

# CHAPTER 3

# Descriptive Statistics: Numerical Measures

#### CONTENTS

STATISTICS IN PRACTICE: SMALL FRY DESIGN

- 3.1 MEASURES OF LOCATION Mean Median Mode Percentiles Quartiles
- MEASURES OF VARIABILITY Range Interquartile Range Variance Standard Deviation Coefficient of Variation
- 3.3 MEASURES OF DISTRIBUTION SHAPE, RELATIVE LOCATION, AND DETECTING OUTLIERS Distribution Shape z-Scores

- Chebyshev's Theorem Empirical Rule **Detecting Outliers**
- 3.4 EXPLORATORY DATA ANALYSIS Five-Number Summary Box Plot
- 3.5 MEASURES OF ASSOCIATION BETWEEN TWO VARIABLES Covariance Interpretation of the Covariance Correlation Coefficient Interpretation of the Correlation Coefficient
- THE WEIGHTED MEAN AND WORKING WITH GROUPED DATA Weighted Mean Grouped Data

# STATISTICS (in PRACTICE

#### **SMALL FRY DESIGN\*** SANTA ANA, CALIFORNIA

Founded in 1997, Small Fry Design is a toy and accessory company that designs and imports products for infants. The company's product line includes teddy bears, mobiles, musical toys, rattles, and security blankets and features high-quality soft toy designs with an emphasis on color, texture, and sound. The products are designed in the United States and manufactured in China.

Small Fry Design uses independent representatives to sell the products to infant furnishing retailers, children's accessory and apparel stores, gift shops, upscale department stores, and major catalog companies. Currently, Small Fry Design products are distributed in more than 1000 retail outlets throughout the United States.

Cash flow management is one of the most critical activities in the day-to-day operation of this company. Ensuring sufficient incoming cash to meet both current and ongoing debt obligations can mean the difference between business success and failure. A critical factor in cash flow management is the analysis and control of accounts receivable. By measuring the average age and dollar value of outstanding invoices, management can predict cash availability and monitor changes in the status of accounts receivable. The company set the following goals: the average age for outstanding invoices should not exceed 45 days, and the dollar value of invoices more than 60 days old should not exceed 5% of the dollar value of all accounts receivable.

In a recent summary of accounts receivable status, the following descriptive statistics were provided for the age of outstanding invoices:

Mean	40 days
Median	35 days
Mode	31 days

<sup>\*</sup>The authors are indebted to John A. McCarthy, President of Small Fry Design, for providing this Statistics in Practice.



Small Fry Design's "King of the Jungle" mobile. © Photo courtesy of Small Fry Design, Inc.

Interpretation of these statistics shows that the mean or average age of an invoice is 40 days. The median shows that half of the invoices remain outstanding 35 days or more. The mode of 31 days, the most frequent invoice age, indicates that the most common length of time an invoice is outstanding is 31 days. The statistical summary also showed that only 3% of the dollar value of all accounts receivable was more than 60 days old. Based on the statistical information, management was satisfied that accounts receivable and incoming cash flow were under control.

In this chapter, you will learn how to compute and interpret some of the statistical measures used by Small Fry Design. In addition to the mean, median, and mode, you will learn about other descriptive statistics such as the range, variance, standard deviation, percentiles, and correlation. These numerical measures will assist in the understanding and interpretation of data.

In Chapter 2 we discussed tabular and graphical presentations used to summarize data. In this chapter, we present several numerical measures that provide additional alternatives for summarizing data.

We start by developing numerical summary measures for data sets consisting of a single variable. When a data set contains more than one variable, the same numerical measures can be computed separately for each variable. However, in the two-variable case, we will also develop measures of the relationship between the variables.

Numerical measures of location, dispersion, shape, and association are introduced. If the measures are computed for data from a sample, they are called **sample statistics**. If the measures are computed for data from a population, they are called **population parameters**. In statistical inference, a sample statistic is referred to as the **point estimator** of the corresponding population parameter. In Chapter 7 we will discuss in more detail the process of point estimation.

In the two chapter appendixes we show how Minitab and Excel can be used to compute many of the numerical measures described in the chapter.

# 3.1

# Measures of Location

#### Mean

Perhaps the most important measure of location is the **mean**, or average value, for a variable. The mean provides a measure of central location for the data. If the data are for a sample, the mean is denoted by  $\bar{x}$ ; if the data are for a population, the mean is denoted by the Greek letter  $\mu$ .

In statistical formulas, it is customary to denote the value of variable x for the first observation by  $x_1$ , the value of variable x for the second observation by  $x_2$ , and so on. In general, the value of variable x for the ith observation is denoted by  $x_i$ . For a sample with n observations, the formula for the sample mean is as follows.

The sample mean  $\bar{x}$  is a sample statistic.

SAMPLE MEAN

$$\bar{x} = \frac{\sum x_i}{n} \tag{3.1}$$

In the preceding formula, the numerator is the sum of the values of the n observations. That is,

$$\sum x_i = x_1 + x_2 + \cdots + x_n$$

The Greek letter  $\Sigma$  is the summation sign.

To illustrate the computation of a sample mean, let us consider the following class size data for a sample of five college classes.

We use the notation  $x_1$ ,  $x_2$ ,  $x_3$ ,  $x_4$ ,  $x_5$  to represent the number of students in each of the five classes.

$$x_1 = 46$$
  $x_2 = 54$   $x_3 = 42$   $x_4 = 46$   $x_5 = 32$ 

Hence, to compute the sample mean, we can write

$$\bar{x} = \frac{\sum x_i}{n} = \frac{x_1 + x_2 + x_3 + x_4 + x_5}{5} = \frac{46 + 54 + 42 + 46 + 32}{5} = 44$$

The sample mean class size is 44 students.

Another illustration of the computation of a sample mean is given in the following situation. Suppose that a college placement office sent a questionnaire to a sample of business school graduates requesting information on monthly starting salaries. Table 3.1 shows the

STUDENTS-HUB.com

TABLE 3.1 MONTHLY STARTING SALARIES FOR A SAMPLE OF 12 BUSINESS SCHOOL GRADUATES



Graduate	Monthly Starting Salary (\$)	Graduate	Monthly Starting Salary (\$)
1	3450	7	3490
2	3550	8	3730
3	3650	9	3540
4	3480	10	3925
5	3355	H W III	3520
6	3310	12	3480

collected data. The mean monthly starting salary for the sample of 12 business college graduates is computed as

$$\bar{x} = \frac{\sum x_i}{n} = \frac{x_1 + x_2 + \dots + x_{12}}{12}$$

$$= \frac{3450 + 3550 + \dots + 3480}{12}$$

$$= \frac{42,480}{12} = 3540$$

Equation (3.1) shows how the mean is computed for a sample with n observations. The formula for computing the mean of a population remains the same, but we use different notation to indicate that we are working with the entire population. The number of observations in a population is denoted by N and the symbol for a population mean is  $\mu$ .

The sample mean  $\bar{x}$  is a point estimator of the population mean  $\mu$ .

POPULATION MEAN

$$u = \frac{\sum x_i}{N} \tag{3.2}$$

## Median

The **median** is another measure of central location. The median is the value in the middle when the data are arranged in ascending order (smallest value to largest value). With an odd number of observations, the median is the middle value. An even number of observations has no single middle value. In this case, we follow convention and define the median as the average of the values for the middle two observations. For convenience the definition of the median is restated as follows.

#### **MEDIAN**

Arrange the data in ascending order (smallest value to largest value).

- (a) For an odd number of observations, the median is the middle value.
- (b) For an even number of observations, the median is the average of the two middle values.

The median is the measure

reported for annual income and property value data

because a few extremely

large incomes or property

values can inflate the mean

In such cases, the median is

the preferred measure of

central location.

of location most often

3.1 Measures of Location

Let us apply this definition to compute the median class size for the sample of five college classes. Arranging the data in ascending order provides the following list.

Because n = 5 is odd, the median is the middle value. Thus the median class size is 46 students. Even though this data set contains two observations with values of 46, each observation is treated separately when we arrange the data in ascending order.

Suppose we also compute the median starting salary for the 12 business college graduates in Table 3.1. We first arrange the data in ascending order.

Because n = 12 is even, we identify the middle two values: 3490 and 3520. The median is the average of these values.

$$Median = \frac{3490 + 3520}{2} = 3505$$

Although the mean is the more commonly used measure of central location, in some situations the median is preferred. The mean is influenced by extremely small and large data values. For instance, suppose that one of the graduates (see Table 3.1) had a starting salary of \$10,000 per month (maybe the individual's family owns the company). If we change the highest monthly starting salary in Table 3.1 from \$3925 to \$10,000 and recompute the mean, the sample mean changes from \$3540 to \$4046. The median of \$3505, however, is unchanged, because \$3490 and \$3520 are still the middle two values. With the extremely high starting salary included, the median provides a better measure of central location than the mean. We can generalize to say that whenever a data set contains extreme values, the median is often the preferred measure of central location.

#### Mode

A third measure of location is the **mode**. The mode is defined as follows.

#### MODE

The mode is the value that occurs with greatest frequency.

To illustrate the identification of the mode, consider the sample of five class sizes. The only value that occurs more than once is 46. Because this value, occurring with a frequency of 2, has the greatest frequency, it is the mode. As another illustration, consider the sample of starting salaries for the business school graduates. The only monthly starting salary that occurs more than once is \$3480. Because this value has the greatest frequency, it is the mode.

Situations can arise for which the greatest frequency occurs at two or more different values. In these instances more than one mode exists. If the data contain exactly two modes, we say that the data are *bimodal*. If data contain more than two modes, we say that the data are *multimodal*. In multimodal cases the mode is almost never reported because listing three or more modes would not be particularly helpful in describing a location for the data.

STUDENTS-HUB.com

#### **Percentiles**

A **percentile** provides information about how the data are spread over the interval from the smallest value to the largest value. For data that do not contain numerous repeated values, the pth percentile divides the data into two parts. Approximately p percent of the observations have values less than the pth percentile; approximately (100 - p) percent of the observations have values greater than the pth percentile. The pth percentile is formally defined as follows.

#### PERCENTILE

The pth percentile is a value such that at least p percent of the observations are less than or equal to this value and at least (100 - p) percent of the observations are greater than or equal to this value.

Colleges and universities frequently report admission test scores in terms of percentiles. For instance, suppose an applicant obtains a raw score of 54 on the verbal portion of an admission test. How this student performed in relation to other students taking the same test may not be readily apparent. However, if the raw score of 54 corresponds to the 70th percentile, we know that approximately 70% of the students scored lower than this individual and approximately 30% of the students scored higher than this individual.

The following procedure can be used to compute the pth percentile.

#### CALCULATING THE pTH PERCENTILE

Step 1. Arrange the data in ascending order (smallest value to largest value).

Step 2. Compute an index i

$$i = \left(\frac{\hat{p}}{100}\right)n$$

where p is the percentile of interest and n is the number of observations.

**Step 3.** (a) If *i is not an integer, round up.* The next integer *greater* than *i* denotes the position of the *p*th percentile.

(b) If i is an integer, the pth percentile is the average of the values in positions i and i + 1.

As an illustration of this procedure, let us determine the 85th percentile for the starting salary data in Table 3.1.

**Step 1.** Arrange the data in ascending order.

3310 3355 3450 3480 3480 3490 3520 3540 3550 3650 3730 3925

Step 2.

$$i = \left(\frac{p}{100}\right)n = \left(\frac{85}{100}\right)12 = 10.2$$

**Step 3.** Because *i* is not an integer, *round up*. The position of the 85th percentile is the next integer greater than 10.2, the 11th position.

Returning to the data, we see that the 85th percentile is the data value in the 11th position, or 3730.

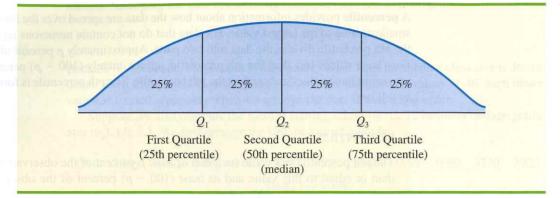
Uploaded By: Haneen

Following these steps

percentiles.

makes it easy to calculate

#### FIGURE 3.1 LOCATION OF THE QUARTILES



As another illustration of this procedure, let us consider the calculation of the 50th percentile for the starting salary data. Applying step 2, we obtain

$$i = \left(\frac{50}{100}\right)12 = 6$$

Because i is an integer, step 3(b) states that the 50th percentile is the average of the sixth and seventh data values; thus the 50th percentile is (3490 + 3520)/2 = 3505. Note that the 50th percentile is also the median.

# Quartiles

Quartiles are just specific percentiles; thus, the steps

for computing percentiles

can be applied directly in

the computation of

quartiles.

It is often desirable to divide data into four parts, so that each part contains approximately one-fourth, or 25% of the observations. Figure 3.1 shows a data distribution divided into four parts. The division points are referred to as the **quartiles** and are defined as

 $Q_1$  = first quartile, or 25th percentile

 $Q_2$  = second quartile, or 50th percentile (also the median)

 $Q_3$  = third quartile, or 75th percentile.

The starting salary data are again arranged in ascending order. We already identified  $Q_2$ , the second quartile (median), as 3505.

3310 3355 3450 3480 3480 3490 3520 3540 3550 3650 3730 3925

The computations of quartiles  $Q_1$  and  $Q_3$  require the use of the rule for finding the 25th and 75th percentiles. These calculations follow.

For  $Q_1$ ,

$$i = \left(\frac{p}{100}\right)n = \left(\frac{25}{100}\right)12 = 3$$

Because i is an integer, step 3(b) indicates that the first quartile, or 25th percentile, is the average of the third and fourth data values; thus,  $Q_1 = (3450 + 3480)/2 = 3465$ . For  $Q_3$ ,

$$i = \left(\frac{p}{100}\right)n = \left(\frac{75}{100}\right)12 = 9$$

Again, because *i* is an integer, step 3(b) indicates that the third quartile, or 75th percentile, is the average of the ninth and tenth data values; thus,  $Q_3 = (3550 + 3650)/2 = 3600$ .

STUDENTS-HUB.com

The quartiles divide the starting salary data into four parts, with each part containing 25% of the observations.

3310 3355 3450 3480 3480 3490 3520 3540 3550 3650 3730 3925 
$$Q_1 = 3465$$
  $Q_2 = 3505$   $Q_3 = 3600$  (Median)

We defined the quartiles as the 25th, 50th, and 75th percentiles. Thus, we computed the quartiles in the same way as percentiles. However, other conventions are sometimes used to compute quartiles, and the actual values reported for quartiles may vary slightly depending on the convention used. Nevertheless, the objective of all procedures for computing quartiles is to divide the data into four equal parts.

#### **NOTES AND COMMENTS**

It is better to use the median than the mean as a measure of central location when a data set contains extreme values. Another measure, sometimes used when extreme values are present, is the *trimmed mean*. It is obtained by deleting a percentage of the smallest and largest values from a data set and then computing the mean of the remaining values. For example, the 5% trimmed mean is obtained by re-

moving the smallest 5% and the largest 5% of the data values and then computing the mean of the remaining values. Using the sample with n = 12 starting salaries, 0.05(12) = 0.6. Rounding this value to 1 indicates that the 5% trimmed mean would remove the 1 smallest data value and the 1 largest data value. The 5% trimmed mean using the 10 remaining observations is 3524.50.

#### **Exercises**

# Methods

- 1. Consider a sample with data values of 10, 20, 12, 17, and 16. Compute the mean and median.
- 2. Consider a sample with data values of 10, 20, 21, 17, 16, and 12. Compute the mean and median.
- 3. Consider a sample with data values of 27, 25, 20, 15, 30, 34, 28, and 25. Compute the 20th, 25th, 65th, and 75th percentiles.
- 4. Consider a sample with data values of 53, 55, 70, 58, 64, 57, 53, 69, 57, 68, and 53. Compute the mean, median, and mode.

# **Applications**

5. The Dow Jones Travel Index reported what business travelers pay for hotel rooms per night in major U.S. cities (*The Wall Street Journal*, January 16, 2004). The average hotel room rates for 20 cities are as follows:



SELF tes

Atlanta	\$163	Minneapolis	\$125
Boston	177	New Orleans	167
Chicago	166	New York	245
Cleveland	126	Orlando	146
Dallas	123	Phoenix	139
Denver	120	Pittsburgh	134
Detroit	144	San Francisco	167
Houston	173	Seattle	162
Los Angeles	160	St. Louis	145
Miami	192	Washington, D.C.	207

- a. What is the mean hotel room rate?
- b. What is the median hotel room rate?
- c. What is the mode?
- d. What is the first quartile?
- e. What is the third quartile?
- 6. The National Association of Colleges and Employers compiled information about annual starting salaries for college graduates by major. The mean starting salary for business administration graduates was \$39,850 (CNNMoney.com, February 15, 2006). Samples with annual starting data for marketing majors and accounting majors follow (data are in thousands):

CD	file
	BASalary

Market	ting Majo	ors						
34.2	45.0	39.5	8.4 37.7	35.8	30.6	35.2	34.2	42.4
Accoun	nting Maj	ors						
33.5	57.1	49.7	40.2	44.2	45.	2	47.8	38.0
53.9	41.1	41.7	40.8	55.5	43.	5	49.1	49.9

- a. Compute the mean, median, and mode of the annual starting salary for both majors.
- b. Compute the first and third quartiles for both majors.
- c. Business administration students with accounting majors generally obtain the highest annual salary after graduation. What do the sample data indicate about the difference between the annual starting salaries for marketing and accounting majors?
- 7. The American Association of Individual Investors conducted an annual survey of discount brokers (*AAII Journal*, January 2003). The commissions charged by 24 discount brokers for two types of trades, a broker-assisted trade of 100 shares at \$50 per share and an online trade of 500 shares at \$50 per share, are shown in Table 3.2.
  - a. Compute the mean, median, and mode for the commission charged on a broker-assisted trade of 100 shares at \$50 per share.
  - b. Compute the mean, median, and mode for the commission charged on an online trade of 500 shares at \$50 per share.
  - c. Which costs more, a broker-assisted trade of 100 shares at \$50 per share or an online trade of 500 shares at \$50 per share?
  - d. Is the cost of a transaction related to the amount of the transaction?

#### TABLE 3.2 COMMISSIONS CHARGED BY DISCOUNT BROKERS



Broker	Broker- Assisted 100 Shares at \$50/Share	Online 500 Shares at \$50/Share	Broker	Assisted 100 Shares at \$50/Share	Online 500 Shares at \$50/Share
Accutrade	30.00	29.95	Merrill Lynch Direct	50.00	29.95
Ameritrade	24.99	10.99	Muriel Siebert	45.00	14.95
Banc of America	54.00	24.95	NetVest	24.00	14.00
Brown & Co.	17.00	5.00	Recom Securities	35.00	12.95
Charles Schwab	55.00	29.95	Scottrade	17.00	7.00
CyberTrader	12.95	9.95	Sloan Securities	39.95	19.95
E*TRADE Securities	49.95	14.95	Strong Investments	55.00	24.95
First Discount	35.00	19.75	TD Waterhouse	45.00	17.95
Freedom Investments	25.00	15.00	T. Rowe Price	50.00	19.95
Harrisdirect	40.00	20.00	Vanguard	48.00	20.00
Investors National	39.00	62.50	Wall Street Discount	29.95	19.95
MB Trading	9.95	10.55	York Securities	40.00	36.00

STUDENTS-HUB.com

8. The cost of consumer purchases such as housing, gasoline, Internet services, tax preparation, and hospitalization were provided in *The Wall Street Journal*, January 2, 2007. Sample data typical of the cost of tax-return preparation by services such as H&R Block are shown here.

CH. 1 451	120	230	110	115	160
le	130	150	105	195	155
	105	360	120	120	140
xCost	100	115	180	235	255

- a. Compute the mean, median, and mode.
- b. Compute the first and third quartiles.
- c. Compute and interpret the 90th percentile.
- J. D. Powers and Associates surveyed cell phone users in order to learn about the minutes
  of cell phone usage per month (Associated Press, June 2002). Minutes per month for a
  sample of 15 cell phone users are shown here.

135	395
830	1180
250	420
245	210
380	105
	830 250 245

- a. What is the mean number of minutes of usage per month?
- . What is the median number of minutes of usage per month?
- c. What is the 85th percentile?
- d. J. D. Powers and Associates reported that the average wireless subscriber plan allows up to 750 minutes of usage per month. What do the data suggest about cell phone subscribers' utilization of their monthly plan?
- 10. A panel of economists provided forecasts of the U.S. economy for the first six months of 2007 (*The Wall Street Journal*, January 2, 2007). The percentage changes in gross domestic product (GDP) forecasted by 30 economists are as follows.



Uploaded By: Haneen

SELF test

2.6	3.1	2.3	2.7	3.4	0.9	2.6	2.8	2.0	2.4
2.7	2.7	2.7	2.9	3.1	2.8	1.7	2.3	2.8	3.5
0.4	2.5	2.2	1.9	1.8	1.1	2.0	2.1	2.5	0.5

- a. What is the minimum forecast for the percentage change in GDP? What is the maximum?
- b. Compute the mean, median, and mode.
- c. Compute the first and third quartiles.
- d. Did the economists provide an optimistic or pessimistic outlook for the U.S. economy? Discuss.
- 11. In automobile mileage and gasoline-consumption testing, 13 automobiles were road tested for 300 miles in both city and highway driving conditions. The following data were recorded for miles-per-gallon performance.

City: 16.2 16.7 15.9 14.4 13.2 15.3 16.8 16.0 16.1 15.3 15.2 15.3 16.2 Highway: 19.4 20.6 18.3 18.6 19.2 17.4 17.2 18.6 19.0 21.1 19.4 18.5 18.7

Use the mean, median, and mode to make a statement about the difference in performance for city and highway driving.

12. Walt Disney Company bought Pixar Animation Studios, Inc., in a deal worth \$7.4 billion (CNNMoney.com, January 24, 2006). A list of the animated movies produced by Disney and Pixar during the previous 10 years follows. The box office revenues are in millions of dollars. Compute the total revenue, the mean, the median, and the quartiles to compare the box office success of the movies produced by both companies. Do the statistics suggest at least one of the reasons Disney was interested in buying Pixar? Discuss.

Disney Movies	Revenue (\$millions)	Pixar Movies	Revenue (\$millions)
Pocahontas	346	Toy Story	362
Hunchback of Notre Dame	325	A Bug's Life	363
Hercules	253	Toy Story 2	485
Mulan	304	Monsters, Inc.	525
Tarzan	448	Finding Nemo	865
Dinosaur	354	The Incredibles	631
The Emperor's New Groove	169		
Lilo & Stitch	273		
Treasure Planet	110		
The Jungle Book 2	136		
Brother Bear	250		
Home on the Range	104		
Chicken Little	249		



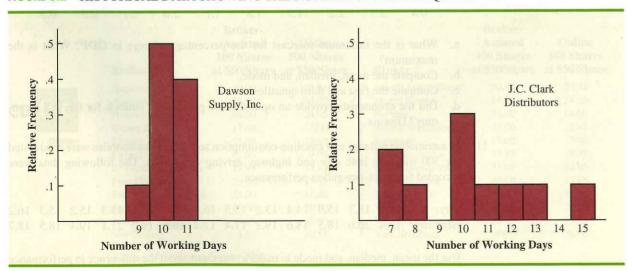
# 3.2

# Measures of Variability

The variability in the delivery time creates uncertainty for production scheduling. Methods in this section help measure and understand variability.

In addition to measures of location, it is often desirable to consider measures of variability, or dispersion. For example, suppose that you are a purchasing agent for a large manufacturing firm and that you regularly place orders with two different suppliers. After several months of operation, you find that the mean number of days required to fill orders is 10 days for both of the suppliers. The histograms summarizing the number of working days required to fill orders from the suppliers are shown in Figure 3.2. Although the mean number of days is 10 for both suppliers, do the two suppliers demonstrate the

FIGURE 3.2 HISTORICAL DATA SHOWING THE NUMBER OF DAYS REQUIRED TO FILL ORDERS



STUDENTS-HUB.com

same degree of reliability in terms of making deliveries on schedule? Note the dispersion, or variability, in delivery times indicated by the histograms. Which supplier would you prefer?

For most firms, receiving materials and supplies on schedule is important. The sevenor eight-day deliveries shown for J.C. Clark Distributors might be viewed favorably; however, a few of the slow 13- to 15-day deliveries could be disastrous in terms of keeping a workforce busy and production on schedule. This example illustrates a situation in which the variability in the delivery times may be an overriding consideration in selecting a supplier. For most purchasing agents, the lower variability shown for Dawson Supply, Inc., would make Dawson the preferred supplier.

We turn now to a discussion of some commonly used measures of variability.

# Range

The simplest measure of variability is the range.

RANGE

Range = Largest value - Smallest value

Let us refer to the data on starting salaries for business school graduates in Table 3.1. The largest starting salary is 3925 and the smallest is 3310. The range is 3925 - 3310 = 615.

Although the range is the easiest of the measures of variability to compute, it is seldom used as the only measure. The reason is that the range is based on only two of the observations and thus is highly influenced by extreme values. Suppose one of the graduates received a starting salary of \$10,000 per month. In this case, the range would be 10,000-3310=6690 rather than 615. This large value for the range would not be especially descriptive of the variability in the data because 11 of the 12 starting salaries are closely grouped between 3310 and 3730.

# Interquartile Range

A measure of variability that overcomes the dependency on extreme values is the **interquartile range (IQR)**. This measure of variability is the difference between the third quartile,  $Q_3$ , and the first quartile,  $Q_1$ . In other words, the interquartile range is the range for the middle 50% of the data.

INTERQUARTILE RANGE  $IQR = Q_3 - Q_1 \tag{3.3}$ 

For the data on monthly starting salaries, the quartiles are  $Q_3 = 3600$  and  $Q_1 = 3465$ . Thus, the interquartile range is 3600 - 3465 = 135.

#### **Variance**

The **variance** is a measure of variability that utilizes all the data. The variance is based on the difference between the value of each observation  $(x_i)$  and the mean. The difference

3.2 Measures of Variability

between each  $x_i$  and the mean ( $\bar{x}$  for a sample,  $\mu$  for a population) is called a *deviation about* the mean. For a sample, a deviation about the mean is written  $(x_i - \bar{x})$ ; for a population, it is written  $(x_i - \mu)$ . In the computation of the variance, the deviations about the mean are squared.

If the data are for a population, the average of the squared deviations is called the *population variance*. The population variance is denoted by the Greek symbol  $\sigma^2$ . For a population of N observations and with  $\mu$  denoting the population mean, the definition of the population variance is as follows.

POPULATION VARIANCE

$$\sigma^2 = \frac{\Sigma (x_i - \mu)^2}{N} \tag{3.4}$$

In most statistical applications, the data being analyzed are for a sample. When we compute a sample variance, we are often interested in using it to estimate the population variance  $\sigma^2$ . Although a detailed explanation is beyond the scope of this text, it can be shown that if the sum of the squared deviations about the sample mean is divided by n-1, and not n, the resulting sample variance provides an unbiased estimate of the population variance. For this reason, the *sample variance*, denoted by  $s^2$ , is defined as follows.

The sample variance  $s^2$  is the estimator of the population variance  $\sigma^2$ .

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$
 (3.5)

To illustrate the computation of the sample variance, we will use the data on class size for the sample of five college classes as presented in Section 3.1. A summary of the data, including the computation of the deviations about the mean and the squared deviations about the mean, is shown in Table 3.3. The sum of squared deviations about the mean is  $\Sigma(x_i - \bar{x})^2 = 256$ . Hence, with n - 1 = 4, the sample variance is

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1} = \frac{256}{4} = 64$$

Note that the units associated with the sample variance often cause confusion. Because the values being summed in the variance calculation,  $(x_i - \bar{x})^2$ , are squared, the units associated with the sample variance are also *squared*. For instance, the sample variance for the class size data is  $s^2 = 64$  (students)<sup>2</sup>. The squared units associated with variance make it difficult to obtain an intuitive understanding and interpretation of the numerical value of the variance. We recommend that you think of the variance as a measure useful in comparing the amount of variability for two or more variables. In a comparison of the variables, the one with the largest variance shows the most variability. Further interpretation of the value of the variance may not be necessary.

The variance is useful in comparing the variability of two or more variables.

STUDENTS-HUB.com

TABLE 3.3 COMPUTATION OF DEVIATIONS AND SQUARED DEVIATIONS ABOUT THE MEAN FOR THE CLASS SIZE DATA

Number of Students in Class (x <sub>i</sub> )	Mean Class Size ( $\bar{x}$ )	Deviation About the Mean $(x_i - \bar{x})$	Squared Deviation About the Mean $(x_i - \bar{x})^2$
46	44	2	4
54	44	10	100
42	44	-2	4
46	44	2	4
32	44	<u>-12</u>	144
	pariesb buttouten	mlugud 0	256
	of the diagraph in case	$\Sigma(x_i - \bar{x})$	$\Sigma (x_i - \bar{x})^2$

As another illustration of computing a sample variance, consider the starting salaries listed in Table 3.1 for the 12 business school graduates. In Section 3.1, we showed that the sample mean starting salary was 3540. The computation of the sample variance  $(s^2 = 27,440.91)$  is shown in Table 3.4.

In Tables 3.3 and 3.4 we show both the sum of the deviations about the mean and the sum of the squared deviations about the mean. For any data set, the sum of the deviations about the mean will *always equal zero*. Note that in Tables 3.3 and 3.4,  $\Sigma(x_i - \bar{x}) = 0$ . The positive deviations and negative deviations cancel each other, causing the sum of the deviations about the mean to equal zero.

TABLE 3.4 COMPUTATION OF THE SAMPLE VARIANCE FOR THE STARTING SALARY DATA

Monthly Salary (x <sub>i</sub> )	Sample Mean $(\bar{x})$	Deviation About the Mean $(x_i - \bar{x})$	Squared Deviation About the Mean $(x_i - \bar{x})^2$
3450	3540	-90	8,100
3550	3540	10	100
3650	3540	110	12,100
3480	3540	-60	3,600
3355	3540	-185	34,225
3310	3540	-230	52,900
3490	3540	-50	2,500
3730	3540	190	36,100
3540	3540	0	0
3925	3540	385	148,225
3520	3540	-20	400
3480	3540	<u>-60</u>	3,600
	- Programme and the second	0,	301,850
	moneyelk krabilet	$\Sigma(x_i-\bar{x})$	$\sum (x_i - \bar{x})^2$
Using equation (	(3.5),		
	$\Sigma(x)$	$\frac{(x^2 - \bar{x})^2}{(x^2 - 1)^2} = \frac{301,850}{11} = 27,440.91$	

## Standard Deviation

The standard deviation is defined to be the positive square root of the variance. Following the notation we adopted for a sample variance and a population variance, we use s to denote the sample standard deviation and  $\sigma$  to denote the population standard deviation. The standard deviation is derived from the variance in the following way.

The sample standard deviation s is the estimator of the population standard deviation  $\sigma$ .

The standard deviation is easier to interpret than the

measured in the same units

variance because the

standard deviation is

as the data.

STANDARD DEVIATION

Sample standard deviation = 
$$s = \sqrt{s^2}$$
 (3.6)

Population standard deviation = 
$$\sigma = \sqrt{\sigma^2}$$
 (3.7)

Recall that the sample variance for the sample of class sizes in five college classes is  $s^2 = 64$ . Thus, the sample standard deviation is  $s = \sqrt{64} = 8$ . For the data on starting salaries, the sample standard deviation is  $s = \sqrt{27,440.91} = 165.65$ .

What is gained by converting the variance to its corresponding standard deviation? Recall that the units associated with the variance are squared. For example, the sample variance for the starting salary data of business school graduates is  $s^2 = 27,440.91$  (dollars)<sup>2</sup>. Because the standard deviation is the square root of the variance, the units of the variance, dollars squared, are converted to dollars in the standard deviation. Thus, the standard deviation of the starting salary data is \$165.65. In other words, the standard deviation is measured in the same units as the original data. For this reason the standard deviation is more easily compared to the mean and other statistics that are measured in the same units as the original data.

Coefficient of Variation

In some situations we may be interested in a descriptive statistic that indicates how large the standard deviation is relative to the mean. This measure is called the coefficient of variation and is usually expressed as a percentage.

The coefficient of variation is a relative measure of variability; it measures the standard deviation relative to the mean.

COEFFICIENT OF VARIATION 
$$\left( \frac{\text{Standard deviation}}{\text{Mean}} \times 100 \right) \%$$
 (3.8)

For the class size data, we found a sample mean of 44 and a sample standard deviation of 8. The coefficient of variation is  $[(8/44) \times 100]\% = 18.2\%$ . In words, the coefficient of variation tells us that the sample standard deviation is 18.2% of the value of the sample mean. For the starting salary data with a sample mean of 3540 and a sample standard deviation of 165.65, the coefficient of variation,  $[(165.65/3540) \times 100]\% = 4.7\%$ , tells us the sample standard deviation is only 4.7% of the value of the sample mean. In general, the coefficient of variation is a useful statistic for comparing the variability of variables that have different standard deviations and different means.

STUDENTS-HUB.com

#### **NOTES AND COMMENTS**

- 1. Statistical software packages and spreadsheets can be used to develop the descriptive statistics presented in this chapter. After the data are entered into a worksheet, a few simple commands can be used to generate the desired output. In Appendixes 3.1 and 3.2, we show how Minitab and Excel can be used to develop descriptive statistics.
- 2. The standard deviation is a commonly used measure of the risk associated with investing in stock and stock funds (BusinessWeek, January 17, 2000). It provides a measure of how monthly returns fluctuate around the long-run average return.
- 3. Rounding the value of the sample mean  $\bar{x}$  and the values of the squared deviations  $(x_i - \bar{x})^2$

may introduce errors when a calculator is used in the computation of the variance and standard deviation. To reduce rounding errors, we recommend carrying at least six significant digits during intermediate calculations. The resulting variance or standard deviation can then be rounded to fewer digits.

4. An alternative formula for the computation of the sample variance is

$$s^2 = \frac{\sum x_i^2 - n\bar{x}^2}{n-1}$$

where 
$$\sum x_i^2 = x_1^2 + x_2^2 + \cdots + x_n^2$$
.

#### **Exercises**

#### Methods

- 13. Consider a sample with data values of 10, 20, 12, 17, and 16. Compute the range and inter-
- 14. Consider a sample with data values of 10, 20, 12, 17, and 16. Compute the variance and standard deviation.
- 15. Consider a sample with data values of 27, 25, 20, 15, 30, 34, 28, and 25. Compute the range, interquartile range, variance, and standard deviation.

# **SELF** test

SELF test

# **Applications**

- 16. A bowler's scores for six games were 182, 168, 184, 190, 170, and 174. Using these data as a sample, compute the following descriptive statistics.
  - a. Range
- c. Standard deviation
- b. Variance
- d. Coefficient of variation
- 17. A home theater in a box is the easiest and cheapest way to provide surround sound for a home entertainment center. A sample of prices is shown here (Consumer Reports Buying Guide, 2004). The prices are for models with a DVD player and for models without a DVD player.

Models with DVD Player	Price	Models without DVD Player	Price
Sony HT-1800DP	\$450	Pioneer HTP-230	\$300
Pioneer HTD-330DV	300	Sony HT-DDW750	300
Sony HT-C800DP	400	Kenwood HTB-306	360
Panasonic SC-HT900	500	RCA RT-2600	290
Panasonic SC-MTI	400	Kenwood HTB-206	300

- a. Compute the mean price for models with a DVD player and the mean price for models without a DVD player. What is the additional price paid to have a DVD player included in a home theater unit?
- b. Compute the range, variance, and standard deviation for the two samples. What does this information tell you about the prices for models with and without a DVD player?

3.3 Measures of Distribution Shape, Relative Location, and Detecting Outliers

97

18. Car rental rates per day for a sample of seven Eastern U.S. cities are as follows (*The Wall Street Journal*, January 16, 2004).

City	Daily Rate	
Boston	\$43	
Atlanta	35	
Miami	34	
New York	58	
Orlando	30	
Pittsburgh	30	
Washington, D.C.	36	

- a. Compute the mean, variance, and standard deviation for the car rental rates.
- b. A similar sample of seven Western U.S. cities showed a sample mean car rental rate of \$38 per day. The variance and standard deviation were 12.3 and 3.5, respectively. Discuss any difference between the car rental rates in Eastern and Western U.S. cities.
- 19. The *Los Angeles Times* regularly reports the air quality index for various areas of Southern California. A sample of air quality index values for Pomona provided the following data: 28, 42, 58, 48, 45, 55, 60, 49, and 50.
  - a. Compute the range and interquartile range.
  - b. Compute the sample variance and sample standard deviation.
  - c. A sample of air quality index readings for Anaheim provided a sample mean of 48.5, a sample variance of 136, and a sample standard deviation of 11.66. What comparisons can you make between the air quality in Pomona and that in Anaheim on the basis of these descriptive statistics?
- 20. The following data were used to construct the histograms of the number of days required to fill orders for Dawson Supply, Inc., and J.C. Clark Distributors (see Figure 3.2).

 Dawson Supply Days for Delivery:
 11
 10
 9
 10
 11
 11
 10
 11
 10
 10

 Clark Distributors Days for Delivery:
 8
 10
 13
 7
 10
 11
 10
 7
 15
 12

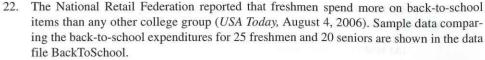
Use the range and standard deviation to support the previous observation that Dawson Supply provides the more consistent and reliable delivery times.

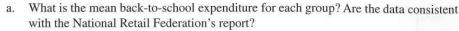
21. How do grocery costs compare across the country? Using a market basket of 10 items including meat, milk, bread, eggs, coffee, potatoes, cereal, and orange juice, *Where to Retire* magazine calculated the cost of the market basket in six cities and in six retirement areas across the country (*Where to Retire*, November/December 2003). The data with market basket cost to the nearest dollar are as follows:

City	Cost	Retirement Area	Cost	
Buffalo, NY	\$33	Biloxi-Gulfport, MS	\$29	
Des Moines, IA	27	Asheville, NC	32	
Hartford, CT	32	Flagstaff, AZ	32	
Los Angeles, CA	38	Hilton Head, SC	34	
Miami, FL	36	Fort Myers, FL	34	
Pittsburgh, PA	32	Santa Fe, NM	31	
		PITCH HARMAN STATE OF THE PARTY		

- Compute the mean, variance, and standard deviation for the sample of cities and the sample of retirement areas.
- b. What observations can be made based on the two samples?

STUDENTS-HUB.com





- b. What is the range for the expenditures in each group?
- c. What is the interquartile range for the expenditures in each group?
- d. What is the standard deviation for expenditures in each group?
- e. Do freshmen or seniors have more variation in back-to-school expenditures?
- 23. Scores turned in by an amateur golfer at the Bonita Fairways Golf Course in Bonita Springs, Florida, during 2005 and 2006 are as follows:

2005 Season	74	78	79	77	75	73	75	77
2006 Season	71	70	75	77	85	80	71	79

- a. Use the mean and standard deviation to evaluate the golfer's performance over the two-year period.
- b. What is the primary difference in performance between 2005 and 2006? What improvement, if any, can be seen in the 2006 scores?
- 24. The following times were recorded by the quarter-mile and mile runners of a university track team (times are in minutes).

Quarter-Mile Times:	.92	.98	1.04	.90	.99
Mile Times:	4.52	4.35	4.60	4.70	4.50

After viewing this sample of running times, one of the coaches commented that the quartermilers turned in the more consistent times. Use the standard deviation and the coefficient of variation to summarize the variability in the data. Does the use of the coefficient of variation indicate that the coach's statement should be qualified?



coffile

BackToSchool

# Measures of Distribution Shape, Relative Location, and Detecting Outliers

We have described several measures of location and variability for data. In addition, it is often important to have a measure of the shape of a distribution. In Chapter 2 we noted that a histogram provides a graphical display showing the shape of a distribution. An important numerical measure of the shape of a distribution is called **skewness**.

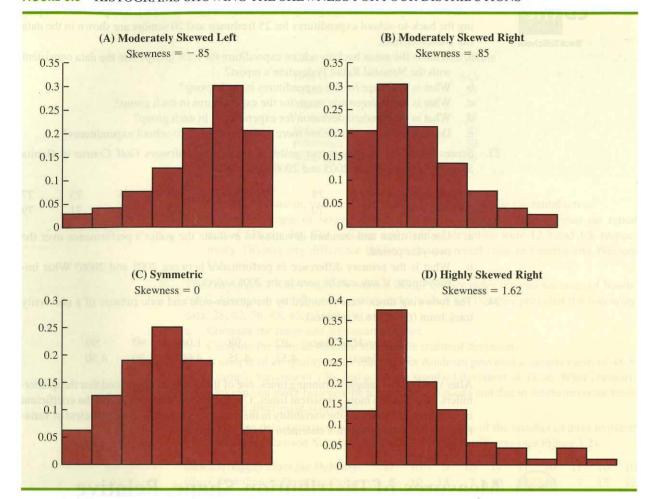
# **Distribution Shape**

Shown in Figure 3.3 are four histograms constructed from relative frequency distributions. The histograms in panels A and B are moderately skewed. The one in panel A is skewed to the left; its skewness is -.85. The histogram in panel B is skewed to the right; its skewness is +.85. The histogram in panel C is symmetric; its skewness is zero. The histogram in panel D is highly skewed to the right; its skewness is 1.62. The formula used to compute skewness is somewhat complex.\* However, the skewness can be easily computed using

Skewness = 
$$\frac{n}{(n-1)(n-2)} \sum \left(\frac{x_i - \bar{x}}{s}\right)^3$$

<sup>\*</sup>The formula for the skewness of sample data:

FIGURE 3.3 HISTOGRAMS SHOWING THE SKEWNESS FOR FOUR DISTRIBUTIONS



statistical software (see Appendixes 3.1 and 3.2). For data skewed to the left, the skewness is negative; for data skewed to the right, the skewness is positive. If the data are symmetric, the skewness is zero.

For a symmetric distribution, the mean and the median are equal. When the data are positively skewed, the mean will usually be greater than the median; when the data are negatively skewed, the mean will usually be less than the median. The data used to construct the histogram in panel D are customer purchases at a women's apparel store. The mean purchase amount is \$77.60 and the median purchase amount is \$59.70. The relatively few large purchase amounts tend to increase the mean, while the median remains unaffected by the large purchase amounts. The median provides the preferred measure of location when the data are highly skewed.

#### z-Scores

In addition to measures of location, variability, and shape, we are also interested in the relative location of values within a data set. Measures of relative location help us determine how far a particular value is from the mean.

By using both the mean and standard deviation, we can determine the relative location of any observation. Suppose we have a sample of *n* observations, with the values denoted

STUDENTS-HUB.com

by  $x_1, x_2, \ldots, x_n$ . In addition, assume that the sample mean,  $\bar{x}$ , and the sample standard deviation, s, are already computed. Associated with each value,  $x_i$ , is another value called its **z-score**. Equation (3.9) shows how the z-score is computed for each  $x_i$ .

z-SCORE  $z_i = \frac{x_i - \bar{x}}{s} \tag{3.9}$  where  $z_i = \text{the } z\text{-score for } x_i$   $\bar{x} = \text{the sample mean}$  s = the sample standard deviation

The z-score is often called the *standardized value*. The z-score,  $z_i$ , can be interpreted as the *number of standard deviations*  $x_i$  *is from the mean*  $\bar{x}$ . For example,  $z_1 = 1.2$  would indicate that  $x_1$  is 1.2 standard deviations greater than the sample mean. Similarly,  $z_2 = -.5$  would indicate that  $x_2$  is .5, or 1/2, standard deviation less than the sample mean. A z-score greater than zero occurs for observations with a value greater than the mean, and a z-score less than zero occurs for observations with a value less than the mean. A z-score of zero indicates that the value of the observation is equal to the mean.

The z-score for any observation can be interpreted as a measure of the relative location of the observation in a data set. Thus, observations in two different data sets with the same z-score can be said to have the same relative location in terms of being the same number of standard deviations from the mean.

The z-scores for the class size data are computed in Table 3.5. Recall the previously computed sample mean,  $\bar{x} = 44$ , and sample standard deviation, s = 8. The z-score of -1.50 for the fifth observation shows it is farthest from the mean; it is 1.50 standard deviations below the mean.

# **Chebyshev's Theorem**

**Chebyshev's theorem** enables us to make statements about the proportion of data values that must be within a specified number of standard deviations of the mean.

TABLE 3.5 z-SCORES FOR THE CLASS SIZE DATA

Number of Students in Class $(x_i)$	Deviation About the Mean $(x_i - \bar{x})$	$\frac{z\text{-Score}}{\left(\frac{x_i - \bar{x}}{s}\right)}$
46	2	2/8 = .25
54	10	10/8 = 1.25
42	-2	-2/8 =25
46	2	2/8 = .25
32	-12	-12/8 = -1.50

#### CHEBYSHEV'S THEOREM

At least  $(1 - 1/z^2)$  of the data values must be within z standard deviations of the mean, where z is any value greater than 1.

Some of the implications of this theorem, with z = 2, 3, and 4 standard deviations, follows

- At least .75, or 75%, of the data values must be within z = 2 standard deviations of the mean.
- At least .89, or 89%, of the data values must be within z = 3 standard deviations of the mean.
- At least .94, or 94%, of the data values must be within z = 4 standard deviations of the mean.

For an example using Chebyshev's theorem, suppose that the midterm test scores for 100 students in a college business statistics course had a mean of 70 and a standard deviation of 5. How many students had test scores between 60 and 80? How many students had test scores between 58 and 82?

For the test scores between 60 and 80, we note that 60 is two standard deviations below the mean and 80 is two standard deviations above the mean. Using Chebyshev's theorem, we see that at least .75, or at least 75%, of the observations must have values within two standard deviations of the mean. Thus, at least 75% of the students must have scored between 60 and 80.

For the test scores between 58 and 82, we see that (58 - 70)/5 = -2.4 indicates 58 is 2.4 standard deviations below the mean and that (82 - 70)/5 = +2.4 indicates 82 is 2.4 standard deviations above the mean. Applying Chebyshev's theorem with z = 2.4, we have

$$\left(1 - \frac{1}{z^2}\right) = \left(1 - \frac{1}{(2.4)^2}\right) = .826$$

At least 82.6% of the students must have test scores between 58 and 82.

# **Empirical Rule**

One of the advantages of Chebyshev's theorem is that it applies to any data set regardless of the shape of the distribution of the data. Indeed, it could be used with any of the distributions in Figure 3.3. In many practical applications, however, data sets exhibit a symmetric mound-shaped or bell-shaped distribution like the one shown in Figure 3.4. When the data are believed to approximate this distribution, the **empirical rule** can be used to determine the percentage of data values that must be within a specified number of standard deviations of the mean.

The empirical rule is based on the normal probability distribution, which will be discussed in Chapter 6. The normal distribution is used extensively throughout the text.

Chebyshev's theorem

not be an integer.

requires z > 1; but z need

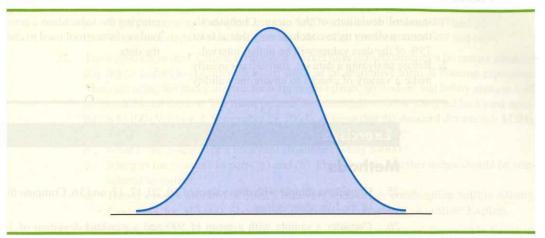
#### EMPIRICAL RULE

For data having a bell-shaped distribution:

- Approximately 68% of the data values will be within one standard deviation of the mean.
- Approximately 95% of the data values will be within two standard deviations of the mean.
- Almost all of the data values will be within three standard deviations of the mean.

# STUDENTS-HUB.com

FIGURE 3.4 A SYMMETRIC MOUND-SHAPED OR BELL-SHAPED DISTRIBUTION



For example, liquid detergent cartons are filled automatically on a production line. Filling weights frequently have a bell-shaped distribution. If the mean filling weight is 16 ounces and the standard deviation is .25 ounces, we can use the empirical rule to draw the following conclusions.

- Approximately 68% of the filled cartons will have weights between 15.75 and 16.25 ounces (within one standard deviation of the mean).
- Approximately 95% of the filled cartons will have weights between 15.50 and 16.50 ounces (within two standard deviations of the mean).
- Almost all filled cartons will have weights between 15.25 and 16.75 ounces (within three standard deviations of the mean).

# **Detecting Outliers**

Sometimes a data set will have one or more observations with unusually large or unusually small values. These extreme values are called **outliers**. Experienced statisticians take steps to identify outliers and then review each one carefully. An outlier may be a data value that has been incorrectly recorded. If so, it can be corrected before further analysis. An outlier may also be from an observation that was incorrectly included in the data set; if so, it can be removed. Finally, an outlier may be an unusual data value that has been recorded correctly and belongs in the data set. In such cases it should remain.

Standardized values (z-scores) can be used to identify outliers. Recall that the empirical rule allows us to conclude that for data with a bell-shaped distribution, almost all the data values will be within three standard deviations of the mean. Hence, in using z-scores to identify outliers, we recommend treating any data value with a z-score less than -3 or greater than +3 as an outlier. Such data values can then be reviewed for accuracy and to determine whether they belong in the data set.

Refer to the z-scores for the class size data in Table 3.5. The z-score of -1.50 shows the fifth class size is farthest from the mean. However, this standardized value is well within the -3 to +3 guideline for outliers. Thus, the z-scores do not indicate that outliers are present in the class size data.

# NOTES AND COMMENTS

 Chebyshev's theorem is applicable for any data set and can be used to state the minimum number of data values that will be within a certain number of standard deviations of the mean. If the data are known to be approximately bell-shaped, more can be said. For instance, the

(continued)

Uploaded By: Haneen

It is a good idea to check

for outliers before making

decisions based on data

analysis. Errors are often

and entering data into the

computer. Outliers should

not necessarily be deleted, but their accuracy and

appropriateness should

be verified.

made in recording data

- empirical rule allows us to say that *approximately* 95% of the data values will be within two standard deviations of the mean; Chebyshev's theorem allows us to conclude only that at least 75% of the data values will be in that interval.
- 2. Before analyzing a data set, statisticians usually make a variety of checks to ensure the validity

of data. In a large study it is not uncommon for errors to be made in recording data values or in entering the values into a computer. Identifying outliers is one tool used to check the validity of the data.

#### **Exercises**

#### Methods

SELF LES

SELF LES

- 25. Consider a sample with data values of 10, 20, 12, 17, and 16. Compute the *z*-score for each of the five observations.
- 26. Consider a sample with a mean of 500 and a standard deviation of 100. What are the *z*-scores for the following data values: 520, 650, 500, 450, and 280?
- 27. Consider a sample with a mean of 30 and a standard deviation of 5. Use Chebyshev's theorem to determine the percentage of the data within each of the following ranges.
  - a. 20 to 40
  - b. 15 to 45
  - c. 22 to 38
  - d. 18 to 42
  - e. 12 to 48
- 28. Suppose the data have a bell-shaped distribution with a mean of 30 and a standard deviation of 5. Use the empirical rule to determine the percentage of data within each of the following ranges.
  - a. 20 to 40
  - b. 15 to 45
  - c. 25 to 35

# **Applications**

- 29. The results of a national survey showed that on average, adults sleep 6.9 hours per night. Suppose that the standard deviation is 1.2 hours.
  - Use Chebyshev's theorem to calculate the percentage of individuals who sleep between 4.5 and 9.3 hours.
  - Use Chebyshev's theorem to calculate the percentage of individuals who sleep between 3.9 and 9.9 hours.
  - c. Assume that the number of hours of sleep follows a bell-shaped distribution. Use the empirical rule to calculate the percentage of individuals who sleep between 4.5 and 9.3 hours per day. How does this result compare to the value that you obtained using Chebyshev's theorem in part (a)?
- 30. The Energy Information Administration reported that the mean retail price per gallon of regular grade gasoline was \$2.30 (Energy Information Administration, February 27, 2006). Suppose that the standard deviation was \$.10 and that the retail price per gallon has a bell-shaped distribution.
  - a. What percentage of regular grade gasoline sold between \$2.20 and \$2.40 per gallon?
  - b. What percentage of regular grade gasoline sold between \$2.20 and \$2.50 per gallon?
  - c. What percentage of regular grade gasoline sold for more than \$2.50 per gallon?
- 31. The national average for the verbal portion of the College Board's Scholastic Aptitude Test (SAT) is 507 (*The World Almanac*, 2006). The College Board periodically rescales the test scores such that the standard deviation is approximately 100. Answer the following questions using a bell-shaped distribution and the empirical rule for the verbal test scores.

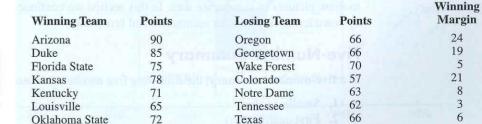
STUDENTS-HUB.com

- a. What percentage of students have an SAT verbal score greater than 607?
- b. What percentage of students have an SAT verbal score greater than 707?
- c. What percentage of students have an SAT verbal score between 407 and 507?
- d. What percentage of students have an SAT verbal score between 307 and 607?
- 32. The high costs in the California real estate market have caused families who cannot afford to buy bigger homes to consider backyard sheds as an alternative form of housing expansion. Many are using the backyard structures for home offices, art studios, and hobby areas as well as for additional storage. The mean price of a customized wooden, shingled backyard structure is \$3100 (Newsweek, September 29, 2003). Assume that the standard deviation is \$1200.
  - a. What is the z-score for a backyard structure costing \$2300?
  - b. What is the z-score for a backyard structure costing \$4900?
  - c. Interpret the *z*-scores in parts (a) and (b). Comment on whether either should be considered an outlier.
  - d. The *Newsweek* article described a backyard shed-office combination built in Albany, California, for \$13,000. Should this structure be considered an outlier? Explain.
- 33. Florida Power & Light (FP&L) Company has enjoyed a reputation for quickly fixing its electric system after storms. However, during the hurricane seasons of 2004 and 2005, a new reality was that the company's historical approach to emergency electric system repairs was no longer good enough (*The Wall Street Journal*, January 16, 2006). Data showing the days required to restore electric service after seven hurricanes during 2004 and 2005 follow.

Hurricane	Days to Restore Service
Charley	13
Frances	12
Jeanne	8
Dennis	3
Katrina	8
Rita	2
Wilma	18

Based on this sample of seven, compute the following descriptive statistics:

- a. Mean, median, and mode
- b. Range and standard deviation
- c. Should Wilma be considered an outlier in terms of the days required to restore electric service?
- d. The seven hurricanes resulted in 10 million service interruptions to customers. Do the statistics show that FP&L should consider updating its approach to emergency electric system repairs? Discuss.
- A sample of 10 NCAA college basketball game scores provided the following data (USA Today, January 26, 2004).







	- 1	ř		
4	Evn	oratory	Data	Analysis
. 4		ioi dioi y	Dulu	

Winning Team	Points	Losing Team	Points	Winning Margin
Purdue	76	Michigan State	70	6
Stanford	77	Southern Cal	67	10
Wisconsin	76	Illinois	56	20

- a. Compute the mean and standard deviation for the points scored by the winning team.
- b. Assume that the points scored by the winning teams for all NCAA games follow a bell-shaped distribution. Using the mean and standard deviation found in part (a), estimate the percentage of all NCAA games in which the winning team scores 84 or more points. Estimate the percentage of NCAA games in which the winning team scores more than 90 points.
- c. Compute the mean and standard deviation for the winning margin. Do the data contain outliers? Explain.
- 35. *Consumer Review* posts reviews and ratings of a variety of products on the Internet. The following is a sample of 20 speaker systems and their ratings (http://www.audioreview.com). The ratings are on a scale of 1 to 5, with 5 being best.

CD	file
	Speakers

Speaker	Rating	Speaker	Rating
Infinity Kappa 6.1	4.00	ACI Sapphire III	4.67
Allison One	4.12	Bose 501 Series	2.14
Cambridge Ensemble II	3.82	DCM KX-212	4.09
Dynaudio Contour 1.3	4.00	Eosone RSF1000	4.17
Hsu Rsch. HRSW12V	4.56	Joseph Audio RM7si	4.88
Legacy Audio Focus	4.32	Martin Logan Aerius	4.26
Mission 73li	4.33	Omni Audio SA 12.3	2.32
PSB 400i	4.50	Polk Audio RT12	4.50
Snell Acoustics D IV	4.64	Sunfire True Subwoofer	4.17
Thiel CS1.5	4.20	Yamaha NS-A636	2.17

- a. Compute the mean and the median.
- b. Compute the first and third quartiles.
- Compute the standard deviation.
- d. The skewness of this data is -1.67. Comment on the shape of the distribution.
- e. What are the *z*-scores associated with Allison One and Omni Audio?
- f. Do the data contain any outliers? Explain.



# **Exploratory Data Analysis**

In Chapter 2 we introduced the stem-and-leaf display as a technique of exploratory data analysis. Recall that exploratory data analysis enables us to use simple arithmetic and easy-to-draw pictures to summarize data. In this section we continue exploratory data analysis by considering five-number summaries and box plots.

# **Five-Number Summary**

In a five-number summary, the following five numbers are used to summarize the data.

- 1. Smallest value
- 2. First quartile  $(Q_1)$
- 3. Median  $(Q_2)$

- **4.** Third quartile  $(Q_3)$
- 5. Largest value

The easiest way to develop a five-number summary is to first place the data in ascending order. Then it is easy to identify the smallest value, the three quartiles, and the largest value. The monthly starting salaries shown in Table 3.1 for a sample of 12 business school graduates are repeated here in ascending order.

3310 3355 3450 3480 3480 3490 3520 3540 3550 3650 3730 3925 
$$Q_1 = 3465$$
  $Q_2 = 3505$   $Q_3 = 3600$  (Median)

The median of 3505 and the quartiles  $Q_1 = 3465$  and  $Q_3 = 3600$  were computed in Section 3.1. Reviewing the data shows a smallest value of 3310 and a largest value of 3925. Thus the five-number summary for the salary data is 3310, 3465, 3505, 3600, 3925. Approximately one-fourth, or 25%, of the observations are between adjacent numbers in a five-number summary.

## **Box Plot**

Box plots provide another way to identify outliers. But

they do not necessarily

identify the same values

as those with a z-score

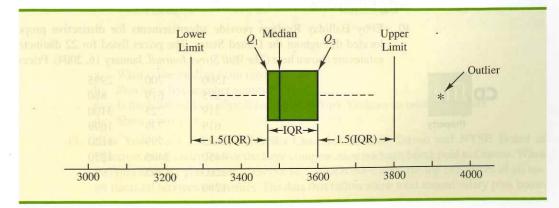
less than -3 or greater

than +3. Either or both procedures may be used.

A **box plot** is a graphical summary of data that is based on a five-number summary. A key to the development of a box plot is the computation of the median and the quartiles,  $Q_1$  and  $Q_3$ . The interquartile range,  $IQR = Q_3 - Q_1$ , is also used. Figure 3.5 is the box plot for the monthly starting salary data. The steps used to construct the box plot follow.

- 1. A box is drawn with the ends of the box located at the first and third quartiles. For the salary data,  $Q_1 = 3465$  and  $Q_3 = 3600$ . This box contains the middle 50% of the data.
- 2. A vertical line is drawn in the box at the location of the median (3505 for the salary data).
- 3. By using the interquartile range, IQR =  $Q_3 Q_1$ , *limits* are located. The limits for the box plot are 1.5(IQR) below  $Q_1$  and 1.5(IQR) above  $Q_3$ . For the salary data, IQR =  $Q_3 Q_1 = 3600 3465 = 135$ . Thus, the limits are 3465 1.5(135) = 3262.5 and 3600 + 1.5(135) = 3802.5. Data outside these limits are considered *outliers*.
- **4.** The dashed lines in Figure 3.5 are called *whiskers*. The whiskers are drawn from the ends of the box to the smallest and largest values *inside the limits* computed in step 3. Thus, the whiskers end at salary values of 3310 and 3730.
- **5.** Finally, the location of each outlier is shown with the symbol \*. In Figure 3.5 we see one outlier, 3925.

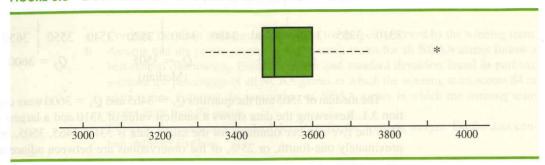
FIGURE 3.5 BOX PLOT OF THE STARTING SALARY DATA WITH LINES SHOWING THE LOWER AND UPPER LIMITS



STUDENTS-HUB.com

In Figure 3.5 we included lines showing the location of the upper and lower limits. These lines were drawn to show how the limits are computed and where they are located for the salary data. Although the limits are always computed, generally they are not drawn on the box plots. Figure 3.6 shows the usual appearance of a box plot for the salary data.

FIGURE 3.6 BOX PLOT OF THE STARTING SALARY DATA



#### **NOTES AND COMMENTS**

- An advantage of the exploratory data analysis
  procedures is that they are easy to use; few numerical calculations are necessary. We simply
  sort the data values into ascending order and
  identify the five-number summary. The box plot
  can then be constructed. It is not necessary to
- compute the mean and the standard deviation for the data.
- 2. In Appendix 3.1, we show how to construct a box plot for the starting salary data using Minitab. The box plot obtained looks just like the one in Figure 3.6, but turned on its side.

#### Exercises

## Methods

- 36. Consider a sample with data values of 27, 25, 20, 15, 30, 34, 28, and 25. Provide the five-number summary for the data.
- 37. Show the box plot for the data in exercise 36.
- 38. Show the five-number summary and the box plot for the following data: 5, 15, 18, 10, 8, 12, 16, 10, 6.
- 39. A data set has a first quartile of 42 and a third quartile of 50. Compute the lower and upper limits for the corresponding box plot. Should a data value of 65 be considered an outlier?

# **Applications**

40. Ebby Halliday Realtors provide advertisements for distinctive properties and estates located throughout the United States. The prices listed for 22 distinctive properties and estates are shown here (*The Wall Street Journal*, January 16, 2004). Prices are in thousands.



SELF test

1500	700	2995
895	619	880
719	725	3100
619	739	1699
625	799	1120
4450	2495	1250
2200	1395	912
1280		

STUDENTS-HUB.com

- a. Provide a five-number summary.
- b. Compute the lower and upper limits.
- c. The highest priced property, \$4,450,000, is listed as an estate overlooking White Rock Lake in Dallas, Texas. Should this property be considered an outlier? Explain.
- d. Should the second highest priced property, listed for \$3,100,000, be considered an outlier? Explain.
- e. Show a box plot.
- 41. Annual sales, in millions of dollars, for 21 pharmaceutical companies follow.

8408	1374	1872	8879	2459	11413	
608	14138	6452	1850	2818	1356	
10498	7478	4019	4341	739	2127	
3653	5794	8305				

- a. Provide a five-number summary.
- b. Compute the lower and upper limits.
- c. Do the data contain any outliers?
- d. Johnson & Johnson's sales are the largest on the list at \$14,138 million. Suppose a data entry error (a transposition) had been made and the sales had been entered as \$41,138 million. Would the method of detecting outliers in part (c) identify this problem and allow for correction of the data entry error?
- e. Show a box plot.
- 42. Major League Baseball payrolls continue to escalate. Team payrolls in millions are as follows (*USA Today* Online Database, March 2006).



-	<b>Team</b>	Payroll	Team	Payroll
-1	Arizona	\$ 62	Milwaukee	\$ 40
1	Atlanta	86	Minnesota	56
J	Baltimore	74	NY Mets	101
I	Boston	124	NY Yankees	208
(	Chi Cubs	87	Oakland	55
(	Chi White Sox	75	Philadelphia	96
(	Cincinnati	62	Pittsburgh	38
(	Cleveland	42	San Diego	63
(	Colorado	48	San Francisco	90
I	Detroit	69	Seattle	88
I	Florida	60	St. Louis	92
I	Houston	77	Tampa Bay	30
I	Kansas City	37	Texas	56
I	_A Angels	98	Toronto	46
	A Dodgers	83	Washington	49

- a. What is the median team payroll?
- b. Provide a five-number summary.
- c. Is the \$208 million payroll for the New York Yankees an outlier? Explain.
- d. Show a box plot.
- 43. New York Stock Exchange (NYSE) Chairman Richard Grasso and NYSE Board of Directors came under fire for the large compensation package being paid to Grasso. When it comes to salary plus bonus, Grasso's \$8.5 million out-earned the top executives of all major financial services companies. The data that follow show total annual salary plus bonus

paid to the top executives of 14 financial services companies (The Wall Street Journal, September 17, 2003). Data are in millions.

Company Aetna AIG Allstate American Express Chubb Cigna Citigroup	\$3.5 6.0 4.1 3.8 2.1 1.0	Company Fannie Mae Federal Home Loan Fleet Boston Freddie Mac Mellon Financial Merrill Lynch Wells Fargo	\$4.3 0.8 1.0 1.2 2.0 7.7 8.0
---	--	--	---

- a. What is the median annual salary plus bonus paid to the top executive of the 14 financial service companies?
- b. Provide a five-number summary.
- Should Grasso's \$8.5 million annual salary plus bonus be considered an outlier for this group of top executives? Explain.
- d. Show a box plot.
- A listing of 46 mutual funds and their 12-month total return percentage is shown in Table 3.6 (Smart Money, February 2004).
  - What are the mean and median return percentages for these mutual funds?
  - What are the first and third quartiles?
  - Provide a five-number summary.
  - Do the data contain any outliers? Show a box plot.

# TABLE 3.6 TWELVE-MONTH RETURN FOR MUTUAL FUNDS

	Return (%)	Mutual Fund	Return (%)
Mutual Fund		Nations Small Company	21.4
Alger Capital Appreciation	23.5	Nations SmallCap Index	24.5
Alger LargeCap Growth	22.8	Nations Strategic Growth	10.4
Alger MidCap Growth	38.3	Nations Value Inv	10.8
Alger SmallCap	41.3	One Group Diversified Equity	10.0
AllianceBernstein Technology	40.6	One Group Diversified Int'l	10.9
Federated American Leaders	15.6	One Group Diversified Mid Cap	15.1
Federated Capital Appreciation	12.4	One Group Equity Income	6.6
Federated Equity-Income	11.5	One Group Int'l Equity Index	13.2
Federated Kaufmann	33.3	One Group Large Cap Growth	13.6
Federated Max-Cap Index	16.0	One Group Large Cap Value	12.8
Federated Stock	16.9	One Group Mid Cap Growth	18.7
Janus Adviser Int'l Growth	10.3	One Group Mid Cap Value	11.4
Janus Adviser Worldwide	3.4	One Group Small Cap Growth	23.6
Janus Enterprise	24.2	PBHG Growth	27.3
Janus High-Yield	12.1	Putnam Europe Equity	20.4
Janus Mercury	20.6	Putnam Int'l Capital Opportunity	36.0
Janus Overseas	11.9	Putnam International Equity	21.:
Janus Worldwide	4.1	Putnam Int'l New Opportunity	26.
Nations Convertible Securities	13.6	Strong Advisor Mid Cap Growth	23.
Nations Int'l Equity	10.7	Strong Growth 20	11.
Nations LargeCap Enhd. Core	13.2	Strong Growth Inv	23.
Nations LargeCap Index	13.5 19.5	Strong Large Cap Growth	14.
Nation MidCap Index	19.5	Conta mellular establish	

# STUDENTS-HUB.com



Thus far we have examined numerical methods used to summarize the data for one variable at a time. Often a manager or decision maker is interested in the relationship between two variables. In this section we present covariance and correlation as descriptive measures of the relationship between two variables.

We begin by reconsidering the application concerning a stereo and sound equipment store in San Francisco as presented in Section 2.4. The store's manager wants to determine the relationship between the number of weekend television commercials shown and the sales at the store during the following week. Sample data with sales expressed in hundreds of dollars are provided in Table 3.7. It shows 10 observations (n = 10), one for each week. The scatter diagram in Figure 3.7 shows a positive relationship, with higher sales (y) associated with a greater number of commercials (x). In fact, the scatter diagram suggests that a straight line could be used as an approximation of the relationship. In the following discussion, we introduce covariance as a descriptive measure of the linear association between two variables.

#### Covariance

For a sample of size n with the observations  $(x_1, y_1), (x_2, y_2)$ , and so on, the sample covariance is defined as follows:

SAMPLE COVARIANCE 
$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$
 (3.10)

This formula pairs each  $x_i$  with a  $y_i$ . We then sum the products obtained by multiplying the deviation of each  $x_i$  from its sample mean  $\bar{x}$  by the deviation of the corresponding  $y_i$  from its sample mean  $\bar{y}$ ; this sum is then divided by n-1.

TABLE 3.7 SAMPLE DATA FOR THE STEREO AND SOUND EQUIPMENT STORE



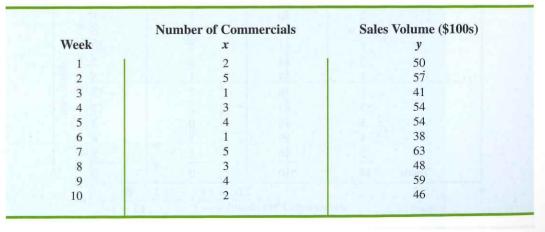
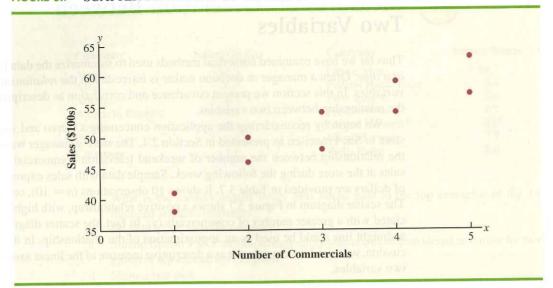


FIGURE 3.7 SCATTER DIAGRAM FOR THE STEREO AND SOUND EQUIPMENT STORE



To measure the strength of the linear relationship between the number of commercials x and the sales volume y in the stereo and sound equipment store problem, we use equation (3.10) to compute the sample covariance. The calculations in Table 3.8 show the computation of  $\sum (x_i - \bar{x})(y_i - \bar{y})$ . Note that  $\bar{x} = 30/10 = 3$  and  $\bar{y} = 510/10 = 51$ . Using equation (3.10), we obtain a sample covariance of

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1} = \frac{99}{9} = 11$$

The formula for computing the covariance of a population of size N is similar to equation (3.10), but we use different notation to indicate that we are working with the entire population.

TABLE 3.8 CALCULATIONS FOR THE SAMPLE COVARIANCE

$x_i$	$y_i$	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
2	50	-1	-1	1
5	57	2	6	12
1	41	-2	-10	20
3	54	0	3	0
4	54	1	3	3
1	38	-2	-13	26
5	63	2	12	24
3	48	0	-3	0
4	59	1	8	8
2	46	-1	5	_5
Totals 30	510		0	99
	Σ	$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} =$	99	

STUDENTS-HUB.com

POPULATION COVARIANCE  $\sigma_{xy} = \frac{\Sigma(x_i - \mu_x)(y_i - \mu_y)}{N} \tag{3.11}$ 

In equation (3.11) we use the notation  $\mu_x$  for the population mean of the variable x and  $\mu_y$  for the population mean of the variable y. The population covariance  $\sigma_{xy}$  is defined for a population of size N.

# Interpretation of the Covariance

To aid in the interpretation of the sample covariance, consider Figure 3.8. It is the same as the scatter diagram of Figure 3.7 with a vertical dashed line at  $\bar{x}=3$  and a horizontal dashed line at  $\bar{y}=51$ . The lines divide the graph into four quadrants. Points in quadrant I correspond to  $x_i$  greater than  $\bar{x}$  and  $y_i$  greater than  $\bar{y}$ , points in quadrant II correspond to  $x_i$  less than  $\bar{x}$  and  $y_i$  greater than  $\bar{y}$ , and so on. Thus, the value of  $(x_i-\bar{x})(y_i-\bar{y})$  must be positive for points in quadrant I, negative for points in quadrant II, positive for points in quadrant IV.

The covariance is a measure of the linear association between two variables.

If the value of  $s_{xy}$  is positive, the points with the greatest influence on  $s_{xy}$  must be in quadrants I and III. Hence, a positive value for  $s_{xy}$  indicates a positive linear association between x and y; that is, as the value of x increases, the value of y increases. If the value of  $s_{xy}$  is negative, however, the points with the greatest influence on  $s_{xy}$  are in quadrants II and IV. Hence, a negative value for  $s_{xy}$  indicates a negative linear association between x and y; that is, as the value of x increases, the value of x decreases. Finally, if the points are evenly distributed across all four quadrants, the value of  $s_{xy}$  will be close to zero, indicating no linear association between x and y. Figure 3.9 shows the values of  $s_{xy}$  that can be expected with three different types of scatter diagrams.

FIGURE 3.8 PARTITIONED SCATTER DIAGRAM FOR THE STEREO AND SOUND EQUIPMENT STORE

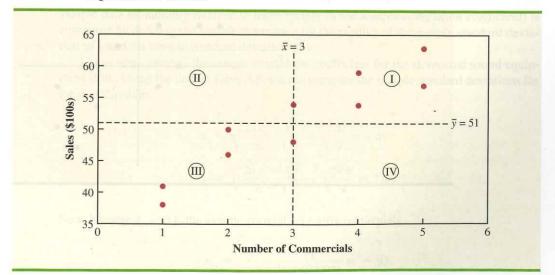
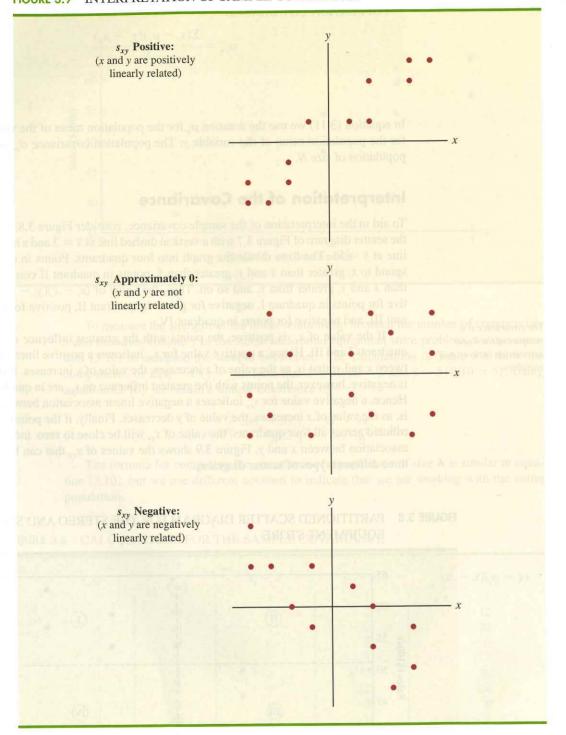


FIGURE 3.9 INTERPRETATION OF SAMPLE COVARIANCE



Referring again to Figure 3.8, we see that the scatter diagram for the stereo and sound equipment store follows the pattern in the top panel of Figure 3.9. As we should expect, the value of the sample covariance indicates a positive linear relationship with  $s_{rv} = 11$ .

From the preceding discussion, it might appear that a large positive value for the covariance indicates a strong positive linear relationship and that a large negative value indicates a strong negative linear relationship. However, one problem with using covariance as a measure of the strength of the linear relationship is that the value of the covariance depends on the units of measurement for x and y. For example, suppose we are interested in the relationship between height x and weight y for individuals. Clearly the strength of the relationship should be the same whether we measure height in feet or inches. Measuring the height in inches, however, gives us much larger numerical values for  $(x_i - \bar{x})$  than when we measure height in feet. Thus, with height measured in inches, we would obtain a larger value for the numerator  $\Sigma(x_i - \bar{x})(y_i - \bar{y})$  in equation (3.10)—and hence a larger covariance—when in fact the relationship does not change. A measure of the relationship between two variables that is not affected by the units of measurement for x and y is the **correlation coefficient**.

#### **Correlation Coefficient**

For sample data, the Pearson product moment correlation coefficient is defined as follows.

PEARSON PRODUCT MOMENT CORRELATION COEFFICIENT: SAMPLE DATA

$$xy = \frac{s_{xy}}{s_x s_y} \tag{3.12}$$

where

 $r_{xy}$  = sample correlation coefficient

 $s_{xy}$  = sample covariance

 $s_x = \text{sample standard deviation of } x$ 

 $s_y = \text{sample standard deviation of } y$ 

Equation (3.12) shows that the Pearson product moment correlation coefficient for sample data (commonly referred to more simply as the *sample correlation coefficient*) is computed by dividing the sample covariance by the product of the sample standard deviation of x and the sample standard deviation of y.

Let us now compute the sample correlation coefficient for the stereo and sound equipment store. Using the data in Table 3.8, we can compute the sample standard deviations for the two variables.

$$s_x = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n-1}} = \sqrt{\frac{20}{9}} = 1.49$$

$$s_y = \sqrt{\frac{\sum(y_i - \bar{y})^2}{n-1}} = \sqrt{\frac{566}{9}} = 7.93$$

Now, because  $s_{xy} = 11$ , the sample correlation coefficient equals

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{11}{(1.49)(7.93)} = +.93$$

The formula for computing the correlation coefficient for a population, denoted by the Greek letter  $\rho_{xy}$  (rho, pronounced "row"), follows.

PEARSON PRODUCT MOMENT CORRELATION COEFFICIENTS
POPULATION DATA

 $\rho_{xy} = \frac{\sigma_{xy}}{\sigma_{x}\sigma_{y}} \tag{3.13}$ 

The sample correlation coefficient  $r_{xy}$  is the estimator of the population correlation coefficient  $\rho_{xy}$ .

where

 $\rho_{xy}$  = population correlation coefficient

 $\sigma_{xy}$  = population covariance

 $\sigma_x$  = population standard deviation for x

 $\sigma_{\rm v}$  = population standard deviation for y

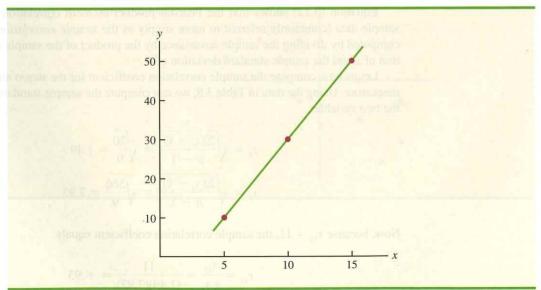
The sample correlation coefficient  $r_{xy}$  provides an estimate of the population correlation coefficient  $\rho_{xy}$ .

# **Interpretation of the Correlation Coefficient**

First let us consider a simple example that illustrates the concept of a perfect positive linear relationship. The scatter diagram in Figure 3.10 depicts the relationship between x and y based on the following sample data.

$x_i$	$y_i$
5	10
10	30
15	50

FIGURE 3.10 SCATTER DIAGRAM DEPICTING A PERFECT POSITIVE LINEAR RELATIONSHIP



STUDENTS-HUB.com

TABLE 3.9 COMPUTATIONS USED IN CALCULATING THE SAMPLE CORRELATION COEFFICIENT

	$x_i$	$y_i$	$x_i - \bar{x}$	$(x_i - \bar{x})^2$	$y_i - \bar{y}$	$(y_i - \bar{y})^2$	$(x_i - \bar{x})(y_i - \bar{y})$
	5	10	-5	25	-20	400	100
	10	30	0	0	0	0	0
	15	50	5	<u>25</u>	_20	400	100
Totals	30	90	0	50	0	800	200

The straight line drawn through each of the three points shows a perfect linear relationship between x and y. In order to apply equation (3.12) to compute the sample correlation we must first compute  $s_{xy}$ ,  $s_x$ , and  $s_y$ . Some of the computations are shown in Table 3.9. Using the results in Table 3.9, we find

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1} = \frac{200}{2} = 100$$

$$s_x = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n - 1}} = \sqrt{\frac{50}{2}} = 5$$

$$s_y = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n - 1}} = \sqrt{\frac{800}{2}} = 20$$

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{100}{5(20)} = 1$$

The correlation coefficient ranges from -1 to +1. Values close to -1 or +1 indicate a strong linear relationship. The closer the correlation is to zero, the weaker the relationship.

Thus, we see that the value of the sample correlation coefficient is 1.

In general, it can be shown that if all the points in a data set fall on a positively sloped straight line, the value of the sample correlation coefficient is +1; that is, a sample correlation coefficient of +1 corresponds to a perfect positive linear relationship between x and y. Moreover, if the points in the data set fall on a straight line having negative slope, the value of the sample correlation coefficient is -1; that is, a sample correlation coefficient of -1 corresponds to a perfect negative linear relationship between x and y.

Let us now suppose that a certain data set indicates a positive linear relationship between x and y but that the relationship is not perfect. The value of  $r_{xy}$  will be less than 1, indicating that the points in the scatter diagram are not all on a straight line. As the points deviate more and more from a perfect positive linear relationship, the value of  $r_{xy}$  becomes smaller and smaller. A value of  $r_{xy}$  equal to zero indicates no linear relationship between x and y, and values of  $r_{xy}$  near zero indicate a weak linear relationship.

For the data involving the stereo and sound equipment store, recall that  $r_{xy} = +.93$ . Therefore, we conclude that a strong positive linear relationship occurs between the number of commercials and sales. More specifically, an increase in the number of commercials is associated with an increase in sales.

In closing, we note that correlation provides a measure of linear association and not necessarily causation. A high correlation between two variables does not mean that changes in one variable will cause changes in the other variable. For example, we may find that the quality rating and the typical meal price of restaurants are positively correlated. However, simply increasing the meal price at a restaurant will not cause the quality rating to increase.

#### **Exercises**

#### Methods

SELF test

45. Five observations taken for two variables follow.

- a. Develop a scatter diagram with x on the horizontal axis.
- b. What does the scatter diagram developed in part (a) indicate about the relationship between the two variables?
- c. Compute and interpret the sample covariance.
- d. Compute and interpret the sample correlation coefficient.
- 46. Five observations taken for two variables follow.

- a. Develop a scatter diagram for these data.
- b. What does the scatter diagram indicate about a relationship between x and y?
- Compute and interpret the sample covariance.
- d. Compute and interpret the sample correlation coefficient.

# **Applications**

47. Nielsen Media Research provides two measures of the television viewing audience: a television program *rating*, which is the percentage of households with televisions watching a program, and a television program *share*, which is the percentage of households watching a program among those with televisions in use. The following data show the Nielsen television ratings and share data for the Major League Baseball World Series over a nine-year period (Associated Press, October 27, 2003).

Rating	19	17	17	14	16	12	15	12	13	
Share	32	28	29	24	26	20	24	20	22	

- a. Develop a scatter diagram with rating on the horizontal axis.
- b. What is the relationship between rating and share? Explain.
- c. Compute and interpret the sample covariance.
- d. Compute the sample correlation coefficient. What does this value tell us about the relationship between rating and share?
- 48. A department of transportation's study on driving speed and mileage for midsize automobiles resulted in the following data.

<b>Driving Speed</b>	30	50	40	55	30	25	60	25	50	55
Mileage	28	25	25	23	30	32	21	35	26	25

Compute and interpret the sample correlation coefficient.

49. PC World provided ratings for 15 notebook PCs (PC World, February 2000). The performance score is a measure of how fast a PC can run a mix of common business applications as compared to a baseline machine. For example, a PC with a performance score of 200 is twice as fast as the baseline machine. A 100-point scale was used to provide an overall rating for each notebook tested in the study. A score in the 90s is exceptional, while one in the 70s is good. Table 3.10 shows the performance scores and the overall ratings for the 15 notebooks.

# STUDENTS-HUB.com

## TABLE 3.10 PERFORMANCE SCORES AND OVERALL RATINGS FOR 15 NOTEBOOK PCs



Notebook	Performance Score	Overall Rating
AMS Tech Roadster 15CTA380	115	67
Compaq Armada M700	191	78
Compaq Prosignia Notebook 150	153	79
Dell Inspiron 3700 C466GT	194	80
Dell Inspiron 7500 R500VT	236	84
Dell Latitude Cpi A366XT	184	76
Enpower ENP-313 Pro	184	77
Gateway Solo 9300LS	216	92
HP Pavilion Notebook PC	185	83
IBM ThinkPad I Series 1480	183	78
Micro Express NP7400	189	77
Micron TransPort NX PII-400	202	78
NEC Versa SX	192	78
Sceptre Soundx 5200	141	73
Sony VAIO PCG-F340	187	77

- a. Compute the sample correlation coefficient.
- b. What does the sample correlation coefficient tell about the relationship between the performance score and the overall rating?
- 50. The Dow Jones Industrial Average (DJIA) and the Standard & Poor's 500 Index (S&P 500) are both used to measure the performance of the stock market. The DJIA is based on the price of stocks for 30 large companies; the S&P 500 is based on the price of stocks for 500 companies. If both the DJIA and S&P 500 measure the performance of the stock market, how are they correlated? The following data show the daily percent increase or daily percent decrease in the DJIA and S&P 500 for a sample of nine days over a three-month period (*The Wall Street Journal*, January 15 to March 10, 2006).



DJIA	.20	.82	99	.04	24	1.01	.30	.55	25
S&P 500	.24	.19	91	.08	33	.87	.36	.83	16

- Show a scatter diagram.
- b. Compute the sample correlation coefficient for these data.
- c. Discuss the association between the DJIA and S&P 500. Do you need to check both before having a general idea about the daily stock market performance?
- The daily high and low temperatures for 12 U.S. cities are as follows (Weather Channel, January 25, 2004).



City	High	Low	City	High	Low
Albany	9	-8	Los Angeles	62	47
Boise	32	26	New Orleans	71	55
Cleveland	21	19	Portland	43	36
Denver	37	10	Providence	18	8
Des Moines	24	16	Raleigh	28	24
Detroit	20	17	Tulsa	55	38

- a. What is the sample mean daily high temperature?
- b. What is the sample mean daily low temperature?
- c. What is the correlation between the high and low temperatures?



# The Weighted Mean and Working with Grouped Data

In Section 3.1, we presented the mean as one of the most important measures of central location. The formula for the mean of a sample with n observations is restated as follows.

$$\bar{x} = \frac{\sum x_i}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$
 (3.14)

In this formula, each  $x_i$  is given equal importance or weight. Although this practice is most common, in some instances, the mean is computed by giving each observation a weight that reflects its importance. A mean computed in this manner is referred to as a **weighted mean**.

# **Weighted Mean**

The weighted mean is computed as follows:

WEIGHTED MEAN	Su. The Dow Jones Industrial Average (DI)	
	$\bar{x} = \frac{2W_i x_i}{\sum_{w}}$	(3.15)
	comment to provide a chorac manufall That to live in the last of the character and the last of the character and the cha	
	y - value of observation i	
	$w_i$ = weight for observation $i$	

When the data are from a sample, equation (3.15) provides the weighted sample mean. When the data are from a population,  $\mu$  replaces  $\bar{x}$  and equation (3.15) provides the weighted population mean.

As an example of the need for a weighted mean, consider the following sample of five purchases of a raw material over the past three months.

Purchase	Cost per Pound (\$)	Number of Pounds
and 12 U.S. and	3.00	1200
2	3.40	500
3	2.80	2750
4	2.90	1000
5	3.25	800

Note that the cost per pound varies from \$2.80 to \$3.40, and the quantity purchased varies from 500 to 2750 pounds. Suppose that a manager asked for information about the mean cost per pound of the raw material. Because the quantities ordered vary, we must use the formula for a weighted mean. The five cost-per-pound data values are  $x_1 = 3.00$ ,  $x_2 = 3.40$ ,  $x_3 = 2.80$ ,  $x_4 = 2.90$ , and  $x_5 = 3.25$ . The weighted mean cost per pound is found by weighting each cost

STUDENTS-HUB.com

by its corresponding quantity. For this example, the weights are  $w_1 = 1200$ ,  $w_2 = 500$ ,  $w_3 = 2750$ ,  $w_4 = 1000$ , and  $w_5 = 800$ . Based on equation (3.15), the weighted mean is calculated as follows:

$$\bar{x} = \frac{1200(3.00) + 500(3.40) + 2750(2.80) + 1000(2.90) + 800(3.25)}{1200 + 500 + 2750 + 1000 + 800}$$
$$= \frac{18,500}{6250} = 2.96$$

Thus, the weighted mean computation shows that the mean cost per pound for the raw material is \$2.96. Note that using equation (3.14) rather than the weighted mean formula would have provided misleading results. In this case, the mean of the five cost-per-pound values is (3.00 + 3.40 + 2.80 + 2.90 + 3.25)/5 = 15.35/5 = \$3.07, which overstates the actual mean cost per pound purchased.

The choice of weights for a particular weighted mean computation depends upon the application. An example that is well known to college students is the computation of a grade point average (GPA). In this computation, the data values generally used are 4 for an A grade, 3 for a B grade, 2 for a C grade, 1 for a D grade, and 0 for an F grade. The weights are the number of credits hours earned for each grade. Exercise 54 at the end of this section provides an example of this weighted mean computations, quantities such as pounds, dollars, or volume are frequently used as weights. In any case, when observations vary in importance, the analyst must choose the weight that best reflects the importance of each observation in the determination of the mean.

Computing a grade point average is a good example of the use of a weighted mean.

# **Grouped Data**

In most cases, measures of location and variability are computed by using the individual data values. Sometimes, however, data are available only in a grouped or frequency distribution form. In the following discussion, we show how the weighted mean formula can be used to obtain approximations of the mean, variance, and standard deviation for **grouped data**.

In Section 2.2 we provided a frequency distribution of the time in days required to complete year-end audits for the public accounting firm of Sanderson and Clifford. The frequency distribution of audit times based on a sample of 20 clients is shown again in Table 3.11. Based on this frequency distribution, what is the sample mean audit time?

To compute the mean using only the grouped data, we treat the midpoint of each class as being representative of the items in the class. Let  $M_i$  denote the midpoint for class i and let  $f_i$  denote the frequency of class i. The weighted mean formula (3.15) is then used with the data values denoted as  $M_i$  and the weights given by the frequencies  $f_i$ . In this case, the denominator of equation (3.15) is the sum of the frequencies, which is the

TABLE 3.11 FREQUENCY DISTRIBUTION OF AUDIT TIMES

Audit Time (days)	Frequency	
10-14	4	
15-19	8	
20-24	5	
25-29	2	
30-34	1	
Total	20	

3.6 The Weighted Mean and Working with Grouped Data

sample size n. That is,  $\Sigma f_i = n$ . Thus, the equation for the sample mean for grouped data is as follows.

SAMPLE MEAN FOR GROUPED DATA

$$\bar{x} = \frac{\sum f_i M_i}{n} \tag{3.16}$$

where

 $M_i$  = the midpoint for class i

 $f_i$  = the frequency for class i

n = the sample size

With the class midpoints,  $M_i$ , halfway between the class limits, the first class of 10-14 in Table 3.11 has a midpoint at (10 + 14)/2 = 12. The five class midpoints and the weighted mean computation for the audit time data are summarized in Table 3.12. As can be seen, the sample mean audit time is 19 days.

To compute the variance for grouped data, we use a slightly altered version of the formula for the variance provided in equation (3.5). In equation (3.5), the squared deviations of the data about the sample mean  $\bar{x}$  were written  $(x_i - \bar{x})^2$ . However, with grouped data, the values are not known. In this case, we treat the class midpoint,  $M_i$ , as being representative of the  $x_i$  values in the corresponding class. Thus, the squared deviations about the sample mean,  $(x_i - \bar{x})^2$ , are replaced by  $(M_i - \bar{x})^2$ . Then, just as we did with the sample mean calculations for grouped data, we weight each value by the frequency of the class,  $f_i$ . The sum of the squared deviations about the mean for all the data is approximated by  $\sum f_i(M_i - \bar{x})^2$ . The term n - 1 rather than n appears in the denominator in order to make the sample variance the estimate of the population variance. Thus, the following formula is used to obtain the sample variance for grouped data.

SAMPLE VARIANCE FOR GROUPED DATA

$$s^2 = \frac{\sum f_i (M_i - \bar{x})^2}{n - 1}$$
 (3.17)

TABLE 3.12 COMPUTATION OF THE SAMPLE MEAN AUDIT TIME FOR GROUPED DATA

Audit Time (days)	Class Midpoint $(M_i)$	Frequency $(f_i)$	$f_i M_i$
10-14	12	4	48
15-19	17	8	136
20-24	22	5	110
25-29	27	2	54
30–34	32	1	_32
		20	380
	Sample mean $\bar{x} = \frac{\sum f_i M_i}{n} =$	$\frac{380}{20} = 19 \text{ days}$	

STUDENTS-HUB.com

TABLE 3.13 COMPUTATION OF THE SAMPLE VARIANCE OF AUDIT TIMES FOR GROUPED DATA (SAMPLE MEAN  $\bar{x} = 19$ )

Audit Time (days)	Class Midpoint $(M_i)$	Frequency $(f_i)$	Deviation $(M_i - \bar{x})$	Squared Deviation $(M_i - \bar{x})^2$	$f_i(M_i - \bar{x})^2$
10-14	12	4	<del>-7</del>	49	196
15-19	17	8	-2	4	32
20-24	22	5	3	9	45
25-29	27	2	8	64	128
30-34	32	1	13	169	169
		20		X I	570
					$\sum f_i (M_i - \bar{x})^2$
	Sa	ample variance $s^2 =$	$\frac{\sum f_i (M_i - \bar{x})^2}{n - 1} = \frac{5}{1}$	$\frac{70}{19} = 30$	

The calculation of the sample variance for audit times based on the grouped data from Table 3.11 is shown in Table 3.13. As can be seen, the sample variance is 30.

The standard deviation for grouped data is simply the square root of the variance for grouped data. For the audit time data, the sample standard deviation is  $s = \sqrt{30} = 5.48$ .

Before closing this section on computing measures of location and dispersion for grouped data, we note that formulas (3.16) and (3.17) are for a sample. Population summary measures are computed similarly. The grouped data formulas for a population mean and variance follow.

POPULATION MEAN FOR GROUPED DATA

$$\mu = \frac{\sum f_i M_i}{N}$$
 (3.18)

POPULATION VARIANCE FOR GROUPED DATA

$$\sigma^2 = \frac{\sum f_i (M_i - \mu)^2}{N}$$
 (3.19)

#### **NOTES AND COMMENTS**

In computing descriptive statistics for grouped data, the class midpoints are used to approximate the data values in each class. As a result, the descriptive statistics for grouped data approximate the descriptive statistics that would result from us-

ing the original data directly. We therefore recommend computing descriptive statistics from the original data rather than from grouped data whenever possible.

#### Exercises

## Methods

52. Consider the following data and corresponding weights.

$x_i$	Weight $(w_i)$
3.2	6
2.0	3
2.5	2
5.0	8

- a. Compute the weighted mean.
- b. Compute the sample mean of the four data values without weighting. Note the difference in the results provided by the two computations.
- 53. Consider the sample data in the following frequency distribution.

SELF	test
------	------

Class	Midpoint	Frequency
3-7	5	4
8-12	10	7
13-17	15	9
18-22	20	5

- a. Compute the sample mean.
- b. Compute the sample variance and sample standard deviation.

# **Applications**



- 54. The grade point average for college students is based on a weighted mean computation. For most colleges, the grades are given the following data values: A (4), B (3), C (2), D (1), and F (0). After 60 credit hours of course work, a student at State University earned 9 credit hours of A, 15 credit hours of B, 33 credit hours of C, and 3 credit hours of D.
  - a. Compute the student's grade point average.
  - b. Students at State University must maintain a 2.5 grade point average for their first 60 credit hours of course work in order to be admitted to the business college. Will this student be admitted?
- 55. Bloomberg Personal Finance (July/August 2001) included the following companies in its recommended investment portfolio. For a portfolio value of \$25,000, the recommended dollar amounts allocated to each stock are shown.

	Portfolio	Estimated Growth Rate	Dividend Yield
Company	(\$)	(%)	(70)
Citigroup	3000	15	1.21
General Electric	5500	14	1.48
Kimberly-Clark	4200	12	1.72
Oracle	3000	25	0.00
Pharmacia	3000	20	0.96
SBC Communications	3800	12	2.48
WorldCom	2500	35	0.00

- a. Using the portfolio dollar amounts as the weights, what is the weighted average estimated growth rate for the portfolio?
- b. What is the weighted average dividend yield for the portfolio?
- 56. A survey of subscribers to Fortune magazine asked the following question: "How many of the last four issues have you read?" Suppose that the following frequency distribution summarizes 500 responses.

Number Read	Frequency	
0	15	
1	10	
2 - 10 - 2 - 10 - 2 - 10 - 10 - 10	40	
3	85	
4	350	
	Total 500	

- a. What is the mean number of issues read by a Fortune subscriber?
- b. What is the standard deviation of the number of issues read?
- 57. The following frequency distribution shows the price per share for the 30 companies in the Dow Jones Industrial Average (*The Wall Street Journal*, January 16, 2006).

Price per Share	Frequency	
\$20-29	7	
\$30-39	6	
\$40-49	6	
\$50-59	3	
\$60-69	4	
\$70–79	3	
\$80-89	1	

Compute the mean price per share and the standard deviation of the price per share for the Dow Jones Industrial Average companies.

# Summary

In this chapter we introduced several descriptive statistics that can be used to summarize the location, variability, and shape of a data distribution. Unlike the tabular and graphical procedures introduced in Chapter 2, the measures introduced in this chapter summarize the data in terms of numerical values. When the numerical values obtained are for a sample, they are called sample statistics. When the numerical values obtained are for a population, they are called population parameters. Some of the notation used for sample statistics and population parameters follow.

In statistical inference, the sample statistic is referred to as the point estimator of the population parameter.

	Sample Statistic	Population Parameter	
Mean	$\bar{x}$	$\mu$	
Variance	$s^2$	$\sigma^2$	
Standard deviation	other bosoms a call or myter	per blomb	
Covariance	$s_{xy}$	$\sigma_{xy}$	
Correlation	which satisfy $r_{xy}$ is essentially	$\rho_{xy}$	

STUDENTS-HUB.com

Key Formulas

125

As measures of central location, we defined the mean, median, and mode. Then the concept of percentiles was used to describe other locations in the data set. Next, we presented the range, interquartile range, variance, standard deviation, and coefficient of variation as measures of variability or dispersion. Our primary measure of the shape of a data distribution was the skewness. Negative values indicate a data distribution skewed to the left. Positive values indicate a data distribution skewed to the right. We then described how the mean and standard deviation could be used, applying Chebyshev's theorem and the empirical rule, to provide more information about the distribution of data and to identify outliers.

In Section 3.4 we showed how to develop a five-number summary and a box plot to provide simultaneous information about the location, variability, and shape of the distribution. In Section 3.5 we introduced covariance and the correlation coefficient as measures of association between two variables. In the final section, we showed how to compute a weighted mean and how to calculate a mean, variance, and standard deviation for grouped data.

The descriptive statistics we discussed can be developed using statistical software packages and spreadsheets. In Appendix 3.1 we show how to develop most of the descriptive statistics introduced in the chapter using Minitab. In Appendix 3.2, we demonstrate the use of Excel for the same purpose.

# Glossary

**Sample statistic** A numerical value used as a summary measure for a sample (e.g., the sample mean,  $\bar{x}$ , the sample variance,  $s^2$ , and the sample standard deviation, s).

**Population parameter** A numerical value used as a summary measure for a population (e.g., the population mean,  $\mu$ , the population variance,  $\sigma^2$ , and the population standard deviation,  $\sigma$ ).

**Point estimator** The sample statistic, such as  $\bar{x}$ ,  $s^2$ , and s, when used to estimate the corresponding population parameter.

Mean A measure of central location computed by summing the data values and dividing by the number of observations.

Median A measure of central location provided by the value in the middle when the data are arranged in ascending order.

Mode A measure of location, defined as the value that occurs with greatest frequency.

**Percentile** A value such that at least p percent of the observations are less than or equal to this value and at least (100 - p) percent of the observations are greater than or equal to this value. The 50th percentile is the median.

Quartiles The 25th, 50th, and 75th percentiles, referred to as the first quartile, the second quartile (median), and third quartile, respectively. The quartiles can be used to divide a data set into four parts, with each part containing approximately 25% of the data.

Range A measure of variability, defined to be the largest value minus the smallest value.

Interquartile range (IQR) A measure of variability, defined to be the difference between

the third and first quartiles.

Variance A measure of variability based on the squared deviations of the data values about the mean.

**Standard deviation** A measure of variability computed by taking the positive square root of the variance.

**Coefficient of variation** A measure of relative variability computed by dividing the standard deviation by the mean and multiplying by 100.

**Skewness** A measure of the shape of a data distribution. Data skewed to the left result in negative skewness; a symmetric data distribution results in zero skewness; and data skewed to the right result in positive skewness.

STUDENTS-HUB.com

**z-score** A value computed by dividing the deviation about the mean  $(x_i - \bar{x})$  by the standard deviation s. A z-score is referred to as a standardized value and denotes the number of standard deviations  $x_i$  is from the mean.

**Chebyshev's theorem** A theorem that can be used to make statements about the proportion of data values that must be within a specified number of standard deviations of the mean.

**Empirical rule** A rule that can be used to compute the percentage of data values that must be within one, two, and three standard deviations of the mean for data that exhibit a bell-shaped distribution.

Outlier An unusually small or unusually large data value.

**Five-number summary** An exploratory data analysis technique that uses five numbers to summarize the data: smallest value, first quartile, median, third quartile, and largest value. **Box plot** A graphical summary of data based on a five-number summary.

**Covariance** A measure of linear association between two variables. Positive values indicate a positive relationship; negative values indicate a negative relationship.

Correlation coefficient A measure of linear association between two variables that takes on values between -1 and +1. Values near +1 indicate a strong positive linear relationship; values near -1 indicate a strong negative linear relationship; and values near zero indicate the lack of a linear relationship.

**Weighted mean** The mean obtained by assigning each observation a weight that reflects its importance.

**Grouped data** Data available in class intervals as summarized by a frequency distribution. Individual values of the original data are not available.

# **Key Formulas**

#### Sample Mean

$$\bar{x} = \frac{\sum x_i}{n} \tag{3.1}$$

#### **Population Mean**

$$\mu = \frac{\sum x_i}{N} \tag{3.2}$$

#### **Interquartile Range**

$$IQR = Q_3 - Q_1 (3.3)$$

#### **Population Variance**

$$\sigma^2 = \frac{\sum (x_i - \mu)^2}{N} \tag{3.4}$$

#### Sample Variance

$$s^2 = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$
 (3.5)

#### Standard Deviation

Sample standard deviation = 
$$s = \sqrt{s^2}$$
 (3.6)

Population standard deviation = 
$$\sigma = \sqrt{\sigma^2}$$
 (3.7)

Supplementary Exercises

#### **Coefficient of Variation**

$$\left(\frac{\text{Standard deviation}}{\text{Mean}} \times 100\right)\%$$
 (3.8)

z-Score

$$z_i = \frac{x_i - \bar{x}}{s} \tag{3.9}$$

Sample Covariance

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$
 (3.10)

**Population Covariance** 

$$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$
 (3.11)

Pearson Product Moment Correlation Coefficient: Sample Data

$$r_{xy} = \frac{s_{xy}}{s_x s_y} \tag{3.12}$$

Pearson Product Moment Correlation Coefficient: Population Data

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \tag{3.13}$$

Weighted Mean

$$\bar{x} = \frac{\sum w_i x_i}{\sum w_i}$$
 (3.15)

Sample Mean for Grouped Data

$$\bar{x} = \frac{\sum f_i M_i}{n} \tag{3.16}$$

Sample Variance for Grouped Data

$$s^2 = \frac{\sum f_i (M_i - \bar{x})^2}{n - 1}$$
 (3.17)

**Population Mean for Grouped Data** 

$$\mu = \frac{\sum f_i M_i}{N} \tag{3.18}$$

**Population Variance for Grouped Data** 

$$\sigma^2 = \frac{\sum f_i (M_i - \mu)^2}{N}$$
 (3.19)

STUDENTS-HUB.com

# **Supplementary Exercises**

58. The following data show the media expenditures (\$ millions) and shipments in millions of barrels (bbls.) for 10 major brands of beer.



Brand	Media Expenditures (\$ millions)	Shipments in bbls (millions)
Budweiser	120.0	36.3
Bud Light	68.7	20.7
Ailler Lite	100.1	15.9
Coors Light	76.6	13.2
Susch	8.7	8.1
Vatural Light	0.1	7.1
Iiller Genuine Draft	21.5	5.6
filler High Life	1.4	4.4
usch Lite	5.3	4.3
filwaukee's Best	1.7	4.3

- a. What is the sample covariance? Does it indicate a positive or negative relationship?
- b. What is the sample correlation coefficient?
- 59. The U.S. Department of Education reports that about 50% of all college students use a student loan to help cover college expenses (National Center for Educational Studies, January 2006). A sample of students who graduated with student loan debt is shown here. The data, in thousands of dollars, show typical amounts of debt upon graduation.

10.1 14.8 5.0 10.2 12.4 12.2 2.0 11.5 17.8 4.0

- a. For those students who use a student loan, what is the mean loan debt upon graduation?
- b. What is the variance? Standard deviation?
- 60. The National Association of Realtors reported the median home price in the United States and the increase in median home price over a five-year period (*The Wall Street Journal*, January 16, 2006). Use the sample home prices shown here to answer the following questions.



995.9	48.8	175.0	263.5	298.0	218.9	209.0
628.3	111.0	212.9	92.6	2325.0	958.0	212.5

- a. What is the sample median home price?
- b. In January 2001, the National Association of Realtors reported a median home price of \$139,300 in the United States. What was the percentage increase in the median home price over the five-year period?
- c. What are the first quartile and the third quartile for the sample data?
- d. Provide a five-number summary for the home prices.
- e. Do the data contain any outliers?
- What is the mean home price for the sample? Why does the National Association of Realtors prefer to use the median home price in its reports?

61. Automobiles traveling on a road with a posted speed limit of 55 miles per hour are checked for speed by a state police radar system. Following is a frequency distribution of speeds.

Speed (miles per hour)	Frequency	
45-49	10	
50-54	40	
55-59	150	
60-64	175	
65-69	75	
70-74	15	
75-79	10	
	Total 475	

- a. What is the mean speed of the automobiles traveling on this road?
- b. Compute the variance and the standard deviation.
- 62. Public transportation and the automobile are two methods an employee can use to get to work each day. Samples of times recorded for each method are shown. Times are in minutes.

Public Transportation:	28	29	32	37	33	25	29	32	41	34
Automobile:	29	31	33	32	34	30	31	32	35	33

- a. Compute the sample mean time to get to work for each method.
- b. Compute the sample standard deviation for each method.
- c. On the basis of your results from parts (a) and (b), which method of transportation should be preferred? Explain.
- d. Develop a box plot for each method. Does a comparison of the box plots support your conclusion in part (c)?
- 63. The days to maturity for a sample of five money market funds are shown here. The dollar amounts invested in the funds are provided. Use the weighted mean to determine the mean number of days to maturity for dollars invested in these five money market funds.

Days to Maturity	Dollar Value (\$ millions)
20	20
12	30
7	10
5	15
6	10

64. Small business owners often look to payroll service companies to handle their employee payroll. Reasons are that small business owners face complicated tax regulations and penalties for employment tax errors are costly. According to the Internal Revenue Service, 26% of all small business employment tax returns contained errors that resulted in a tax penalty to the owner (*The Wall Street Journal*, January 30, 2006). The tax penalties for a sample of 20 small business owners follow:

820	270	450	1010	890	700	1350	350	300	1200
390	730	2040	230	640	350	420	270	370	620

- a. What is the mean tax penalty for improperly filed employment tax returns?
- b. What is the standard deviation?

c. Is the highest penalty, \$2040, an outlier?

- d. What are some of the advantages of a small business owner hiring a payroll service company to handle employee payroll services, including the employment tax returns?
- 65. The following data show the trailing 52-week primary share earnings and book values as reported by 10 companies (*The Wall Street Journal*, March 13, 2000).

	Company	Book Value	Earnings
	Am Elec	25.21	2.69
	Columbia En	23.20	3.01
w har-come	Con Ed	25.19	3.13
	Duke Energy	20.17	2.25
	Edison Int'l	13.55	1.79
	Enron Cp.	7.44	1.27
	Peco	13.61	3.15
	Pub Sv Ent	21.86	3.29
	Southn Co.	8.77	1.86
	Unicom	23.22	2.74

- a. Develop a scatter diagram for the data with book value on the x-axis.
- b. What is the sample correlation coefficient, and what does it tell you about the relationship between the earnings per share and the book value?
- 66. Dividend yield is the annual dividend per share a company pays divided by the current market price per share expressed as a percentage. A sample of 10 large companies provided the following dividend yield data (*The Wall Street Journal*, January 16, 2004).

Company	Yield %	Company	Yield %
Altria Group	5.0	General Motors	3.7
American Express	0.8	JPMorgan Chase	3.5
Caterpillar	1.8	McDonald's	1.6
Eastman Kodak	1.9	United Technology	1.5
ExxonMobil	2.5	Wal-Mart Stores	0.7

- a. What are the mean and median dividend yields?
- b. What are the variance and standard deviation?
- c. Which company provides the highest dividend yield?
- d. What is the z-score for McDonald's? Interpret this z-score.
- e. What is the z-score for General Motors? Interpret this z-score.
- f. Based on z-scores, do the data contain any outliers?
- 67. A forecasting technique referred to as moving averages uses the average or mean of the most recent *n* periods to forecast the next value for time series data. With a three-period moving average, the most recent three periods of data are used in the forecast computation. Consider a product with the following demand for the first three months of the current year: January (800 units), February (750 units), and March (900 units).
  - a. What is the three-month moving average forecast for April?
  - b. A variation of this forecasting technique is called weighted moving averages. The weighting allows the more recent time series data to receive more weight or more importance in the computation of the forecast. For example, a weighted three-month moving average might give a weight of 3 to data one month old, a weight of 2 to data two months old, and a weight of 1 to data three months old. Use the data given to provide a three-month weighted moving average forecast for April.



68. The U.S. Census Bureau provides statistics on family life in the United States, including the age at the time of first marriage, current marital status, and size of household (http://www.census.gov, March 20, 2006). The following data show the age at the time of first marriage for a sample of men and a sample of women.



Men	26	23	28	25	27	30	26	35	28
	21	24	27	29	30	27	32	27	25
Women	20	28	23	30	24	29	26	25	
	22	22	25	23	27	26	19		

- a. Determine the median age at the time of first marriage for men and women.
- b. Compute the first and third quartiles for both men and women.
- c. Twenty-five years ago the median age at the time of first marriage was 25 for men and 22 for women. What insight does this information provide about the decision of when to marry among young people today?
- Road & Track provided the following sample of the tire ratings and load-carrying capacity
  of automobiles tires.

Tire Rating	Load-Carrying Capacity
75	853
82	1047
85	1135
87	1201
88	1235
91	1356
92	1389
93	1433
105	2039

- . Develop a scatter diagram for the data with tire rating on the x-axis.
- b. What is the sample correlation coefficient, and what does it tell you about the relationship between tire rating and load-carrying capacity?
- 70. According to the 2003 Annual Consumer Spending Survey, the average monthly Bank of America Visa credit card charge was \$1838 (U.S. Airways Attaché Magazine, December 2003). A sample of monthly credit card charges provides the following data.



236	1710	1351	825	7450
316	4135	1333	1584	387
991	3396	170	1428	1688

- a. Compute the mean and median.
- b. Compute the first and third quartiles.
- c. Compute the range and interquartile range.
- d. Compute the variance and standard deviation.
- e. The skewness measure for these data is 2.12. Comment on the shape of this distribution. Is it the shape you would expect? Why or why not?
- f. Do the data contain outliers?

# Case Problem 1 Pelican Stores

Pelican Stores, a division of National Clothing, is a chain of women's apparel stores operating throughout the country. The chain recently ran a promotion in which discount

STUDENTS-HUB.com

coupons were sent to customers of other National Clothing stores. Data collected for a sample of 100 in-store credit card transactions at Pelican Stores during one day while the promotion was running are contained in the file named PelicanStores. Table 3.14 shows a portion of the data set. The proprietary card method of payment refers to charges made using a National Clothing charge card. Customers who made a purchase using a discount coupon are referred to as promotional customers and customers who made a purchase but did not use a discount coupon are referred to as regular customers. Because the promotional coupons were not sent to regular Pelican Stores customers, management considers the sales made to people presenting the promotional coupons as sales it would not otherwise make. Of course, Pelican also hopes that the promotional customers will continue to shop at its stores.

Most of the variables shown in Table 3.14 are self-explanatory, but two of the variables require some clarification.

Items The total number of items purchased

Net Sales The total amount (\$) charged to the credit card

Pelican's management would like to use this sample data to learn about its customer base and to evaluate the promotion involving discount coupons.

# **Managerial Report**

Use the methods of descriptive statistics presented in this chapter to summarize the data and comment on your findings. At a minimum, your report should include the following:

- 1. Descriptive statistics on net sales and descriptive statistics on net sales by various classifications of customers.
- 2. Descriptive statistics concerning the relationship between age and net sales.





Customer	Type of Customer	Items	Net Sales	Method of Payment	Gender	Marital Status	Age
1	Regular	1	39.50	Discover	Male	Married	32
2	Promotional	1	102.40	Proprietary Card	Female	Married	36
3	Regular	1	22.50	Proprietary Card	Female	Married	32
4	Promotional	5	100.40	Proprietary Card	Female	Married	28
5	Regular	2	54.00	MasterCard	Female	Married	34
6	Regular	1	44.50	MasterCard	Female	Married	44
7	Promotional	2	78.00	Proprietary Card	Female	Married	30
8	Regular	1	22.50	Visa	Female	Married	40
9	Promotional	2	56.52	Proprietary Card	Female	Married	46
10	Regular	1	44.50	Proprietary Card	Female	Married	36
WIRA sump	Hint drive		orbs, Witter	Caller With Lovie 485	STIPS-E-18		
3.0	70511	*	A Limb	in alternation of	HONE F		
96	Regular	- ais	39.50	MasterCard	Female	Married	44
97	Promotional	9	253.00	Proprietary Card	Female	Married	30
98	Promotional	10	287.59	Proprietary Card	Female	Married	52
99	Promotional	2	47.60	Proprietary Card	Female	Married	30
100	Promotional	1	28.44	Proprietary Card	Female	Married	44

# Case Problem 2 Motion Picture Industry

The motion picture industry is a competitive business. More than 50 studios produce a total of 300 to 400 new motion pictures each year, and the financial success of each motion picture varies considerably. The opening weekend gross sales, the total gross sales, the number of theaters the movie was shown in, and the number of weeks the motion picture was in the top 60 for gross sales are common variables used to measure the success of a motion picture. Data collected for a sample of 100 motion pictures produced in 2005 are contained in the file named Movies. Table 3.15 shows the data for the first 10 motion pictures in the file.

# **Managerial Report**

Use the numerical methods of descriptive statistics presented in this chapter to learn how these variables contribute to the success of a motion picture. Include the following in your report.

- 1. Descriptive statistics for each of the four variables along with a discussion of what the descriptive statistics tell us about the motion picture industry.
- 2. What motion pictures, if any, should be considered high-performance outliers? Explain.
- **3.** Descriptive statistics showing the relationship between total gross sales and each of the other variables. Discuss.

TABLE 3.15 PERFORMANCE DATA FOR 10 MOTION PICTURES



Motion Picture	Opening Weekend Gross Sales (\$ millions)	Total Gross Sales (\$ millions)	Number of Theaters	Weeks in Top 60
Coach Carter	29.17	67.25	2574	16
Ladies in Lavender	0.15	6.65	119	22
Batman Begins	48.75	205.28	3858	18
Unleashed	10.90	24.47	1962	8
Pretty Persuasion	0.06	0.23	24	4
Fever Pitch	12.40	42.01	3275	- 14
Harry Potter and the Goblet of Fire	102.69	287.18	3858	13
Monster-in-Law	23.11	82.89	3424	16
White Noise	24.11	55.85	2279	7
Mr. and Mrs. Smith	50.34	186.22	3451	21

# Case Problem 3 Business Schools of Asia-Pacific



The pursuit of a higher education degree in business is now international. A survey shows that more and more Asians choose the Master of Business Administration degree route to corporate success. As a result, the number of applicants for MBA courses at Asia-Pacific schools continues to increase.

Across the region, thousands of Asians show an increasing willingness to temporarily shelve their careers and spend two years in pursuit of a theoretical business qualification. Courses in these schools are notoriously tough and include economics, banking, marketing, behavioral sciences, labor relations, decision making, strategic thinking, business law, and more. The data set in Table 3.16 shows some of the characteristics of the leading Asia-Pacific business schools.

STUDENTS-HUB.com

# 3LE 3.16 DATA FOR 25 ASIA-PACIFIC BUSINESS SCHOOLS

Business School	Full-Time Enrollment	Students per Faculty	Local Tuition (\$)	Foreign Tuition (\$)	Age	%Foreign	GMAT	English Test	Work Experience	Starting Salary (\$)
Melbourne Business School	200	5	24,420	29,600	28	47	Yes	No	Yes	71,400
University of New South Wales (Sydney)	228	4	19,993	32,582	29	28	Yes	No	Yes	65,200
Indian Institute of Management (Ahmedabad)	392	5	4,300	4,300	22	0	No	No	No	7,100
Chinese University of Hong Kong	06	5	11,140	11,140	29	10	Yes	No	No	31,000
International University of Japan (Niigata)	126	4	33,060	33,060	28	09	Yes	Yes	No	87,000
Asian Institute of Management (Manila)	389	5	7,562	000,6	25	50	Yes	No	Yes	22,800
Indian Institute of Management (Bangalore)	380	5	3,935	16,000	23		Yes	No	No	7,500
National University of Singapore	147	9	6,146	7,170	29	51	Yes	Yes	Yes	43,300
Indian Institute of Management (Calcutta)	463	8	2,880	16,000	23	0	No	No	No	7,400
Australian National University (Canberra)	42	2	20,300	20,300	30	80	Yes	Yes	Yes	46,600
Nanyang Technological University (Singapore)	50	5	8,500	8,500	32	20	Yes	No	Yes	49,300
University of Queensland (Brisbane)	138	17	16,000	22,800	32	26	No	No	Yes	49,600
Hong Kong University of Science and Technology	09	2	11,513	11,513	26	37	Yes	No	Yes	34,000
Macquarie Graduate School of Management (Sydney)	12	∞	17,172	19,778	34	27	No	No	Yes	60,100
Chulalongkorn University (Bangkok)	200	7	17,355	17,355	25	9	Yes	No	Yes	17,600
Monash Mt. Eliza Business School (Melbourne)	350	13	16,200	22,500	30	30	Yes	Yes	Yes	52,500
Asian Institute of Management (Bangkok)	300	10	18,200	18,200	29	06	No	Yes	Yes	25,000
University of Adelaide	20	19	16,426	23,100	30	10	No	No	Yes	000,99
Massey University (Palmerston North, New Zealand)	30	15	13,106	21,625	37	35	No	Yes	Yes	41,400
Royal Melbourne Institute of Technology Business	id.				ì		A A L			o (
Graduate School	30	7	13,880	17,765	32	30	No	Yes	Yes	48,900
Jamnalal Bajaj Institute of Management Studies (Mumbia)	240	6	1,000	1,000	24	0	No	No	Yes	7,000
Curtin Institute of Technology (Perth)	86	15	9,475	19,097	29	43	Yes	No	Yes	55,000
Lahore University of Management Sciences	70	14	11,250	26,300	23	2.5	No	No	No	7,500
Universiti Sains Malaysia (Penang)	30	5	2,260	2,260	32	15	No	Yes	Yes	16,000
De La Salle University (Manila)	44	17	3,300	3,600	28	3.5	Yes	No	Yes	13,100

# **Managerial Report**

Use the methods of descriptive statistics to summarize the data in Table 3.16. Discuss your findings.

- 1. Include a summary for each variable in the data set. Make comments and interpretations based on maximums and minimums, as well as the appropriate means and proportions. What new insights do these descriptive statistics provide concerning Asia-Pacific business schools?
- 2. Summarize the data to compare the following:
  - a. Any difference between local and foreign tuition costs.
  - b. Any difference between mean starting salaries for schools requiring and not requiring work experience.
  - Any difference between starting salaries for schools requiring and not requiring English tests.
- **3.** Do starting salaries appear to be related to tuition?
- **4.** Present any additional graphical and numerical summaries that will be beneficial in communicating the data in Table 3.16 to others.

# Appendix 3.1 Descriptive Statistics Using Minitab

In this appendix, we describe how to use Minitab to develop descriptive statistics. Table 3.1 listed the starting salaries for 12 business school graduates. Panel A of Figure 3.11 shows the descriptive statistics obtained by using Minitab to summarize these data. Definitions of the headings in panel A follow.

number of data values N\* number of missing data values Mean SE Mean standard error of mean standard deviation StDev minimum data value Minimum Q1 first quartile Median median Q3 third quartile Maximum maximum data value

The label SE Mean refers to the *standard error of the mean*. It is computed by dividing the standard deviation by the square root of N. The interpretation and use of this measure are discussed in Chapter 7 when we introduce the topics of sampling and sampling distributions.

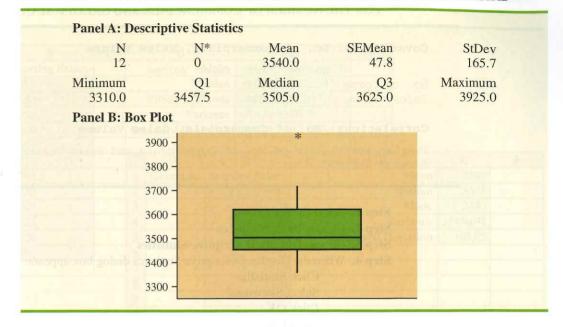
Although the numerical measures of range, interquartile range, variance, and coefficient of variation do not appear on the Minitab output, these values can be easily computed from the results in Figure 3.11 as follows.

Range = Maximum - Minimum IQR = 
$$Q_3 - Q_1$$
 Variance =  $(StDev)^2$  Coefficient of Variation =  $(StDev/Mean) \times 100$ 

Finally, note that Minitab's quartiles  $Q_1 = 3457.5$  and  $Q_3 = 3625$  are slightly different from the quartiles  $Q_1 = 3465$  and  $Q_3 = 3600$  computed in Section 3.1. The different

STUDENTS-HUB.com

#### FIGURE 3.11 DESCRIPTIVE STATISTICS AND BOX PLOT PROVIDED BY MINITAB





conventions\* used to identify the quartiles explain this variation. Hence, the values of  $Q_1$  and  $Q_3$  provided by one convention may not be identical to the values of  $Q_1$  and  $Q_3$  provided by another convention. Any differences tend to be negligible, however, and the results provided should not mislead the user in making the usual interpretations associated with quartiles.

Let us now see how the statistics in Figure 3.11 are generated. The starting salary data are in column C2 of a Minitab worksheet. The following steps can then be used to generate the descriptive statistics.

- Step 1. Select the Stat menu
- Step 2. Choose Basic Statistics
- Step 3. Choose Display Descriptive Statistics
- Step 4. When the Display Descriptive Statistics dialog box appears:

Enter C2 in the Variables box

Click OK

Panel B of Figure 3.11 is a box plot provided by Minitab. The box drawn from the first to third quartiles contains the middle 50% of the data. The line within the box locates the median. The asterisk indicates an outlier at 3925.

The following steps generate the box plot shown in panel B of Figure 3.11.

- Step 1. Select the Graph menu
- Step 2. Choose Boxplot
- Step 3. Select Simple and click OK
- Step 4. When the Boxplot-One Y, Simple dialog box appears:

Enter C2 in the Graph variables box

Click OK

The skewness measure does not appear as part of Minitab's standard descriptive statistics output. However, we can include it in the descriptive statistics display by following these steps.

<sup>\*</sup>With the n observations arranged in ascending order (smallest value to largest value), Minitab uses the positions given by (n + 1)/4 and 3(n + 1)/4 to locate  $Q_1$  and  $Q_3$ , respectively. When a position is fractional, Minitab interpolates between the two adjacent ordered data values to determine the corresponding quartile.

# FIGURE 3.12 COVARIANCE AND CORRELATION PROVIDED BY MINITAB FOR THE NUMBER OF COMMERCIALS AND SALES DATA

Covariances: No. of Commercials, Sales Volume

No. of Comme Sales Volume

No. of Comme 2.22222

Sales Volume 11.00000

Correlations: No. of Commercials, Sales Volume

Pearson correlation of No. of Commercials and Sales Volume = 0.930 P-Value = 0.000

62.88889

Step 1. Select the Stat menu

Step 2. Choose Basic Statistics

Step 3. Choose Display Descriptive Statistics

Step 4. When the Display Descriptive Statistics dialog box appears:

Click Statistics

Select Skewness

Click OK

Click OK

The skewness measure of 1.09 will then appear in your worksheet.



Figure 3.12 shows the covariance and correlation output that Minitab provided for the stereo and sound equipment store data in Table 3.7. In the covariance portion of the figure, No. of Comme denotes the number of weekend television commercials and Sales Volume denotes the sales during the following week. The value in column No. of Comme and row Sales Volume, 11, is the sample covariance as computed in Section 3.5. The value in column No. of Comme and row No. of Comme, 2.22222, is the sample variance for the number of commercials, and the value in column Sales Volume and row Sales Volume, 62.88889, is the sample variance for sales. The sample correlation coefficient, 0.930, is shown in the correlation portion of the output. Note: The interpretation of the *p*-value = 0.000 is discussed in Chapter 9.

Let us now describe how to obtain the information in Figure 3.12. We entered the data for the number of commercials into column C2 and the data for sales volume into column C3 of a Minitab worksheet. The steps necessary to generate the covariance output in Figure 3.12 follow.

Step 1. Select the Stat menu

Step 2. Choose Basic Statistics

Step 3. Choose Covariance

**Step 4.** When the Covariance dialog box appears:

Enter C2 C3 in the Variables box

Click OK

To obtain the correlation output in Figure 3.12, only one change is necessary in the steps for obtaining the covariance. In step 3, the **Correlation** option is selected.

# Appendix 3.2 Descriptive Statistics Using Excel

Excel can be used to generate the descriptive statistics discussed in this chapter. We show how Excel can be used to generate several measures of location and variability for a single variable and to generate the covariance and correlation coefficient as measures of association between two variables.

STUDENTS-HUB.com

FIGURE 3.13 USING EXCEL FUNCTIONS FOR COMPUTING THE MEAN, MEDIAN, MODE, VARIANCE, AND STANDARD DEVIATION

1	A	В	C	H.B	D	reżynici	inn Covar	E	(File	F			
1	Graduate	Starting Salary			- 1	Mean	=AVERAC	GE(B2:I	B13)				
2	1	3450			M	edian	=MEDIAN	N(B2:B)	13)				
3	2	3550			1	Mode	=MODE(E	2:B13)					
4	3	3650			Var	iance	=VAR(B2:	B13)					
5	4	3480		Sta	ındard Devi	ation	=STDEV(	B2:B13	)				
6	5	3355			~								
7	6	3310		4	A		В	C		D		E	F
8	7	3490		1	Graduate	Start	ing Salary				Mean	3540	
9	8	3730		2	- 1		3450				Median	3505	
10	9	3540		3	2		3550				Mode	3480	
11	10	3925		4	3		3650			7	ariance	27440.91	
12	11	3520		5	4		3480		Stan	ndard D	eviation	165.65	
13	12	3480		6	5		3355						
14				7	6		3310						
				8	7		3490						
				9	8		3730						
				10	9		3540						
				11	10		3925						
				12	11		3520						
				13	12		3480						
				14									

# **Using Excel Functions**



Excel provides functions for computing the mean, median, mode, sample variance, and sample standard deviation. We illustrate the use of these Excel functions by computing the mean, median, mode, sample variance, and sample standard deviation for the starting salary data in Table 3.1. Refer to Figure 3.13 as we describe the steps involved. The data are entered in column B.

Excel's AVERAGE function can be used to compute the mean by entering the following formula into cell E1:

=AVERAGE(B2:B13)

Similarly, the formulas =MEDIAN(B2:B13), =MODE(B2:B13), =VAR(B2:B13), and =STDEV(B2:B13) are entered into cells E2:E5, respectively, to compute the median, mode, variance, and standard deviation. The worksheet in the foreground shows that the values computed using the Excel functions are the same as we computed earlier in the chapter.

Excel also provides functions that can be used to compute the covariance and correlation coefficient. You must be careful when using these functions because the covariance function treats the data as a population and the correlation function treats the data as a sample. Thus, the result obtained using Excel's covariance function must be adjusted to provide the sample covariance. We show here how these functions can be used to compute the sample covariance and the sample correlation coefficient for the stereo and sound equipment store data in Table 3.7. Refer to Figure 3.14 as we present the steps involved.



# FIGURE 3.14 USING EXCEL FUNCTIONS FOR COMPUTING COVARIANCE AND CORRELATION

W	A	В	C	D			E		F		G		
1	Week	Commercials	Sales		Po	pulation				311,C2:C11)			
2	1	2	50			Sample	Correlation	=CORR	EL(B2	:B11,C2:C11)			
3	2	5	57				ILPA JED						~
4	3	1	41		A	A	В	C	D	E		F	G
5	4	3	54		1	Week	Commercials	Sales		Population (		9.90	
6	5	4	54		2	1	2	50		Sample Correlation		0.93	
7	6	1	38		3	2	5	57					
8	7	5	63		4	3	1	41					
9	8	3	48		5	4	3	54					
10	9	4	59		6	5	4	54					
11	10	2	46		7	6	1	38					
12					8	7	5	63					
					9	8	3	48					
					10	9	4	59					
					11	10	2	46					
					12								

Excel's covariance function, COVAR, can be used to compute the population covariance by entering the following formula into cell F1:

#### =COVAR(B2:B11,C2:C11)

Similarly, the formula =CORREL(B2:B11,C2:C11) is entered into cell F2 to compute the sample correlation coefficient. The worksheet in the foreground shows the values computed using the Excel functions. Note that the value of the sample correlation coefficient (.93) is the same as computed using equation (3.12). However, the result provided by the Excel COVAR function, 9.9, was obtained by treating the data as a population. Thus, we must adjust the Excel result of 9.9 to obtain the sample covariance. The adjustment is rather simple. First, note that the formula for the population covariance, equation (3.11), requires dividing by the total number of observations in the data set. But the formula for the sample covariance, equation (3.10), requires dividing by the total number of observations minus 1. So, to use the Excel result of 9.9 to compute the sample covariance, we simply multiply 9.9 by n/(n-1). Because n=10, we obtain

$$s_{xy} = \left(\frac{10}{9}\right)9.9 = 11$$

Thus, the sample covariance for the stereo and sound equipment data is 11.

# **Using Excel's Descriptive Statistics Tool**

As we already demonstrated, Excel provides statistical functions to compute descriptive statistics for a data set. These functions can be used to compute one statistic at a time (e.g., mean, variance, etc.). Excel also provides a variety of Data Analysis Tools. One of these tools, called Descriptive Statistics, allows the user to compute a variety of

# STUDENTS-HUB.com

#### FIGURE 3.15 EXCEL'S DESCRIPTIVE STATISTICS TOOL OUTPUT

A	A	В	C	D	E	F
1	Graduate	Starting Salary		Starting Sala	ry	
2	1	3450				
3	2	3550		Mean	3540	
4	3	3650		Standard Error	47.82	
5	4	3480		Median	3505	
6	5	3355		Mode	3480	
7	6	3310		<b>Standard Deviation</b>	165.65	
8	7	3490		Sample Variance	27440.91	
9	8	3730	98	Kurtosis	1.7189	
10	9	3540		Skewness	1.0911	3
11	10	3925		Range	615	
12	11	3520		Minimum	3310	
13	12	3480		Maximum	3925	
14				Sum	42480	
15				Count	12	
16						



descriptive statistics at once. We show here how it can be used to compute descriptive statistics for the starting salary data in Table 3.1. Refer to Figure 3.15 as we describe the steps involved.

- Step 1. Click the Data tab on the Ribbon
- Step 2. In the Analysis group, click Data Analysis
- **Step 3.** When the Data Analysis dialog box appears:

**Choose Descriptive Statistics** 

Click OK

**Step 4.** When the Descriptive Statistics dialog box appears:

Enter B1:B13 in the Input Range box

**Select Grouped By Columns** 

Select Labels in First Row

Select Output Range

Enter D1 in the **Output Range** box (to identify the upper left-hand corner of the section of the worksheet where the descriptive statistics will appear)

Select Summary statistics

Click OK

Cells D1:E15 of Figure 3.15 show the descriptive statistics provided by Excel. The boldface entries are the descriptive statistics we covered in this chapter. The descriptive statistics that are not boldface are either covered subsequently in the text or discussed in more advanced texts.