{\mu}_Y$ and $\tilde{\mu}_Y$ by picking the estimator with the smallest variance. If $\hat{\mu}_Y$ has a smaller variance than $\tilde{\mu}_Y$, then $\hat{\mu}_Y$ is said to be more efficient than $\tilde{\mu}_Y$. The terminology “efficiency” stems from the notion that if $\hat{\mu}_Y$ has a smaller variance than $\tilde{\mu}_Y$, then it uses the information in the data more efficiently than does $\tilde{\mu}_Y$.

Bias, consistency, and efficiency are summarized in Key Concept 3.2.

Properties of \bar{Y}

How does \bar{Y} fare as an estimator of μ_Y when judged by the three criteria of bias, consistency, and efficiency?

Bias and consistency. The sampling distribution of \bar{Y} has already been examined in Sections 2.5 and 2.6. As shown in Section 2.5, $E(\bar{Y}) = \mu_Y$, so \bar{Y} is an unbiased estimator of μ_Y . Similarly, the law of large numbers (Key Concept 2.6) states that $\bar{Y} \xrightarrow{p} \mu_Y$; that is, \bar{Y} is consistent.

Efficiency. What can be said about the efficiency of \bar{Y} ? Because efficiency entails a comparison of estimators, we need to specify the estimator or estimators to which \bar{Y} is to be compared.

We start by comparing the efficiency of \bar{Y} to the estimator Y_1 . Because Y_1, \dots, Y_n are i.i.d., the mean of the sampling distribution of Y_1 is $E(Y_1) = \mu_Y$; thus Y_1 is an unbiased estimator of μ_Y . Its variance is $\text{var}(Y_1) = \sigma_Y^2$. From Section 2.5, the variance of \bar{Y} is σ_Y^2/n . Thus, for $n \geq 2$, the variance of \bar{Y} is less than the variance of Y_1 ; that is, \bar{Y} is a more efficient estimator than Y_1 , so, according to the criterion of efficiency, \bar{Y} should be used instead of Y_1 . The estimator Y_1 might strike you as an obviously poor estimator—why would you go to the trouble of collecting a sample of n observations only to throw away all but the first?—and the concept of efficiency provides a formal way to show that \bar{Y} is a more desirable estimator than Y_1 .

What about a less obviously poor estimator? Consider the weighted average in which the observations are alternately weighted by $\frac{1}{2}$ and $\frac{3}{2}$:

$$\tilde{Y} = \frac{1}{n} \left(\frac{1}{2}Y_1 + \frac{3}{2}Y_2 + \frac{1}{2}Y_3 + \frac{3}{2}Y_4 + \cdots + \frac{1}{2}Y_{n-1} + \frac{3}{2}Y_n \right), \quad (3.1)$$

where the number of observations n is assumed to be even for convenience. The mean of \tilde{Y} is μ_Y , and its variance is $\text{var}(\tilde{Y}) = 1.25 \sigma_Y^2/n$ (Exercise 3.11). Thus \tilde{Y} is

Efficiency of \bar{Y} : \bar{Y} Is BLUE

KEY CONCEPT

3.3

Let $\hat{\mu}_Y$ be an estimator of μ_Y that is a weighted average of Y_1, \dots, Y_n ; that is, $\hat{\mu}_Y = (1/n) \sum_{i=1}^n a_i Y_i$, where a_1, \dots, a_n are nonrandom constants. If $\hat{\mu}_Y$ is unbiased, then $\text{var}(\bar{Y}) < \text{var}(\hat{\mu}_Y)$ unless $\hat{\mu}_Y = \bar{Y}$. Thus \bar{Y} is the Best Linear Unbiased Estimator (BLUE); that is, \bar{Y} is the most efficient estimator of μ_Y among all unbiased estimators that are weighted averages of Y_1, \dots, Y_n .

unbiased, and because $\text{var}(\tilde{Y}) \rightarrow 0$ as $n \rightarrow \infty$, \tilde{Y} is consistent. However, \tilde{Y} has a larger variance than \bar{Y} . Thus \bar{Y} is more efficient than \tilde{Y} .

The estimators \bar{Y} , Y_1 , and \tilde{Y} have a common mathematical structure: They are weighted averages of Y_1, \dots, Y_n . The comparisons in the previous two paragraphs show that the weighted averages Y_1 and \tilde{Y} have larger variances than \bar{Y} . In fact, these conclusions reflect a more general result: \bar{Y} is the most efficient estimator of *all* unbiased estimators that are weighted averages of Y_1, \dots, Y_n . Said differently, \bar{Y} is the **Best Linear Unbiased Estimator (BLUE)**; that is, it is the most efficient (best) estimator among all estimators that are unbiased and are linear functions of Y_1, \dots, Y_n . This result is stated in Key Concept 3.3 and is proved in Chapter 5.

\bar{Y} is the least squares estimator of μ_Y . The sample average \bar{Y} provides the best fit to the data in the sense that the average squared differences between the observations and \bar{Y} are the smallest of all possible estimators.

Consider the problem of finding the estimator m that minimizes

$$\sum_{i=1}^n (Y_i - m)^2, \quad (3.2)$$

which is a measure of the total squared gap or distance between the estimator m and the sample points. Because m is an estimator of $E(Y)$, you can think of it as a prediction of the value of Y_i , so the gap $Y_i - m$ can be thought of as a prediction mistake. The sum of squared gaps in Expression (3.2) can be thought of as the sum of squared prediction mistakes.

The estimator m that minimizes the sum of squared gaps $Y_i - m$ in Expression (3.2) is called the **least squares estimator**. One can imagine using trial and error to solve the least squares problem: Try many values of m until you are satisfied that you have the value that makes Expression (3.2) as small as possible. Alternatively, as is done in Appendix 3.2, you can use algebra or calculus to show that choosing $m = \bar{Y}$ minimizes the sum of squared gaps in Expression (3.2), so that \bar{Y} is the least squares estimator of μ_Y .

Off the Mark!

In 2009, India's general elections, also referred to as the national elections, was the largest democratic election in the world until the Indian general elections 2014 held from April 7, 2014. Shortly before the general elections, pollsters predicted a close fight between the coalition parties—the United Progressive Alliance (UPA) and the National Democratic Alliance (NDA). Psephologists envisaged that while the UPA might have had the upper hand, the NDA could not be written off. They predicted that the UPA would get between 201 and 235 seats in the 14th Lok Sabha (the lower house of India's bicameral Parliament) and the NDA between 165 and 186 seats. The actual results were surprising: UPA got 262 seats, while NDA could only manage to get 157 seats.

What could be the possible reasons for opinion polls being wide off the mark? In countries that do not have a homogenous population, such as India, caste, religion, and geographies influence electoral outcomes greatly. Vulnerable sections of the population may be afraid to disclose their actual preference. Political polls have since become much more sophisticated and adjust for sampling bias, but they still can make mistakes. If opinion polls do not randomly select samples across various locations and sections of people, they may still not hit the mark.

Source: Atul Thakur, "Why Opinion Polls Are Often Wide off the Mark," *The Times of India*, April 13, 2014.

The Importance of Random Sampling

We have assumed that Y_1, \dots, Y_n are i.i.d. draws, such as those that would be obtained from simple random sampling. This assumption is important because non-random sampling can result in \bar{Y} being biased. Suppose that to estimate the monthly national unemployment rate, a statistical agency adopts a sampling scheme in which interviewers survey working-age adults sitting in city parks at 10 a.m. on the second Wednesday of the month. Because most employed people are at work at that hour (not sitting in the park!), the unemployed are overly represented in the sample, and an estimate of the unemployment rate based on this sampling plan would be biased. This bias arises because this sampling scheme overrepresents, or oversamples, the unemployed members of the population. This example is fictitious, but the "Off the Mark!" box gives a real-world example of biases introduced by sampling that is not entirely random.

It is important to design sample selection schemes in a way that minimizes bias. Appendix 3.1 includes a discussion of what the Bureau of Labor Statistics actually does when it conducts the U.S. Current Population Survey (CPS), the survey it uses to estimate the monthly U.S. unemployment rate.

3.2 Hypothesis Tests Concerning the Population Mean

Many hypotheses about the world around us can be phrased as yes/no questions. Do the mean hourly earnings of recent U.S. college graduates equal \$20 per hour? Are mean earnings the same for male and female college graduates? Both these questions embody specific hypotheses about the population distribution of earnings. The statistical challenge is to answer these questions based on a sample of evidence. This section describes **hypothesis tests** concerning the population mean (Does the population mean of hourly earnings equal \$20?). Hypothesis tests involving two populations (Are mean earnings the same for men and women?) are taken up in Section 3.4.

Null and Alternative Hypotheses

The starting point of statistical hypotheses testing is specifying the hypothesis to be tested, called the **null hypothesis**. Hypothesis testing entails using data to compare the null hypothesis to a second hypothesis, called the **alternative hypothesis**, that holds if the null does not.

The null hypothesis is that the population mean, $E(Y)$, takes on a specific value, denoted $\mu_{Y,0}$. The null hypothesis is denoted H_0 and thus is

$$H_0: E(Y) = \mu_{Y,0}. \quad (3.3)$$

For example, the conjecture that, on average in the population, college graduates earn \$20 per hour constitutes a null hypothesis about the population distribution of hourly earnings. Stated mathematically, if Y is the hourly earnings of a randomly selected recent college graduate, then the null hypothesis is that $E(Y) = 20$; that is, $\mu_{Y,0} = 20$ in Equation (3.3).

The alternative hypothesis specifies what is true if the null hypothesis is not. The most general alternative hypothesis is that $E(Y) \neq \mu_{Y,0}$, which is called a **two-sided alternative hypothesis** because it allows $E(Y)$ to be either less than or greater than $\mu_{Y,0}$. The two-sided alternative is written as

$$H_1: E(Y) \neq \mu_{Y,0} \text{ (two-sided alternative)}. \quad (3.4)$$

One-sided alternatives are also possible, and these are discussed later in this section.

The problem facing the statistician is to use the evidence in a randomly selected sample of data to decide whether to accept the null hypothesis H_0 or to reject it in favor of the alternative hypothesis H_1 . If the null hypothesis is “accepted,” this does not mean that the statistician declares it to be true; rather, it is accepted tentatively with the recognition that it might be rejected later based on additional evidence. For this reason, statistical hypothesis testing can be posed as either rejecting the null hypothesis or failing to do so.

The p -Value

In any given sample, the sample average \bar{Y} will rarely be exactly equal to the hypothesized value $\mu_{Y,0}$. Differences between \bar{Y} and $\mu_{Y,0}$ can arise because the true mean, in fact, does not equal $\mu_{Y,0}$ (the null hypothesis is false) or because the true mean equals $\mu_{Y,0}$ (the null hypothesis is true) but \bar{Y} differs from $\mu_{Y,0}$ because of random sampling. It is impossible to distinguish between these two possibilities with certainty. Although a sample of data cannot provide conclusive evidence about the null hypothesis, it is possible to do a probabilistic calculation that permits testing the null hypothesis in a way that accounts for sampling uncertainty. This calculation involves using the data to compute the p -value of the null hypothesis.

The **p -value**, also called the **significance probability**, is the probability of drawing a statistic at least as adverse to the null hypothesis as the one you actually computed in your sample, assuming the null hypothesis is correct. In the case at hand, the p -value is the probability of drawing \bar{Y} at least as far in the tails of its distribution under the null hypothesis as the sample average you actually computed.

For example, suppose that, in your sample of recent college graduates, the average wage is \$22.64. The p -value is the probability of observing a value of \bar{Y} at least as different from \$20 (the population mean under the null hypothesis) as the observed value of \$22.64 by pure random sampling variation, assuming that the null hypothesis is true. If this p -value is small (say, 0.1%), then it is very unlikely that this sample would have been drawn if the null hypothesis is true; thus it is reasonable to conclude that the null hypothesis is not true. By contrast, if this p -value is large (say, 40%), then it is quite likely that the observed sample average of \$22.64 could have arisen just by random sampling variation if the null hypothesis is true; accordingly, the evidence against the null hypothesis is weak in this probabilistic sense, and it is reasonable not to reject the null hypothesis.

To state the definition of the p -value mathematically, let \bar{Y}^{act} denote the value of the sample average actually computed in the data set at hand, and let \Pr_{H_0} denote the probability computed under the null hypothesis (that is, computed assuming that $E(Y) = \mu_{Y,0}$). The p -value is

$$p\text{-value} = \Pr_{H_0} [|\bar{Y} - \mu_{Y,0}| > |\bar{Y}^{act} - \mu_{Y,0}|]. \quad (3.5)$$

That is, the p -value is the area in the tails of the distribution of \bar{Y} under the null hypothesis beyond $\mu_{Y,0} \pm |\bar{Y}^{act} - \mu_{Y,0}|$. If the p -value is large, then the observed value \bar{Y}^{act} is consistent with the null hypothesis, but if the p -value is small, it is not.

To compute the p -value, it is necessary to know the sampling distribution of \bar{Y} under the null hypothesis. As discussed in Section 2.6, when the sample size is small, this distribution is complicated. However, according to the central limit theorem, when the sample size is large, the sampling distribution of \bar{Y} is well approximated by a normal distribution. Under the null hypothesis the mean of this normal distribution is $\mu_{Y,0}$, so under the null hypothesis \bar{Y} is distributed $N(\mu_{Y,0}, \sigma_{\bar{Y}}^2)$, where $\sigma_{\bar{Y}}^2 = \sigma_Y^2/n$.

This large-sample normal approximation makes it possible to compute the p -value without needing to know the population distribution of Y , as long as the sample size is large. The details of the calculation, however, depend on whether σ_Y^2 is known.

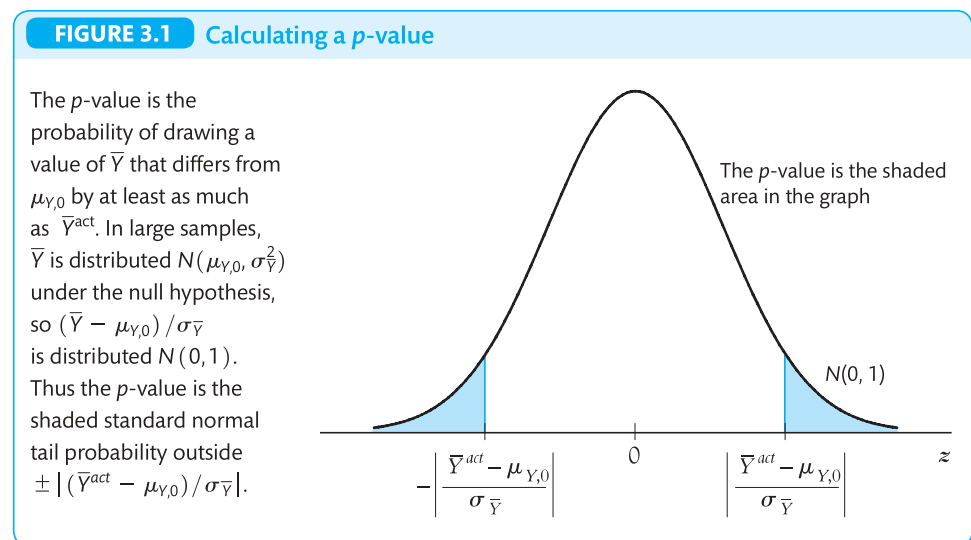
Calculating the p -Value When σ_Y Is Known

The calculation of the p -value when σ_Y is known is summarized in Figure 3.1. If the sample size is large, then under the null hypothesis the sampling distribution of \bar{Y} is $N(\mu_{Y,0}, \sigma_{\bar{Y}}^2)$, where $\sigma_{\bar{Y}}^2 = \sigma_Y^2/n$. Thus, under the null hypothesis, the standardized version of \bar{Y} , $(\bar{Y} - \mu_{Y,0})/\sigma_{\bar{Y}}$, has a standard normal distribution. The p -value is the probability of obtaining a value of \bar{Y} farther from $\mu_{Y,0}$ than \bar{Y}^{act} under the null hypothesis or, equivalently, it is the probability of obtaining $(\bar{Y} - \mu_{Y,0})/\sigma_{\bar{Y}}$ greater than $(\bar{Y}^{act} - \mu_{Y,0})/\sigma_{\bar{Y}}$ in absolute value. This probability is the shaded area shown in Figure 3.1. Written mathematically, the shaded tail probability in Figure 3.1 (that is, the p -value) is

$$p\text{-value} = \Pr_{H_0}\left(\left|\frac{\bar{Y} - \mu_{Y,0}}{\sigma_{\bar{Y}}}\right| > \left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_{\bar{Y}}}\right|\right) = 2\Phi\left(-\left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{\sigma_{\bar{Y}}}\right|\right), \quad (3.6)$$

where Φ is the standard normal cumulative distribution function. That is, the p -value is the area in the tails of a standard normal distribution outside $\pm |\bar{Y}^{act} - \mu_{Y,0}|/\sigma_{\bar{Y}}$.

The formula for the p -value in Equation (3.6) depends on the variance of the population distribution, σ_Y^2 . In practice, this variance is typically unknown. [An exception is when Y_i is binary, so that its distribution is Bernoulli, in which case the variance is determined by the null hypothesis; see Equation (2.7) and Exercise 3.2.] Because in general σ_Y^2 must be estimated before the p -value can be computed, we now turn to the problem of estimating σ_Y^2 .



The Sample Variance, Sample Standard Deviation, and Standard Error

The sample variance, s_Y^2 , is an estimator of the population variance, σ_Y^2 ; the sample standard deviation, s_Y , is an estimator of the population standard deviation, σ_Y ; and the standard error of the sample average, \bar{Y} , is an estimator of the standard deviation of the sampling distribution of \bar{Y} .

The sample variance and standard deviation. The **sample variance**, s_Y^2 , is

$$s_Y^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2. \quad (3.7)$$

The **sample standard deviation**, s_Y , is the square root of the sample variance.

The formula for the sample variance is much like the formula for the population variance. The population variance, $E(Y - \mu_Y)^2$, is the average value of $(Y - \mu_Y)^2$ in the population distribution. Similarly, the sample variance is the sample average of $(Y_i - \mu_Y)^2$, $i = 1, \dots, n$, with two modifications: First, μ_Y is replaced by \bar{Y} , and second, the average uses the divisor $n - 1$ instead of n .

The reason for the first modification—replacing μ_Y by \bar{Y} —is that μ_Y is unknown and thus must be estimated; the natural estimator of μ_Y is \bar{Y} . The reason for the second modification—dividing by $n - 1$ instead of by n —is that estimating μ_Y by \bar{Y} introduces a small downward bias in $(Y_i - \bar{Y})^2$. Specifically, as is shown in Exercise 3.18, $E[(Y_i - \bar{Y})^2] = [(n-1)/n]\sigma_Y^2$. Thus $E\sum_{i=1}^n (Y_i - \bar{Y})^2 = nE[(Y_i - \bar{Y})^2] = (n-1)\sigma_Y^2$. Dividing by $n - 1$ in Equation (3.7) instead of n corrects for this small downward bias, and as a result s_Y^2 is unbiased.

Dividing by $n - 1$ in Equation (3.7) instead of n is called a **degrees of freedom** correction: Estimating the mean uses up some of the information—that is, uses up 1 “degree of freedom”—in the data, so that only $n - 1$ degrees of freedom remain.

Consistency of the sample variance. The sample variance is a consistent estimator of the population variance:

$$s_Y^2 \xrightarrow{p} \sigma_Y^2. \quad (3.8)$$

In other words, the sample variance is close to the population variance with high probability when n is large.

The result in Equation (3.9) is proven in Appendix 3.3 under the assumptions that Y_1, \dots, Y_n are i.i.d. and Y_i has a finite fourth moment; that is, $E(Y_i^4) < \infty$. Intuitively, the reason that s_Y^2 is consistent is that it is a sample average, so s_Y^2 obeys the law of large numbers. For s_Y^2 to obey the law of large numbers in Key Concept 2.6, $(Y_i - \mu_Y)^2$ must have finite variance, which in turn means that $E(Y_i^4)$ must be finite; in other words, Y_i must have a finite fourth moment.

The Standard Error of \bar{Y}

KEY CONCEPT

3.4

The standard error of \bar{Y} is an estimator of the standard deviation of \bar{Y} . The standard error of \bar{Y} is denoted $SE(\bar{Y})$ or $\hat{\sigma}_{\bar{Y}}$. When Y_1, \dots, Y_n are i.i.d.,

$$SE(\bar{Y}) = \hat{\sigma}_{\bar{Y}} = s_Y / \sqrt{n}. \quad (3.9)$$

The standard error of \bar{Y} . Because the standard deviation of the sampling distribution of \bar{Y} is $\sigma_{\bar{Y}} = \sigma_Y / \sqrt{n}$, Equation (3.9) justifies using s_Y / \sqrt{n} as an estimator of $\sigma_{\bar{Y}}$. The estimator of $\sigma_{\bar{Y}}, s_Y / \sqrt{n}$, is called the **standard error of \bar{Y}** and is denoted $SE(\bar{Y})$ or $\hat{\sigma}_{\bar{Y}}$. The standard error of \bar{Y} is summarized as in Key Concept 3.4.

When Y_1, \dots, Y_n are i.i.d. draws from a Bernoulli distribution with success probability p , the formula for the variance of \bar{Y} simplifies to $p(1-p)/n$ (see Exercise 3.2). The formula for the standard error also takes on a simple form that depends only on \bar{Y} and n : $SE(\bar{Y}) = \sqrt{\bar{Y}(1-\bar{Y})/n}$.

Calculating the p -Value When σ_Y Is Unknown

Because s_Y^2 is a consistent estimator of σ_Y^2 , the p -value can be computed by replacing $\sigma_{\bar{Y}}$ in Equation (3.6) by the standard error, $SE(\bar{Y}) = \hat{\sigma}_{\bar{Y}}$. That is, when σ_Y is unknown and Y_1, \dots, Y_n are i.i.d., the p -value is calculated using the formula

$$p\text{-value} = 2\Phi\left(-\left|\frac{\bar{Y}^{act} - \mu_{Y,0}}{SE(\bar{Y})}\right|\right). \quad (3.10)$$

The t -Statistic

The standardized sample average $(\bar{Y} - \mu_{Y,0}) / SE(\bar{Y})$ plays a central role in testing statistical hypotheses and has a special name, the **t -statistic** or **t -ratio**:

$$t = \frac{\bar{Y} - \mu_{Y,0}}{SE(\bar{Y})}. \quad (3.11)$$

In general, a **test statistic** is a statistic used to perform a hypothesis test. The t -statistic is an important example of a test statistic.

Large-sample distribution of the t -statistic. When n is large, s_Y^2 is close to σ_Y^2 with high probability. Thus the distribution of the t -statistic is approximately the same as the distribution of $(\bar{Y} - \mu_{Y,0}) / \sigma_{\bar{Y}}$, which in turn is well approximated by the standard normal distribution when n is large because of the central limit theorem (Key Concept 2.7). Accordingly, under the null hypothesis,

$$t \text{ is approximately distributed } N(0, 1) \text{ for large } n. \quad (3.12)$$

The formula for the p -value in Equation (3.10) can be rewritten in terms of the t -statistic. Let t^{act} denote the value of the t -statistic actually computed:

$$t^{act} = \frac{\bar{Y}^{act} - \mu_{Y,0}}{SE(\bar{Y})}. \quad (3.13)$$

Accordingly, when n is large, the p -value can be calculated using

$$p\text{-value} = 2\Phi(-|t^{act}|). \quad (3.14)$$

As a hypothetical example, suppose that a sample of $n = 200$ recent college graduates is used to test the null hypothesis that the mean wage, $E(Y)$, is \$20 per hour. The sample average wage is $\bar{Y}^{act} = \$22.64$, and the sample standard deviation is $s_Y = \$18.14$. Then the standard error of \bar{Y} is $s_Y/\sqrt{n} = 18.14/\sqrt{200} = 1.28$. The value of the t -statistic is $t^{act} = (22.64 - 20)/1.28 = 2.06$. From Appendix Table 1, the p -value is $2\Phi(-2.06) = 0.039$, or 3.9%. That is, assuming the null hypothesis to be true, the probability of obtaining a sample average at least as different from the null as the one actually computed is 3.9%.

Hypothesis Testing with a Prespecified Significance Level

When you undertake a statistical hypothesis test, you can make two types of mistakes: You can incorrectly reject the null hypothesis when it is true, or you can fail to reject the null hypothesis when it is false. Hypothesis tests can be performed without computing the p -value if you are willing to specify in advance the probability you are willing to tolerate of making the first kind of mistake—that is, of incorrectly rejecting the null hypothesis when it is true. If you choose a prespecified probability of rejecting the null hypothesis when it is true (for example, 5%), then you will reject the null hypothesis if and only if the p -value is less than 0.05. This approach gives preferential treatment to the null hypothesis, but in many practical situations, this preferential treatment is appropriate.

Hypothesis tests using a fixed significance level. Suppose it has been decided that the hypothesis will be rejected if the p -value is less than 5%. Because the area under the tails of the standard normal distribution outside ± 1.96 is 5%, this gives a simple rule:

$$\text{Reject } H_0 \text{ if } |t^{act}| > 1.96. \quad (3.15)$$

That is, reject if the absolute value of the t -statistic computed from the sample is greater than 1.96. If n is large enough, then under the null hypothesis the t -statistic has a $N(0, 1)$ distribution. Thus the probability of erroneously rejecting the null hypothesis (rejecting the null hypothesis when it is, in fact, true) is 5%.

This framework for testing statistical hypotheses has some specialized terminology, summarized in Key Concept 3.5. The significance level of the test in

The Terminology of Hypothesis Testing

KEY CONCEPT

3.5

A statistical hypothesis test can make two types of mistakes: a **type I error**, in which the null hypothesis is rejected when in fact it is true; and a **type II error**, in which the null hypothesis is not rejected when in fact it is false. The prespecified rejection probability of a statistical hypothesis test when the null hypothesis is true—that is, the prespecified probability of a type I error—is the **significance level** of the test. The **critical value** of the test statistic is the value of the statistic for which the test just rejects the null hypothesis at the given significance level. The set of values of the test statistic for which the test rejects the null hypothesis is the **rejection region**, and the set of values of the test statistic for which it does not reject the null hypothesis is the **acceptance region**. The probability that the test actually incorrectly rejects the null hypothesis when it is true is the **size of the test**, and the probability that the test correctly rejects the null hypothesis when the alternative is true is the **power of the test**.

The **p -value** is the probability of obtaining a test statistic, by random sampling variation, at least as adverse to the null hypothesis value as is the statistic actually observed, assuming that the null hypothesis is correct. Equivalently, the p -value is the smallest significance level at which you can reject the null hypothesis.

Equation (3.15) is 5%, the critical value of this two-sided test is 1.96, and the rejection region is the values of the t -statistic outside ± 1.96 . If the test rejects at the 5% significance level, the population mean μ_Y is said to be statistically significantly different from $\mu_{Y,0}$ at the 5% significance level.

Testing hypotheses using a prespecified significance level does not require computing p -values. In the previous example of testing the hypothesis that the mean earnings of recent college graduates is \$20 per hour, the t -statistic was 2.06. This value exceeds 1.96, so the hypothesis is rejected at the 5% level. Although performing the test with a 5% significance level is easy, reporting only whether the null hypothesis is rejected at a prespecified significance level conveys less information than reporting the p -value.

What significance level should you use in practice? This is a question of active debate. Historically, statisticians and econometricians have used a 5% significance level. If you were to test many statistical hypotheses at the 5% level, you would incorrectly reject the null, on average, once in 20 cases. Whether this is a small number depends on how you look at it. If only a small fraction of all null hypotheses tested are, in fact, false, then among those tests that reject, the probability of the null actually being false can be small (Exercise 3.22). This probability—the fraction of incorrect rejections among all rejections—is called the false positive rate. The false positive rate can have great practical importance. For example, for newly reported statistically

KEY CONCEPT

3.6

Testing the Hypothesis $E(Y) = \mu_{Y,0}$ Against the Alternative $E(Y) \neq \mu_{Y,0}$

1. Compute the standard error of \bar{Y} , $SE(\bar{Y})$ [Equation (3.8)].
2. Compute the t -statistic [Equation (3.13)].
3. Compute the p -value [Equation (3.14)]. Reject the hypothesis at the 5% significance level if the p -value is less than 0.05 (equivalently, if $|t^{act}| > 1.96$).

significant findings of effective medical treatments, it is the fraction for which the treatment is in fact ineffective. Concern that the false positive rate can be high when the 5% significance level is used has led some statisticians to recommend using instead a 0.5% significance level when reporting new results (Benjamin et al., 2017). Similar concerns can apply in a legal setting, where justice might aim to keep the fraction of false convictions low. Using a 0.5% significance level leads to two-sided rejection when the t -statistic exceeds 2.81 in absolute value. In such cases, a p -value between 0.05 and 0.005 can be viewed as suggestive, but not conclusive, evidence against the null that merits further investigation.

The choice of significance level requires judgment and depends on the application. In some economic applications, a false positive might be less of a problem than in a medical context, where the false positive could lead to patients receiving ineffective treatments. In such cases, a 5% significance level could be appropriate.

Whatever the significance level, it is important to keep in mind that p -values are designed for tests of a null hypothesis, so they, like t -statistics, are useful only when the null hypothesis itself is of interest. This section uses the example of earnings. Even though many interns are unpaid, nobody thinks that, on average, workers earn nothing at all, so the null hypothesis that earnings are zero is economically uninteresting and not worth testing. In contrast, the null hypothesis that the mean earnings of men and of women are the same is interesting and of societal importance, and that null hypothesis is examined in Section 3.4.

Key Concept 3.6 summarizes hypothesis tests for the population mean against the two-sided alternative.

One-Sided Alternatives

In some circumstances, the alternative hypothesis might be that the mean exceeds $\mu_{Y,0}$. For example, one hopes that education helps in the labor market, so the relevant alternative to the null hypothesis that earnings are the same for college graduates and non-college graduates is not just that their earnings differ, but rather that graduates earn more than nongraduates. This is called a **one-sided alternative hypothesis** and can be written

$$H_1: E(Y) > \mu_{Y,0} \text{ (one-sided alternative)}. \quad (3.16)$$

The general approach to computing p -values and to hypothesis testing is the same for one-sided alternatives as it is for two-sided alternatives, with the modification that only large positive values of the t -statistic reject the null hypothesis rather than values that are large in absolute value. Specifically, to test the one-sided hypothesis in Equation (3.16), construct the t -statistic in Equation (3.13). The p -value is the area under the standard normal distribution to the right of the calculated t -statistic. That is, the p -value, based on the $N(0, 1)$ approximation to the distribution of the t -statistic, is

$$p\text{-value} = \Pr_{H_0}(Z > t^{act}) = 1 - \Phi(t^{act}). \quad (3.17)$$

The $N(0, 1)$ critical value for a one-sided test with a 5% significance level is 1.64. The rejection region for this test is all values of the t -statistic exceeding 1.64.

The one-sided hypothesis in Equation (3.16) concerns values of μ_Y exceeding $\mu_{Y,0}$. If instead the alternative hypothesis is that $E(Y) < \mu_{Y,0}$, then the discussion of the previous paragraph applies except that the signs are switched; for example, the 5% rejection region consists of values of the t -statistic less than -1.64 .

3.3 Confidence Intervals for the Population Mean

Because of random sampling error, it is impossible to learn the exact value of the population mean of Y using only the information in a sample. However, it is possible to use data from a random sample to construct a set of values that contains the true population mean μ_Y with a certain prespecified probability. Such a set is called a **confidence set**, and the prespecified probability that μ_Y is contained in this set is called the **confidence level**. The confidence set for μ_Y turns out to be all the possible values of the mean between a lower and an upper limit, so that the confidence set is an interval, called a **confidence interval**.

Here is one way to construct a 95% confidence set for the population mean. Begin by picking some arbitrary value for the mean; call it $\mu_{Y,0}$. Test the null hypothesis that $\mu_Y = \mu_{Y,0}$ against the alternative that $\mu_Y \neq \mu_{Y,0}$ by computing the t -statistic; if its absolute value is less than 1.96, this hypothesized value $\mu_{Y,0}$ is not rejected at the 5% level, so write down this nonrejected value $\mu_{Y,0}$. Now pick another arbitrary value of $\mu_{Y,0}$ and test it; if you cannot reject it, write down this value on your list. Do this again and again; indeed, do so for all possible values of the population mean. Continuing this process yields the set of all values of the population mean that cannot be rejected at the 5% level by a two-sided hypothesis test.

This list is useful because it summarizes the set of hypotheses you can and cannot reject (at the 5% level) based on your data: If someone walks up to you with a specific number in mind, you can tell him whether his hypothesis is rejected or not simply by looking up his number on your handy list. A bit of clever reasoning shows that this set of values has a remarkable property: The probability that it contains the true value of the population mean is 95%.

KEY CONCEPT

Confidence Intervals for the Population Mean

3.7

A 95% two-sided confidence interval for μ_Y is an interval constructed so that it contains the true value of μ_Y in 95% of all possible random samples. When the sample size n is large, 90%, 95%, and 99% confidence intervals for μ_Y are:

$$90\% \text{ confidence interval for } \mu_Y = \{\bar{Y} \pm 1.64SE(\bar{Y})\},$$

$$95\% \text{ confidence interval for } \mu_Y = \{\bar{Y} \pm 1.96SE(\bar{Y})\}, \text{ and}$$

$$99\% \text{ confidence interval for } \mu_Y = \{\bar{Y} \pm 2.58SE(\bar{Y})\}.$$

The clever reasoning goes like this: Suppose the true value of μ_Y is 21.5 (although we do not know this). Then \bar{Y} has a normal distribution centered on 21.5, and the t -statistic testing the null hypothesis $\mu_Y = 21.5$ has a $N(0, 1)$ distribution. Thus, if n is large, the probability of rejecting the null hypothesis $\mu_Y = 21.5$ at the 5% level is 5%. But because you tested all possible values of the population mean in constructing your set, in particular you tested the true value, $\mu_Y = 21.5$. In 95% of all samples, you will correctly accept 21.5; this means that in 95% of all samples, your list will contain the true value of μ_Y . Thus the values on your list constitute a 95% confidence set for μ_Y .

This method of constructing a confidence set is impractical, for it requires you to test all possible values of μ_Y as null hypotheses. Fortunately, there is a much easier approach. According to the formula for the t -statistic in Equation (3.13), a trial value of $\mu_{Y,0}$ is rejected at the 5% level if it is more than 1.96 standard errors away from \bar{Y} . Thus the set of values of μ_Y that are not rejected at the 5% level consists of those values within $\pm 1.96SE(\bar{Y})$ of \bar{Y} ; that is, a 95% confidence interval for μ_Y is $\bar{Y} - 1.96SE(\bar{Y}) \leq \mu_Y \leq \bar{Y} + 1.96SE(\bar{Y})$. Key Concept 3.7 summarizes this approach.

As an example, consider the problem of constructing a 95% confidence interval for the mean hourly earnings of recent college graduates using a hypothetical random sample of 200 recent college graduates where $\bar{Y} = \$22.64$ and $SE(\bar{Y}) = 1.28$. The 95% confidence interval for mean hourly earnings is $22.64 \pm 1.96 \times 1.28 = 22.64 \pm 2.51 = (\$20.13, \$25.15)$.

This discussion so far has focused on two-sided confidence intervals. One could instead construct a one-sided confidence interval as the set of values of μ_Y that cannot be rejected by a one-sided hypothesis test. Although one-sided confidence intervals have applications in some branches of statistics, they are uncommon in applied econometric analysis.

Coverage probabilities. The **coverage probability** of a confidence interval for the population mean is the probability, computed over all possible random samples, that it contains the true population mean.

3.4 Comparing Means from Different Populations

Do recent male and female college graduates earn the same amount on average? Answering this question involves comparing the means of two different population distributions. This section summarizes how to test hypotheses and how to construct confidence intervals for the difference in the means from two different populations.

Hypothesis Tests for the Difference Between Two Means

To illustrate a **test for the difference between two means**, let μ_w be the mean hourly earnings in the population of women recently graduated from college, and let μ_m be the population mean for recently graduated men. Consider the null hypothesis that mean earnings for these two populations differ by a certain amount, say, d_0 . Then the null hypothesis and the two-sided alternative hypothesis are

$$H_0: \mu_m - \mu_w = d_0 \text{ vs. } H_1: \mu_m - \mu_w \neq d_0. \quad (3.18)$$

The null hypothesis that men and women in these populations have the same mean earnings corresponds to H_0 in Equation (3.18) with $d_0 = 0$.

Because these population means are unknown, they must be estimated from samples of men and women. Suppose we have samples of n_m men and n_w women drawn at random from their populations. Let the sample average annual earnings be \bar{Y}_m for men and \bar{Y}_w for women. Then an estimator of $\mu_m - \mu_w$ is $\bar{Y}_m - \bar{Y}_w$.

To test the null hypothesis that $\mu_m - \mu_w = d_0$ using $\bar{Y}_m - \bar{Y}_w$, we need to know the sampling distribution of $\bar{Y}_m - \bar{Y}_w$. Recall that \bar{Y}_m is, according to the central limit theorem, approximately distributed $N(\mu_m, \sigma_m^2/n_m)$, where σ_m^2 is the population variance of earnings for men. Similarly, \bar{Y}_w is approximately distributed $N(\mu_w, \sigma_w^2/n_w)$, where σ_w^2 is the population variance of earnings for women. Also, recall from Section 2.4 that a weighted average of two normal random variables is itself normally distributed. Because \bar{Y}_m and \bar{Y}_w are constructed from different randomly selected samples, they are independent random variables. Thus $\bar{Y}_m - \bar{Y}_w$ is distributed $N[\mu_m - \mu_w, (\sigma_m^2/n_m) + (\sigma_w^2/n_w)]$.

If σ_m^2 and σ_w^2 are known, then this approximate normal distribution can be used to compute p -values for the test of the null hypothesis that $\mu_m - \mu_w = d_0$. In practice, however, these population variances are typically unknown, so they must be estimated. As before, they can be estimated using the sample variances, s_m^2 and s_w^2 , where s_m^2 is defined as in Equation (3.7), except that the statistic is computed only for

the men in the sample, and s_w^2 is defined similarly for the women. Thus the standard error of $\bar{Y}_m - \bar{Y}_w$ is

$$SE(\bar{Y}_m - \bar{Y}_w) = \sqrt{\frac{s_m^2}{n_m} + \frac{s_w^2}{n_w}}. \quad (3.19)$$

For a simplified version of Equation (3.19) when Y is a Bernoulli random variable, see Exercise 3.15.

The t -statistic for testing the null hypothesis is constructed analogously to the t -statistic for testing a hypothesis about a single population mean, by subtracting the null hypothesized value of $\mu_m - \mu_w$ from the estimator $\bar{Y}_m - \bar{Y}_w$ and dividing the result by the standard error of $\bar{Y}_m - \bar{Y}_w$:

$$t = \frac{(\bar{Y}_m - \bar{Y}_w) - d_0}{SE(\bar{Y}_m - \bar{Y}_w)} \quad (t\text{-statistic for comparing two means}). \quad (3.20)$$

If both n_m and n_w are large, then this t -statistic has a standard normal distribution when the null hypothesis is true.

Because the t -statistic in Equation (3.20) has a standard normal distribution under the null hypothesis when n_m and n_w are large, the p -value of the two-sided test is computed exactly as it was in the case of a single population. That is, the p -value is computed using Equation (3.14).

To conduct a test with a prespecified significance level, simply calculate the t -statistic in Equation (3.20), and compare it to the appropriate critical value. For example, the null hypothesis is rejected at the 5% significance level if the absolute value of the t -statistic exceeds 1.96.

If the alternative is one-sided rather than two-sided (that is, if the alternative is that $\mu_m - \mu_w > d_0$), then the test is modified as outlined in Section 3.2. The p -value is computed using Equation (3.17), and a test with a 5% significance level rejects when $t > 1.64$.

Confidence Intervals for the Difference Between Two Population Means

The method for constructing confidence intervals summarized in Section 3.3 extends to constructing a confidence interval for the difference between the means, $d = \mu_m - \mu_w$. Because the hypothesized value d_0 is rejected at the 5% level if $|t| > 1.96$, d_0 will be in the confidence set if $|t| \leq 1.96$. But $|t| \leq 1.96$ means that the estimated difference, $\bar{Y}_m - \bar{Y}_w$, is less than 1.96 standard errors away from d_0 . Thus the 95% two-sided confidence interval for d consists of those values of d within ± 1.96 standard errors of $\bar{Y}_m - \bar{Y}_w$:

$$\begin{aligned} &95\% \text{ confidence interval for } d = \mu_m - \mu_w \text{ is} \\ &(\bar{Y}_m - \bar{Y}_w) \pm 1.96SE(\bar{Y}_m - \bar{Y}_w). \end{aligned} \quad (3.21)$$

With these formulas in hand, the box “Social Class or Education? Childhood Circumstances and Adult Earnings Revisited” contains an empirical investigation of differences in earnings of different households in the United Kingdom.

3.5 Differences-of-Means Estimation of Causal Effects Using Experimental Data

Recall from Section 1.2 that a randomized controlled experiment randomly selects subjects (individuals or, more generally, entities) from a population of interest, then randomly assigns them either to a treatment group, which receives the experimental treatment, or to a control group, which does not receive the treatment. The difference between the sample means of the treatment and control groups is an estimator of the causal effect of the treatment.

The Causal Effect as a Difference of Conditional Expectations

The causal effect of a treatment is the expected effect on the outcome of interest of the treatment as measured in an ideal randomized controlled experiment. This effect can be expressed as the difference of two conditional expectations. Specifically, the **causal effect** on Y of treatment level x is the difference in the conditional expectations, $E(Y|X = x) - E(Y|X = 0)$, where $E(Y|X = x)$ is the expected value of Y for the treatment group (which receives treatment level $X = x$) in an ideal randomized controlled experiment and $E(Y|X = 0)$ is the expected value of Y for the control group (which receives treatment level $X = 0$). In the context of experiments, the causal effect is also called the **treatment effect**. If there are only two treatment levels (that is, if the treatment is binary), then we can let $X = 0$ denote the control group and $X = 1$ denote the treatment group. If the treatment is binary, then the causal effect (that is, the treatment effect) is $E(Y|X = 1) - E(Y|X = 0)$ in an ideal randomized controlled experiment.

Estimation of the Causal Effect Using Differences of Means

If the treatment in a randomized controlled experiment is binary, then the causal effect can be estimated by the difference in the sample average outcomes between the treatment and control groups. The hypothesis that the treatment is ineffective is equivalent to the hypothesis that the two means are the same, which can be tested using the t -statistic for comparing two means, given in Equation (3.20). A 95% confidence interval for the difference in the means of the two groups is a 95% confidence interval for the causal effect, so a 95% confidence interval for the causal effect can be constructed using Equation (3.21).

A well-designed, well-run experiment can provide a compelling estimate of a causal effect. For this reason, randomized controlled experiments are commonly conducted in some fields, such as medicine. In economics, however, experiments tend to be expensive, difficult to administer, and, in some cases, ethically questionable, so they are used less often. For this reason, econometricians sometimes study “natural

Social Class or Education? Childhood Circumstances and Adult Earnings Revisited

The box in Chapter 2 “The Distribution of Adulthood Earnings in the United Kingdom by Childhood Socioeconomic Circumstances” suggests that when an individual’s father has a “routine” occupation, the individual, as an adult, goes on to live in a household with lower average income.

Are there any other factors that affect it? Yes, it is possible that there are relevant intermediate factors like education. It is generally hypothesized and observed that more education is associated with greater income, which will allow individuals to increase their contribution to household income.

Table 3.1 breaks down the differences in mean household income for individuals according to their father’s NS-SEC occupation type, and considers these differences for selected highest level of educational qualification. These categories include those with no qualifications, those whose highest qualification level is GCSE (exams generally taken at age 16), those whose highest educational qualification is A-Level (exams generally taken at age 18), and those with an undergraduate degree or higher. For simplicity, only

individuals whose father’s NS-SEC occupational category was either the highest (“higher”) or the lowest (“routine”) are included in this analysis.

The data shows that, as expected, within both groups according to the NS-SEC of a father’s occupation, those with higher qualifications are part of households with higher total income. The income gap between those with qualifications of at least one degree and those with no qualifications stands at £1467.38 where the father’s NS-SEC category is higher, and at a comparable £1527.98 where the father’s NS-SEC category is routine.

It is interesting to note the differences between mean income by the father’s occupational categorization ($Y_h - Y_r$) for each of the educational groupings. For instance, individuals with no qualifications whose father’s NS-SEC job categorization was higher are part of households with a mean income of £2223.13 while for the classification routine this value stood at £1842.98. This implies a difference in means of £380.15, with a standard error of $\sqrt{2115.12^2/1129 + 1487.29^2/6383} = £65.64$ with

TABLE 3.1 Differences in Household Income According to Childhood Socioeconomic Circumstances, Grouped by Level of Highest Qualification

Qualification	Father’s NS-SEC = Higher			Father’s NS-SEC = Routine			Difference, Higher vs. Routine		
	Y_h	s_h	n_h	Y_r	s_r	n_r	$Y_h - Y_r$	$SE(Y_h - Y_r)$	95% Confidence Interval for d
None	£2,223.13	£2,115.12	1129	£1,842.98	£1,487.29	6383	£380.15	£65.64	£251.38 £508.93
GCSE/O-Level	£2,837.18	£1,819.73	1962	£2,596.93	£1,738.47	4042	£240.25	£49.35	£143.49 £337.00
A-Level	£3,045.99	£2,451.81	1216	£2,745.70	£1,912.50	1169	£300.30	£89.85	£124.11 £476.49
Undergraduate degree or more	£3,690.51	£2,743.55	4359	£3,370.96	£2,443.58	2505	£319.55	£64.11	£193.86 £445.23
All categories	£3,215.71	£2,497.73	8666	£2405.45	£1,886.86	14099	£810.25	£31.18	£749.13 £871.38

Source: Understanding Society.

a 95% confidence interval of (£251.38, £508.93). It is worth noting the difference in income, pooling these educational categories together, between those whose father's NS-SEC categorization is "higher" and those where this categorization is lower is £810.25. The results in the table suggest a difference in composition by educational attainment of these groupings according to the father's NS-SEC category. When broken down in this way, however, the estimated difference for every qualification level is substantially lower than £810.25. All of these estimated differences are significantly different from zero.

This empirical analysis suggests that levels of education do play some part in explaining the

difference in household income according to the socioeconomic status of the father. However, does this analysis tell us the full story? Are individuals with higher levels of education likely to be in households with more than one earner? Does the difference in household income arise from an individual's own contribution to household income or, if the individual is cohabiting, also from her or his partner's contribution to household income? Is this relationship affected by changing patterns of educational attainment that are correlated with age? We will examine questions such as these further once we have introduced the basics of multivariate regression in later chapters.

experiments," also called quasi-experiments, in which some event unrelated to the treatment or subject characteristics has the effect of assigning different treatments to different subjects *as if* they had been part of a randomized controlled experiment. The box "A Way to Increase Voter Turnout" provides an example of such a quasi-experiment that yielded some surprising conclusions.

3.6 Using the t -Statistic When the Sample Size Is Small

In Sections 3.2 through 3.5, the t -statistic is used in conjunction with critical values from the standard normal distribution for hypothesis testing and for the construction of confidence intervals. The use of the standard normal distribution is justified by the central limit theorem, which applies when the sample size is large. When the sample size is small, the standard normal distribution can provide a poor approximation to the distribution of the t -statistic. If, however, the population distribution is itself normally distributed, then the exact distribution (that is, the finite-sample distribution; see Section 2.6) of the t -statistic testing the mean of a single population is the Student t distribution with $n - 1$ degrees of freedom, and critical values can be taken from the Student t distribution.

A Way to Increase Voter Turnout

Apathy among citizens toward political participation, especially in voting, has been noted in the United Kingdom and other democratic countries. This kind of behavior is generally seen in economies where people have greater mobility, maintain an intensive work culture, and work for private corporate entities. Apart from these, there could be other dominant factors that have had a negative impact on the citizens' willingness to participate in elections—politicians failing to keep their promises, inappropriately using public funds.

In 2005, during the campaign period before the general election, a study was conducted in a Manchester constituency in the United Kingdom. The constituency's voter turnout rate in the 2001 general election had been 48.6%, while the national average had been 59.4%. Thus, voter participation in this constituency was far below the national average. For the experiment, three groups (two treatment groups and one control group) were randomly selected out of the registered voters from whom landline numbers could be obtained. One of the treatment groups was exposed to strong canvassing in the form of telephone calls, and the other treatment group was exposed to strong canvassing in the form of door-to-door visits. No contacts were made with the control group. The callers and the door-to-door canvassers were given instructions to ask respondents three questions, namely, whether the respondents thought voting is important, whether the respondents intended to vote, and whether they would vote by post. The conversations were informal and the main objective of this exercise was to persuade citizens to vote, by focusing on the importance

of voting. The callers and canvassers were also advised to respond to any concerns of the voters regarding the voting process.

The researchers got interesting results from the elections. The participation rate was 55.1% in the group, which was exposed to canvassing. The participation rate for the treatment group, which was treated with telephone calls, was 55%. Both these rates had a difference with the control group, which was not exposed to any experiment. Further calculations using suitable methodologies gave estimates of the effects of canvassing and telephone calls. 6.7% and 7.3% were the estimates of the two. The overall experiment was a success as the two interventions done on the two treatments groups by a non-partisan source had impacts that were statistically significant.

This exercise illustrated that citizens can be nudged to participate in elections by creating awareness through personal contacts. In yet another democracy, India, the 2014 general election saw a record voter turnout. A top Election Commission official has said that the Election Commission's efforts to increase voters' awareness and their registration has helped the process.

Sources: 1. Alice Moseley, Corinne Wales, Gerry Stoker, Graham Smith, Liz Richardson, Peter John, and Sarah Cotterill, "Nudge, Nudge, Think, Think Experimenting with Ways to Change Civic Behaviour," *Bloomsbury Academic*, March 2013. 2. "Lok Sabha Polls 2014: Country Records Highest Voter Turnout since Independence," *The Economic Times*, May 13, 2014.

The t -Statistic and the Student t Distribution

The t -statistic testing the mean. Consider the t -statistic used to test the hypothesis that the mean of Y is $\mu_{Y,0}$, using data Y_1, \dots, Y_n . The formula for this statistic is given by Equation (3.10), where the standard error of \bar{Y} is given by Equation (3.8). Substitution of the latter expression into the former yields the formula for the t -statistic:

$$t = \frac{\bar{Y} - \mu_{Y,0}}{\sqrt{s_Y^2/n}}, \quad (3.22)$$

where s_Y^2 is given in Equation (3.7).

As discussed in Section 3.2, under general conditions the t -statistic has a standard normal distribution if the sample size is large and the null hypothesis is true [see Equation (3.12)]. Although the standard normal approximation to the t -statistic is reliable for a wide range of distributions of Y if n is large, it can be unreliable if n is small. The exact distribution of the t -statistic depends on the distribution of Y , and it can be very complicated. There is, however, one special case in which the exact distribution of the t -statistic is relatively simple: If Y_1, \dots, Y_n are i.i.d. draws from a normal distribution, then the t -statistic in Equation (3.22) has a Student t distribution with $n - 1$ degrees of freedom. (The mathematics behind this result is provided in Sections 18.4 and 19.4.)

If the population distribution is normally distributed, then critical values from the Student t distribution can be used to perform hypothesis tests and to construct confidence intervals. As an example, consider a hypothetical problem in which $t^{act} = 2.15$ and $n = 8$, so that the degrees of freedom is $n - 1 = 7$. From Appendix Table 2, the 5% two-sided critical value for the t_7 distribution is 2.36. Because the t -statistic is smaller in absolute value than the critical value ($2.15 < 2.36$), the null hypothesis would not be rejected at the 5% significance level against the two-sided alternative. The 95% confidence interval for μ_Y , constructed using the t_7 distribution, would be $\bar{Y} \pm 2.36SE(\bar{Y})$. This confidence interval is wider than the confidence interval constructed using the standard normal critical value of 1.96.

The t -statistic testing differences of means. The t -statistic testing the difference of two means, given in Equation (3.20), does not have a Student t distribution, even if the population distribution of Y is normal. (The Student t distribution does not apply here because the variance estimator used to compute the standard error in Equation (3.19) does not produce a denominator in the t -statistic with a chi-squared distribution.)

A modified version of the differences-of-means t -statistic, based on a different standard error formula—the “pooled” standard error formula—has an exact Student t distribution when Y is normally distributed; however, the pooled standard error formula applies only in the special case that the two groups have the same variance or that each group has the same number of observations (Exercise 3.21). Adopt the

notation of Equation (3.19) so that the two groups are denoted as m and w . The pooled variance estimator is

$$s_{pooled}^2 = \frac{1}{n_m + n_w - 2} \left[\sum_{\substack{i=1 \\ \text{group } m}}^{n_m} (Y_i - \bar{Y}_m)^2 + \sum_{\substack{i=1 \\ \text{group } w}}^{n_w} (Y_i - \bar{Y}_w)^2 \right], \quad (3.23)$$

where the first summation is for the observations in group m and the second summation is for the observations in group w . The pooled standard error of the difference in means is $SE_{pooled}(\bar{Y}_m - \bar{Y}_w) = s_{pooled} \times \sqrt{1/n_m + 1/n_w}$, and the pooled t -statistic is computed using Equation (3.20), where the standard error is the pooled standard error, $SE_{pooled}(\bar{Y}_m - \bar{Y}_w)$.

If the population distribution of Y in group m is $N(\mu_m, \sigma_m^2)$, if the population distribution of Y in group w is $N(\mu_w, \sigma_w^2)$, and if the two group variances are the same (that is, $\sigma_m^2 = \sigma_w^2$), then under the null hypothesis the t -statistic computed using the pooled standard error has a Student t distribution with $n_m + n_w - 2$ degrees of freedom.

The drawback of using the pooled variance estimator s_{pooled}^2 is that it applies only if the two population variances are the same (assuming $n_m \neq n_w$). If the population variances are different, the pooled variance estimator is biased and inconsistent. If the population variances are different but the pooled variance formula is used, the null distribution of the pooled t -statistic is not a Student t distribution, even if the data are normally distributed; in fact, it does not even have a standard normal distribution in large samples. Therefore, the pooled standard error and the pooled t -statistic should not be used unless you have a good reason to believe that the population variances are the same.

Use of the Student t Distribution in Practice

For the problem of testing the mean of Y , the Student t distribution is applicable if the underlying population distribution of Y is normal. For economic variables, however, normal distributions are the exception (for example, see the boxes in Chapter 2 “The Distribution of Adulthood Earnings in the United Kingdom” and “The Unpegging of the Swiss Franc”). Even if the data are not normally distributed, the normal approximation to the distribution of the t -statistic is valid if the sample size is large. Therefore, inferences—hypothesis tests and confidence intervals—about the mean of a distribution should be based on the large-sample normal approximation.

When comparing two means, any economic reason for two groups having different means typically implies that the two groups also could have different variances. Accordingly, the pooled standard error formula is inappropriate, and the correct standard error formula, which allows for different group variances, is as given in Equation (3.19). Even if the population distributions are normal, the t -statistic computed using the standard error formula in Equation (3.19) does not have a Student

t distribution. In practice, therefore, inferences about differences in means should be based on Equation (3.19), used in conjunction with the large-sample standard normal approximation.

Even though the Student t distribution is rarely applicable in economics, some software uses the Student t distribution to compute p -values and confidence intervals. In practice, this does not pose a problem because the difference between the Student t distribution and the standard normal distribution is negligible if the sample size is large. For $n > 15$, the difference in the p -values computed using the Student t and standard normal distributions never exceeds 0.01; for $n > 80$, the difference never exceeds 0.002. In most modern applications, and in all applications in this text, the sample sizes are in the hundreds or thousands, large enough for the difference between the Student t distribution and the standard normal distribution to be negligible.

3.7 Scatterplots, the Sample Covariance, and the Sample Correlation

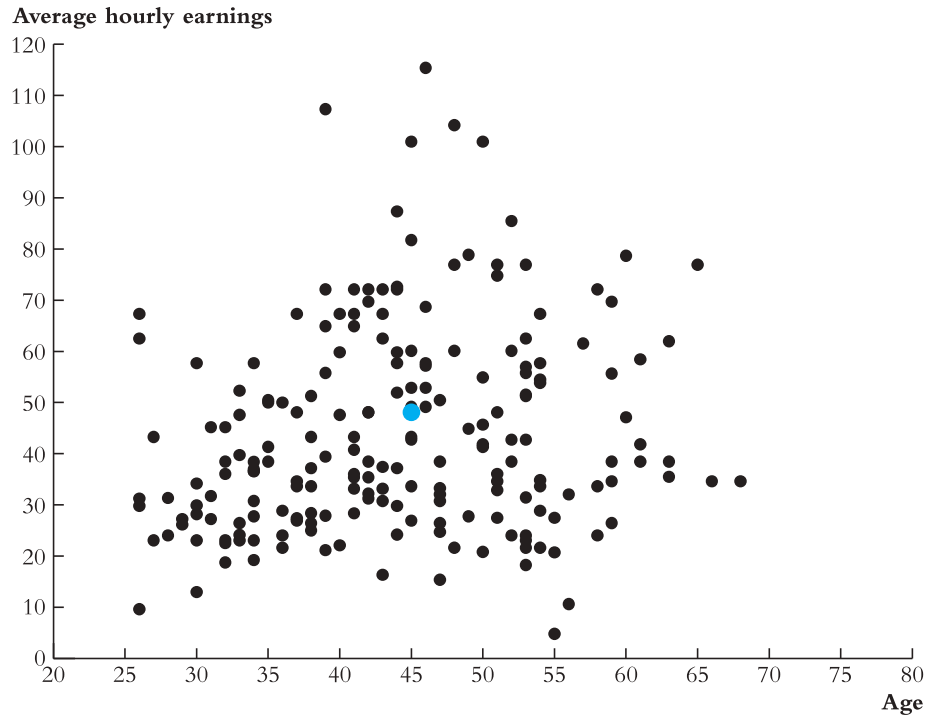
What is the relationship between age and earnings? This question, like many others, relates one variable, X (age), to another, Y (earnings). This section reviews three ways to summarize the relationship between variables: the scatterplot, the sample covariance, and the sample correlation coefficient.

Scatterplots

A **scatterplot** is a plot of n observations on X_i and Y_i , in which each observation is represented by the point (X_i, Y_i) . For example, Figure 3.2 is a scatterplot of age (X) and hourly earnings (Y) for a sample of 200 managers in the information industry from the March 2016 CPS. Each dot in Figure 3.2 corresponds to an (X, Y) pair for one of the observations. For example, one of the workers in this sample is 45 years old and earns \$49.15 per hour; this worker's age and earnings are indicated by the highlighted dot in Figure 3.2. The scatterplot shows a positive relationship between age and earnings in this sample: Older workers tend to earn more than younger workers. This relationship is not exact, however, and earnings could not be predicted perfectly using only a person's age.

Sample Covariance and Correlation

The covariance and correlation were introduced in Section 2.3 as two properties of the joint probability distribution of the random variables X and Y . Because the population distribution is unknown, in practice we do not know the population covariance or correlation. The population covariance and correlation can, however, be estimated by taking a random sample of n members of the population and collecting the data (X_i, Y_i) , $i = 1, \dots, n$.

FIGURE 3.2 Scatterplot of Average Hourly Earnings vs. Age

Each point in the plot represents the age and average earnings of one of the 200 workers in the sample. The highlighted dot corresponds to a 45-year-old worker who earns \$49.15 per hour. The data are for computer and information systems managers from the March 2016 CPS.

The sample covariance and correlation are estimators of the population covariance and correlation. Like the estimators discussed previously in this chapter, they are computed by replacing a population mean (the expectation) with a sample mean. The **sample covariance**, denoted s_{XY} , is

$$s_{XY} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}). \quad (3.24)$$

Like the sample variance, the average in Equation (3.24) is computed by dividing by $n-1$ instead of n ; here, too, this difference stems from using \bar{X} and \bar{Y} to estimate the respective population means. When n is large, it makes little difference whether division is by n or $n-1$.

The **sample correlation coefficient**, or **sample correlation**, is denoted r_{XY} and is the ratio of the sample covariance to the sample standard deviations:

$$r_{XY} = \frac{s_{XY}}{s_X s_Y}. \quad (3.25)$$

The sample correlation measures the strength of the linear association between X and Y in a sample of n observations. Like the population correlation, the sample correlation is unit free and lies between -1 and 1 : $|r_{XY}| \leq 1$.

The sample correlation equals 1 if $X_i = Y_i$ for all i and equals -1 if $X_i = -Y_i$ for all i . More generally, the correlation is ± 1 if the scatterplot is a straight line. If the line slopes upward, then there is a positive relationship between X and Y and the correlation is 1 . If the line slopes down, then there is a negative relationship and the correlation is -1 . The closer the scatterplot is to a straight line, the closer the correlation is to ± 1 . A high correlation coefficient does not necessarily mean that the line has a steep slope; rather, it means that the points in the scatterplot fall very close to a straight line.

Consistency of the sample covariance and correlation. Like the sample variance, the sample covariance is consistent. That is,

$$s_{XY} \xrightarrow{p} \sigma_{XY}. \quad (3.26)$$

In other words, in large samples the sample covariance is close to the population covariance with high probability.

The proof of the result in Equation (3.26) under the assumption that (X_i, Y_i) are i.i.d. and that X_i and Y_i have finite fourth moments is similar to the proof in Appendix 3.3 that the sample covariance is consistent and is left as an exercise (Exercise 3.20).

Because the sample variance and sample covariance are consistent, the sample correlation coefficient is consistent; that is, $r_{XY} \xrightarrow{p} \text{corr}(X_i, Y_i)$.

Example. As an example, consider the data on age and earnings in Figure 3.2. For these 200 workers, the sample standard deviation of age is $s_A = 9.57$ years, and the sample standard deviation of earnings is $s_E = \$19.93$ per hour. The sample covariance between age and earnings is $s_{AE} = 91.51$ (the units are years \times dollars per hour, not readily interpretable). Thus the sample correlation coefficient is $r_{AE} = 91.51 / (9.57 \times 19.93) = 0.48$. The correlation of 0.48 means that there is a positive relationship between age and earnings, but as is evident in the scatterplot, this relationship is far from perfect.

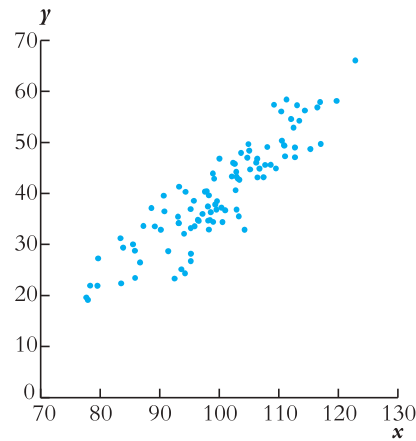
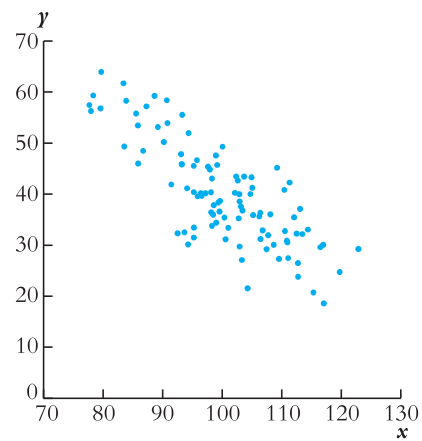
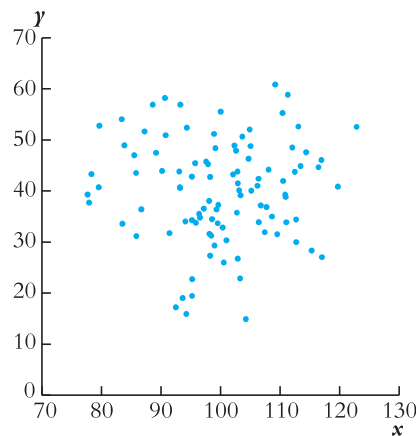
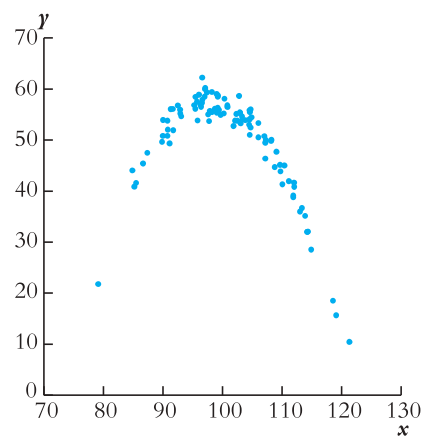
To verify that the correlation does not depend on the units of measurement, suppose that earnings had been reported in cents, in which case the sample standard deviation of earnings is 1993¢ per hour and the covariance between age and earnings is 9151 (units are years \times cents per hour); then the correlation is $9151 / (9.57 \times 1993) = 0.48$, or 48% .

Figure 3.3 gives additional examples of scatterplots and correlation. Figure 3.3a shows a strong positive linear relationship between these variables, and the sample correlation is 0.9 .

Figure 3.3b shows a strong negative relationship with a sample correlation of -0.8 . Figure 3.3c shows a scatterplot with no evident relationship, and the sample

FIGURE 3.3 Scatterplots for Four Hypothetical Data Sets

The scatterplots in Figures 3.3a and 3.3b show strong linear relationships between X and Y . In Figure 3.3c, X is independent of Y and the two variables are uncorrelated. In Figure 3.3d, the two variables also are uncorrelated even though they are related nonlinearly.

**(a)** Correlation = +0.9**(b)** Correlation = -0.8**(c)** Correlation = 0.0**(d)** Correlation = 0.0 (quadratic)

correlation is 0. Figure 3.3d shows a clear relationship: As X increases, Y initially increases but then decreases. Despite this discernable relationship between X and Y , the sample correlation is 0; the reason is that for these data small values of Y are associated with *both* large and small values of X .

This final example emphasizes an important point: The correlation coefficient is a measure of *linear* association. There is a relationship in Figure 3.3d, but it is not linear.

Summary

1. The sample average, \bar{Y} , is an estimator of the population mean, μ_Y . When Y_1, \dots, Y_n are i.i.d.,
 - a. the sampling distribution of \bar{Y} has mean μ_Y and variance $\sigma_{\bar{Y}}^2 = \sigma_Y^2/n$;
 - b. \bar{Y} is unbiased;
 - c. by the law of large numbers, \bar{Y} is consistent; and
 - d. by the central limit theorem, \bar{Y} has an approximately normal sampling distribution when the sample size is large.
2. The t -statistic is used to test the null hypothesis that the population mean takes on a particular value. If n is large, the t -statistic has a standard normal sampling distribution when the null hypothesis is true.
3. The t -statistic can be used to calculate the p -value associated with the null hypothesis. The p -value is the probability of drawing a statistic at least as adverse to the null hypothesis as the one you actually computed in your sample, assuming the null hypothesis is correct. A small p -value is evidence that the null hypothesis is false.
4. A 95% confidence interval for μ_Y is an interval constructed so that it contains the true value of μ_Y in 95% of all possible samples.
5. Hypothesis tests and confidence intervals for the difference in the means of two populations are conceptually similar to tests and intervals for the mean of a single population.
6. The sample correlation coefficient is an estimator of the population correlation coefficient and measures the linear relationship between two variables—that is, how well their scatterplot is approximated by a straight line.

Key Terms

estimator (105)	p -value (significance probability) (110)
estimate (105)	sample variance (112)
bias (106)	sample standard deviation (112)
consistency (106)	degrees of freedom (112)
efficiency (106)	standard error of \bar{Y} (113)
BLUE (Best Linear Unbiased Estimator) (107)	t -statistic (113)
least squares estimator (107)	t -ratio (113)
hypothesis tests (109)	test statistic (113)
null hypothesis (109)	type I error (115)
alternative hypothesis (109)	type II error (115)
two-sided alternative hypothesis (109)	significance level (115)
	critical value (115)

rejection region (115)	test for the difference between two means (119)
acceptance region (115)	causal effect (121)
size of a test (115)	treatment effect (121)
power of a test (115)	scatterplot (127)
one-sided alternative hypothesis (116)	sample covariance (128)
confidence set (117)	sample correlation coefficient (sample correlation) (128)
confidence level (117)	
confidence interval (117)	
coverage probability (118)	

MyLab Economics Can Help You Get a Better Grade

MyLab Economics If your exam were tomorrow, would you be ready? For each chapter, **MyLab Economics** Practice Tests and Study Plan help you prepare for your exams. You can also find the Exercises and all Review the Concepts Questions available now in **MyLab Economics**. To see how it works, turn to the **MyLab Economics** spread on the inside front cover of this text and then go to www.pearson.com/mylab/economics.

For additional Empirical Exercises and Data Sets, log on to the Companion Website at www.pearsonglobaleditions.com.

Review the Concepts

- 3.1 Explain the difference between an unbiased estimator and a consistent estimator.
- 3.2 What is meant by the efficiency of an estimator? Which estimator is known as BLUE?
- 3.3 A population distribution has a mean of 15 and a variance of 10. Determine the mean and variance of \bar{Y} from an i.i.d. sample from this population for (a) $n = 5$; (b) $n = 500$; and (c) $n = 5000$. Relate your answers to the law of large numbers.
- 3.4 What is the difference between standard error and standard deviation? How is the standard error of the sample mean calculated?
- 3.5 What is the difference between a null hypothesis and an alternative hypothesis? Among size, significance level, and power? Between a one-sided alternative hypothesis and a two-sided alternative hypothesis?
- 3.6 Why does a confidence interval contain more information than the result of a single hypothesis test?
- 3.7 What is a scatterplot? What statistical features of a dataset can be represented using a scatterplot diagram?
- 3.8 Sketch a hypothetical scatterplot for a sample of size 10 for two random variables with a population correlation of (a) 1.0; (b) -1.0 ; (c) 0.9; (d) -0.5 ; and (e) 0.0.

Exercises

- 3.1** In a population, $\mu_Y = 75$ and $\sigma_Y^2 = 45$. Use the central limit theorem to answer the following questions:
- In a random sample of size $n = 50$, find $\Pr(\bar{Y} < 73)$.
 - In a random sample of size $n = 90$, find $\Pr(76 < \bar{Y} < 77)$.
 - In a random sample of size $n = 120$, find $\Pr(\bar{Y} > 69)$.
- 3.2** Let Y be a Bernoulli random variable with success probability $\Pr(Y = 1) = p$, and let Y_1, \dots, Y_n be i.i.d. draws from this distribution. Let \hat{p} be the fraction of successes (1s) in this sample.
- Show that $\hat{p} = \bar{Y}$.
 - Show that \hat{p} is an unbiased estimator of p .
 - Show that $\text{var}(\hat{p}) = p(1 - p)/n$.
- 3.3** In a poll of 500 likely voters, 270 responded that they would vote for the candidate from the democratic party, while 230 responded that they would vote for the candidate from the republican party. Let p denote the fraction of all likely voters who preferred the democratic candidate at the time of the poll, and let \hat{p} be the fraction of survey respondents who preferred the democratic candidate.
- Use the poll results to estimate p .
 - Use the estimator of the variance of \hat{p} , $\hat{p}(1 - \hat{p})/n$, to calculate the standard error of your estimator.
 - What is the p -value for the test of $H_0: p = 0.5$, vs. $H_1: p \neq 0.5$?
 - What is the p -value for the test of $H_0: p = 0.5$, vs. $H_1: p > 0.5$?
 - Why do the results from (c) and (d) differ?
 - Did the poll contain statistically significant evidence that the democratic candidate was ahead of the republican candidate at the time of the poll? Explain.
- 3.4** Using the data in Exercise 3.3:
- Construct a 95% confidence interval for p .
 - Construct a 99% confidence interval for p .
 - Why is the interval in (b) wider than the interval in (a)?
 - Without doing any additional calculations, test the hypothesis $H_0: p = 0.50$ vs. $H_1: p \neq 0.50$ at the 5% significance level.
- 3.5** A survey of 1000 registered voters is conducted, and the voters are asked to choose between candidate A and candidate B. Let p denote the fraction of voters in the population who prefer candidate A, and let \hat{p} denote the fraction of voters in the sample who prefer candidate A.
- You are interested in the competing hypotheses $H_0: p = 0.4$ vs. $H_1: p \neq 0.4$. Suppose that you decide to reject H_0 if $|\hat{p} - 0.4| > 0.01$.

- i. What is the size of this test?
 - ii. Compute the power of this test if $p = 0.45$.
 - b.** In the survey, $\hat{p} = 0.44$.
 - i. Test $H_0: p = 0.4$ vs. $H_1: p \neq 0.4$ using a 10% significance level.
 - ii. Test $H_0: p = 0.4$ vs. $H_1: p < 0.4$ using a 10% significance level.
 - iii. Construct a 90% confidence interval for p .
 - iv. Construct a 99% confidence interval for p .
 - v. Construct a 60% confidence interval for p .
 - c.** Suppose that the survey is carried out 30 times, using independently selected voters in each survey. For each of these 30 surveys, a 90% confidence interval for p is constructed.
 - i. What is the probability that the true value of p is contained in all 30 of these confidence intervals?
 - ii. How many of these confidence intervals do you expect to contain the true value of p ?
 - d.** In survey jargon, the “margin of error” is $1.96 \times SE(\hat{p})$; that is, it is half the length of the 95% confidence interval. Suppose you want to design a survey that has a margin of error of at most 0.5%. That is, you want $\Pr(|\hat{p} - p| > 0.005 \leq 0.005) \leq 0.005$. How large should n be if the survey uses simple random sampling?
- 3.6** Let Y_1, \dots, Y_n be i.i.d. draws from a distribution with mean μ . A test of $H_0: \mu = 10$ vs. $H_1: \mu \neq 10$ using the usual t -statistic yields a p -value of 0.07.
- a.** Does the 90% confidence interval contain $\mu = 10$? Explain.
 - b.** Can you determine if $\mu = 8$ is contained in the 95% confidence interval? Explain.
- 3.7** In a given population, 50% of the likely voters are women. A survey using a simple random sample of 1000 landline telephone numbers finds 55% women. Is there evidence that the survey is biased? Explain.
- 3.8** A new version of the SAT is given to 1500 randomly selected high school seniors. The sample mean test score is 1230, and the sample standard deviation is 145. Construct a 95% confidence interval for the population mean test score for high school seniors.
- 3.9** Suppose that a plant manufactures integrated circuits with a mean life of 1000 hours and a standard deviation of 100 hours. An inventor claims to have developed an improved process that produces integrated circuits with a longer mean life and the same standard deviation. The plant manager randomly selects 50 integrated circuits produced by the process. She says that she will believe the inventor’s claim if the sample mean life of the integrated circuits

is greater than 1100 hours; otherwise, she will conclude that the new process is no better than the old process. Let μ denote the mean of the new process. Consider the null and alternative hypotheses $H_0: \mu = 1000$ vs. $H_1: \mu > 1000$.

- a. What is the size of the plant manager's testing procedure?
 - b. Suppose the new process is in fact better and has a mean integrated circuit life of 1150 hours. What is the power of the plant manager's testing procedure?
 - c. What testing procedure should the plant manager use if she wants the size of her test to be 1%?
- 3.10** Suppose a new standardized test is given to 150 randomly selected third-grade students in Amsterdam. The sample average score \bar{Y} on the test is 42 points, and the sample standard deviation, s_Y , is 6 points.
- a. The authors plan to administer the test to all third-grade students in Amsterdam. Construct a 99% confidence interval for the mean score of all third graders in Amsterdam.
 - b. Suppose the same test is given to 300 randomly selected third graders from Rotterdam, producing a sample average of 48 points and sample standard deviation of 10 points. Construct a 95% confidence interval for the difference in mean scores between Rotterdam and Amsterdam.
 - c. Can you conclude with a high degree of confidence that the population means for Rotterdam and Amsterdam students are different? (What is the standard error of the difference in the two sample means? What is the p -value of the test of no difference in means versus some difference?)
- 3.11** Consider the estimator \tilde{Y} , defined in Equation (3.1). Show that (a) $E(\tilde{Y}) = \mu_Y$ and (b) $\text{var}(\tilde{Y}) = 1.25\sigma_Y^2/n$.
- 3.12** To investigate possible gender discrimination in a British firm, a sample of 120 men and 150 women with similar job descriptions are selected at random. A summary of the resulting monthly salaries follows:

	Average Salary (\bar{Y})	Standard Deviation (s_Y)	n
Men	£8200	£450	120
Women	£7900	£520	150

- a. What do these data suggest about wage differences in the firm? Do they represent statistically significant evidence that average wages of men and women are different? (To answer this question, first, state the null and alternative hypotheses; second, compute the relevant t -statistic; third, compute the p -value associated with the t -statistic; and, finally, use the p -value to answer the question.)
- b. Do these data suggest that the firm is guilty of gender discrimination in its compensation policies? Explain.

3.13 Data on fifth-grade test scores (reading and mathematics) for 400 school districts in Brussels yield average score $\bar{Y} = 712.1$ and standard deviation $s_Y = 23.2$.

- Construct a 90% confidence interval for the mean test score in the population.
- When the districts were divided into districts with small classes (< 20 students per teacher) and large classes (≥ 20 students per teacher), the following results were found:

Class Size	Average Salary (\bar{Y})	Standard Deviation (s_Y)	n
Small	721.8	24.4	150
Large	710.9	20.6	250

Is there statistically significant evidence that the districts with smaller classes have higher average test scores? Explain.

3.14 Values of height in inches (X) and weight in pounds (Y) are recorded from a sample of 200 male college students. The resulting summary statistics are $\bar{X} = 71.2$ in., $\bar{Y} = 164$ lb, $s_X = 1.9$ in., $s_Y = 16.4$ lb, $s_{XY} = 22.54$ in. \times lb, and $r_{XY} = 0.8$. Convert these statistics to the metric system (meters and kilograms).

3.15 Y_a and Y_b are Bernoulli random variables from two different populations, denoted a and b . Suppose $E(Y_a) = p_a$ and $E(Y_b) = p_b$. A random sample of size n_a is chosen from population a , with a sample average denoted \hat{p}_a , and a random sample of size n_b is chosen from population b , with a sample average denoted \hat{p}_b . Suppose the sample from population a is independent of the sample from population b .

- Show that $E(\hat{p}_a) = p_a$ and $\text{var}(\hat{p}_a) = p_a(1 - p_a) / n_a$. Show that $E(\hat{p}_b) = p_b$ and $\text{var}(\hat{p}_b) = p_b(1 - p_b) / n_b$.

- Show that $\text{var}(\hat{p}_a - \hat{p}_b) = \frac{p_a(1 - p_a)}{n_a} + \frac{p_b(1 - p_b)}{n_b}$.
(Hint: Remember that the samples are independent.)

- Suppose n_a and n_b are large. Show that a 95% confidence interval for $p_a - p_b$ is given by $(\hat{p}_a - \hat{p}_b) \pm 1.96 \sqrt{\frac{\hat{p}_a(1 - \hat{p}_a)}{n_a} + \frac{\hat{p}_b(1 - \hat{p}_b)}{n_b}}$.
How would you construct a 90% confidence interval for $p_a - p_b$?

3.16 Assume that grades on a standardized test are known to have a mean of 500 for students in Europe. The test is administered to 600 randomly selected students in Ukraine; in this sample, the mean is 508, and the standard deviation (s) is 75.

- Construct a 95% confidence interval for the average test score for Ukrainian students.

- b.** Is there statistically significant evidence that Ukrainian students perform differently than other students in Europe?
- c.** Another 500 students are selected at random from Ukraine. They are given a 3-hour preparation course before the test is administered. Their average test score is 514, with a standard deviation of 65.
 - i. Construct a 95% confidence interval for the change in average test score associated with the prep course.
 - ii. Is there statistically significant evidence that the prep course helped?
- d.** The original 600 students are given the prep course and then are asked to take the test a second time. The average change in their test scores is 7 points, and the standard deviation of the change is 40 points.
 - i. Construct a 95% confidence interval for the change in average test scores.
 - ii. Is there statistically significant evidence that students will perform better on their second attempt, after taking the prep course?
 - iii. Students may have performed better in their second attempt because of the prep course or because they gained test-taking experience in their first attempt. Describe an experiment that would quantify these two effects.

3.17 Read the box “Social Class or Education? Childhood Circumstances and Adult Earnings Revisited” in Section 3.5.

- a.** Construct a 95% confidence interval for the difference in the household earnings of people whose father NS-SEC classification was higher between those with no educational qualifications and those with an undergraduate degree or more.
- b.** Construct a 95% confidence interval for the difference in the household earnings of people whose father NS-SEC classification was routine between those with no educational qualifications and those with an undergraduate degree or more.
- c.** Construct a 95% confidence interval for the difference between your answers calculated in parts **a** and **b**.

3.18 This exercise shows that the sample variance is an unbiased estimator of the population variance when Y_1, \dots, Y_n are i.i.d. with mean μ_Y and variance σ_Y^2 .

- a.** Use Equation (2.32) to show that

$$E(Y_i - \bar{Y})^2 = \text{var}(Y_i) - 2\text{cov}(Y_i, \bar{Y}) + \text{var}(\bar{Y}).$$
- b.** Use Equation (2.34) to show that $\text{cov}(\bar{Y}, Y_i) = \sigma_Y^2/n$.
- c.** Use the results in (a) and (b) to show that $E(s_Y^2) = \sigma_Y^2$.

- 3.19 a.** \bar{Y} is an unbiased estimator of μ_Y . Is \bar{Y}^2 an unbiased estimator of μ_Y^2 ?
- b.** \bar{Y} is a consistent estimator of μ_Y . Is \bar{Y}^2 a consistent estimator of μ_Y^2 ?
- 3.20** Suppose (X_i, Y_i) are i.i.d. with finite fourth moments. Prove that the sample covariance is a consistent estimator of the population covariance; that is, $s_{XY} \xrightarrow{P} \sigma_{XY}$, where s_{XY} is defined in Equation (3.24). (*Hint:* Use the strategy of Appendix 3.3.)
- 3.21** Show that the pooled standard error $[SE_{pooled}(\bar{Y}_m - \bar{Y}_w)]$ given following Equation (3.23) equals the usual standard error for the difference in means in Equation (3.19) when the two group sizes are the same ($n_m = n_w$).
- 3.22** Suppose $Y_i \sim i.i.d.N(\mu_Y, \sigma_Y^2)$ for $i = 1, \dots, n$. With σ_Y^2 known, the t -statistic for testing $H_0: \mu_Y = 0$ vs. $H_1: \mu_Y > 0$ is $t = (\bar{Y} - 0)/SE(\bar{Y})$, where $SE(\bar{Y}) = \sigma_Y/\sqrt{n}$. Suppose $\sigma_Y = 10$ and $n = 100$, so that $SE(\bar{Y}) = 1$. Using a test with a size of 5%, the null hypothesis is rejected if $t > 1.64$.
- a.** Suppose $\mu_Y = 0$, so the null hypothesis is true. What is the probability that the null hypothesis is rejected?
- b.** Suppose $\mu_Y = 2$, so the alternative hypothesis is true. What is the probability that the null hypothesis is rejected?
- c.** Suppose that in 90% of cases the data are drawn from a population where the null is true ($\mu_Y = 0$) and in 10% of cases the data come from a population where the alternative is true and $\mu_Y = 2$. Your data came from either the first or the second population, but you don't know which.
- You compute the t -statistic. What is the probability that $t > 1.64$ —that is, that you reject the null hypothesis?
 - Suppose you reject the null hypothesis; that is, $t > 1.64$. What is the probability that the sample data were drawn from the $\mu_Y = 0$ population?
- d.** It is hard to discover a new effective drug. Suppose 90% of new drugs are ineffective and only 10% are effective. Let Y denote the drop in the level of a specific blood toxin for a patient taking a new drug. If the drug is ineffective, $\mu_Y = 0$ and $\sigma_Y = 10$; if the drug is effective, $\mu_Y = 2$ and $\sigma_Y = 10$.
- A new drug is tested on a random sample of $n = 100$ patients, data are collected, and the resulting t -statistic is found to be greater than 1.64. What is the probability that the drug is ineffective (i.e., what is the false positive rate for the test using $t > 1.64$)?
 - Suppose the one-sided test uses instead the 0.5% significance level. What is the probability that the drug is ineffective (i.e., what is the false positive rate)?

Empirical Exercises

- E3.1** On the text website, <http://www.pearsonglobaleditions.com>, you will find the data file **CPS96_15**, which contains an extended version of the data set used in Table 3.1 of the text for the years 1996 and 2015. It contains data on full-time workers, ages 25–34, with a high school diploma or a B.A./B.S. as their highest degree. A detailed description is given in **CPS96_15_Description**, available on the website. Use these data to complete the following.
- a.**
 - i. Compute the sample mean for average hourly earnings (*AHE*) in 1996 and 2015.
 - ii. Compute the sample standard deviation for *AHE* in 1996 and 2015.
 - iii. Construct a 95% confidence interval for the population means of *AHE* in 1996 and 2015.
 - iv. Construct a 95% confidence interval for the change in the population means of *AHE* between 1996 and 2015.
 - b.** In 2015, the value of the Consumer Price Index (CPI) was 237.0. In 1996, the value of the CPI was 156.9. Repeat (a), but use *AHE* measured in real 2015 dollars (\$2015); that is, adjust the 1996 data for the price inflation that occurred between 1996 and 2015.
 - c.** If you were interested in the change in workers' purchasing power from 1996 to 2015, would you use the results from (a) or (b)? Explain.
 - d.** Using the data for 2015:
 - i. Construct a 95% confidence interval for the mean of *AHE* for high school graduates.
 - ii. Construct a 95% confidence interval for the mean of *AHE* for workers with a college degree.
 - iii. Construct a 95% confidence interval for the difference between the two means.
 - e.** Repeat (d) using the 1996 data expressed in \$2015.
 - f.** Using appropriate estimates, confidence intervals, and test statistics, answer the following questions:
 - i. Did real (inflation-adjusted) wages of high school graduates increase from 1996 to 2015?
 - ii. Did real wages of college graduates increase?
 - iii. Did the gap between earnings of college and high school graduates increase? Explain.
 - g.** Table 3.1 presents information on the gender gap for college graduates. Prepare a similar table for high school graduates, using the 1996 and 2015 data. Are there any notable differences between the results for high school and college graduates?

E3.2 A consumer is given the chance to buy a baseball card for \$1, but he declines the trade. If the consumer is now given the baseball card, will he be willing to sell it for \$1? Standard consumer theory suggests yes, but behavioral economists have found that “ownership” tends to increase the value of goods to consumers. That is, the consumer may hold out for some amount more than \$1 (for example, \$1.20) when selling the card, even though he was willing to pay only some amount less than \$1 (for example, \$0.88) when buying it. Behavioral economists call this phenomenon the “endowment effect.” John List investigated the endowment effect in a randomized experiment involving sports memorabilia traders at a sports-card show. Traders were randomly given one of two sports collectibles, say good A or good B, that had approximately equal market value.³ Those receiving good A were then given the option of trading good A for good B with the experimenter; those receiving good B were given the option of trading good B for good A with the experimenter. Data from the experiment and a detailed description can be found on the text website, <http://www.pearsonglobaleditions.com>, in the files **Sportscards** and **Sportscards_Description**.⁴

- a.
 - i. Suppose that, absent any endowment effect, all the subjects prefer good A to good B. What fraction of the experiment’s subjects would you expect to trade the good that they were given for the other good? (*Hint:* Because of random assignment of the two treatments, approximately 50% of the subjects received good A, and 50% received good B.)
 - ii. Suppose that, absent any endowment effect, 50% of the subjects prefer good A to good B, and the other 50% prefer good B to good A. What fraction of the subjects would you expect to trade the good they were given for the other good?
 - iii. Suppose that, absent any endowment effect, $X\%$ of the subjects prefer good A to good B, and the other $(100 - X)\%$ prefer good B to good A. Show that you would expect 50% of the subjects to trade the good they were given for the other good.
- b. Using the sports-card data, what fraction of the subjects traded the good they were given? Is the fraction significantly different from 50%? Is there evidence of an endowment effect? (*Hint:* Review Exercises 3.2 and 3.3.)
- c. Some have argued that the endowment effect may be present but that it is likely to disappear as traders gain more trading experience. Half of the experimental subjects were dealers, and the other half were nondealers. Dealers have more experience than nondealers. Repeat (b) for dealers and nondealers. Is there a significant difference in their behavior?

³Good A was a ticket stub from the game in which Cal Ripken, Jr., set the record for consecutive games played, and good B was a souvenir from the game in which Nolan Ryan won his 300th game.

⁴These data were provided by Professor John List of the University of Chicago and were used in his paper “Does Market Experience Eliminate Market Anomalies,” *Quarterly Journal of Economics*, 2003, 118(1): 41–71.

Is the evidence consistent with the hypothesis that the endowment effect disappears as traders gain more experience? (*Hint:* Review Exercise 3.15.)

APPENDIX

3.1 The U.S. Current Population Survey

Each month the U.S. Census Bureau and the U.S. Bureau of Labor Statistics conduct the Current Population Survey (CPS), which provides data on labor force characteristics of the population, including the levels of employment, unemployment, and earnings. Approximately 54,000 U.S. households are surveyed each month. The sample is chosen by randomly selecting addresses from a database of addresses from the most recent decennial census augmented with data on new housing units constructed after the last census. The exact random sampling scheme is rather complicated (first, small geographical areas are randomly selected; then housing units within these areas are randomly selected); details can be found in the *Handbook of Labor Statistics* and on the Bureau of Labor Statistics website (www.bls.gov).

The survey conducted each March is more detailed than those in other months and asks questions about earnings during the previous year. The statistics in Tables 2.4 and 3.1 were computed using the March surveys. The CPS earnings data are for full-time workers, defined to be persons employed more than 35 hours per week for at least 48 weeks in the previous year.

More details on the data can be found in the replication materials for this chapter, available at <http://www.pearsonglobaleditions.com>.

APPENDIX

3.2 Two Proofs That \bar{Y} Is the Least Squares Estimator of μ_Y

This appendix provides two proofs, one using calculus and one not, that \bar{Y} minimizes the sum of squared prediction mistakes in Equation (3.2)—that is, that \bar{Y} is the least squares estimator of $E(Y)$.

Calculus Proof

To minimize the sum of squared prediction mistakes, take its derivative and set it to 0:

$$\frac{d}{dm} \sum_{i=1}^n (Y_i - m)^2 = -2 \sum_{i=1}^n (Y_i - m) = -2 \sum_{i=1}^n Y_i + 2nm = 0. \quad (3.27)$$

Solving for the final equation for m shows that $\sum_{i=1}^n (Y_i - m)^2$ is minimized when $m = \bar{Y}$.

Noncalculus Proof

The strategy is to show that the difference between the least squares estimator and \bar{Y} must be 0, from which it follows that \bar{Y} is the least squares estimator. Let $d = \bar{Y} - m$, so that $m = \bar{Y} - d$. Then $(Y_i - m)^2 = (Y_i - [\bar{Y} - d])^2 = ([Y_i - \bar{Y}] + d)^2 = (Y_i - \bar{Y})^2 + 2d(Y_i - \bar{Y}) + d^2$. Thus the sum of squared prediction mistakes [Equation (3.2)] is

$$\sum_{i=1}^n (Y_i - m)^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 + 2d \sum_{i=1}^n (Y_i - \bar{Y}) + nd^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 + nd^2, \quad (3.28)$$

where the second equality uses the fact that $\sum_{i=1}^n (Y_i - \bar{Y}) = 0$. Because both terms in the final line of Equation (3.28) are nonnegative and because the first term does not depend on d , $\sum_{i=1}^n (Y_i - m)^2$ is minimized by choosing d to make the second term, nd^2 , as small as possible. This is done by setting $d = 0$ —that is, by setting $m = \bar{Y}$ —so that \bar{Y} is the least squares estimator of $E(Y)$.

APPENDIX

3.3 A Proof That the Sample Variance Is Consistent

This appendix uses the law of large numbers to prove that the sample variance, s_Y^2 , is a consistent estimator of the population variance, σ_Y^2 , as stated in Equation (3.9), when Y_1, \dots, Y_n are i.i.d. and $E(Y_i^4) < \infty$.

First, consider a version of the sample variance that uses n instead of $n - 1$ as a divisor:

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2 &= \frac{1}{n} \sum_{i=1}^n Y_i^2 - 2\bar{Y} \frac{1}{n} \sum_{i=1}^n Y_i + \bar{Y}^2 \\ &= \frac{1}{n} \sum_{i=1}^n Y_i^2 - \bar{Y}^2 \\ &\xrightarrow{p} (\sigma_Y^2 + \mu_Y^2) - \mu_Y^2 \\ &= \sigma_Y^2, \end{aligned} \quad (3.29)$$

where the first equality uses $(Y_i - \bar{Y})^2 = Y_i^2 - 2\bar{Y}Y_i + \bar{Y}^2$ and the second uses $\frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y}$.

The convergence in the third line follows from (i) applying the law of large numbers to $\frac{1}{n} \sum_{i=1}^n Y_i^2 \xrightarrow{p} E(Y^2)$ (which follows because Y_i^2 are i.i.d. and have finite variance because $E(Y_i^4)$ is finite), (ii) recognizing that $E(Y_i^2) = \sigma_Y^2 + \mu_Y^2$ (Key Concept 2.3), and (iii) noting $\bar{Y} \xrightarrow{p} \mu_Y$, so that $\bar{Y}^2 \xrightarrow{p} \mu_Y^2$. Finally, $s_Y^2 = \left(\frac{n}{n-1}\right) \left(\frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2\right) \xrightarrow{p} \sigma_Y^2$ follows from Equation (3.29) and $\left(\frac{n}{n-1}\right) \rightarrow 1$.

Linear Regression with One Regressor

The superintendent of an elementary school district must decide whether to hire additional teachers, and she wants your advice. Hiring the teachers will reduce the number of students per teacher (the student–teacher ratio) by two but will increase the district’s expenses. So she asks you: If she cuts class sizes by two, what will the effect be on student performance, as measured by scores on standardized tests?

Now suppose a father tells you that his family wants to move to a town with a good school system. He is interested in a specific school district: Test scores for this district are not publicly available, but the father knows its class size, based on the district’s student–teacher ratio. So he asks you: if he tells you the district’s class size, could you predict that district’s standardized test scores?

These two questions are clearly related: They both pertain to the relation between class size and test scores. Yet they are different. To answer the superintendent’s question, you need an estimate of the causal effect of a change in one variable (the student–teacher ratio, X) on another (test scores, Y). To answer the father’s question, you need to know how X relates to Y , on average, across school districts so you can use this relation to predict Y given X in a specific district.

These two questions are examples of two different types of questions that arise in econometrics. The first type of questions pertains to **causal inference**: using data to estimate the effect on an outcome of interest of an intervention that changes the value of another variable. The second type of questions concerns **prediction**: using the observed value of some variable to predict the value of another variable.

This chapter introduces the linear regression model relating one variable, X , to another, Y . This model postulates a linear relationship between X and Y . Just as the mean of Y is an unknown characteristic of the population distribution of Y , the intercept and slope of the line relating X and Y are unknown characteristics of the population joint distribution of X and Y . The econometric problem is to estimate the intercept and slope using a sample of data on these two variables.

Like the differences in means, linear regression is a statistical procedure that can be used for causal inference and for prediction. The two uses, however, place different requirements on the data. Section 3.5 explained how a difference in mean outcomes between a treatment and a control group estimates the causal effect of the treatment when the treatment is randomly assigned in an experiment. When X is continuous, computing differences-in-means no longer works because there are many values X can take on, not just two. If, however, we make the additional assumption that the relation between X and Y is linear, then if X is randomly assigned, we can use linear regression to estimate the causal effect on Y of an intervention that changes X . Even if X is not randomly assigned,