

CHAPTER 2



Descriptive Statistics: Tabular and Graphical Presentations

CONTENTS

STATISTICS IN PRACTICE:
COLGATE-PALMOLIVE COMPANY

2.1 SUMMARIZING
CATEGORICAL DATA
Frequency Distribution
Relative Frequency and Percent
Frequency Distributions
Bar Charts and Pie Charts

2.2 SUMMARIZING
QUANTITATIVE DATA
Frequency Distribution
Relative Frequency and Percent
Frequency Distributions

Dot Plot
Histogram
Cumulative Distributions
Ogive

2.3 EXPLORATORY DATA
ANALYSIS: THE STEM-AND-
LEAF DISPLAY

2.4 CROSSTABULATIONS AND
SCATTER DIAGRAMS
Crosstabulation
Simpson's Paradox
Scatter Diagram and Trendline



STATISTICS *in* PRACTICE

COLGATE-PALMOLIVE COMPANY*

NEW YORK, NEW YORK

The Colgate-Palmolive Company started as a small soap and candle shop in New York City in 1806. Today, Colgate-Palmolive employs more than 40,000 people working in more than 200 countries and territories around the world. Although best known for its brand names of Colgate, Palmolive, Ajax, and Fab, the company also markets Mennen, Hill's Science Diet, and Hill's Prescription Diet products.

The Colgate-Palmolive Company uses statistics in its quality assurance program for home laundry detergent products. One concern is customer satisfaction with the quantity of detergent in a carton. Every carton in each size category is filled with the same amount of detergent by weight, but the volume of detergent is affected by the density of the detergent powder. For instance, if the powder density is on the heavy side, a smaller volume of detergent is needed to reach the carton's specified weight. As a result, the carton may appear to be under-filled when opened by the consumer.

To control the problem of heavy detergent powder, limits are placed on the acceptable range of powder density. Statistical samples are taken periodically, and the density of each powder sample is measured. Data summaries are then provided for operating personnel so that corrective action can be taken if necessary to keep the density within the desired quality specifications.

A frequency distribution for the densities of 150 samples taken over a one-week period and a histogram are shown in the accompanying table and figure. Density levels above .40 are unacceptably high. The frequency distribution and histogram show that the operation is meeting its quality guidelines with all of the densities less than or equal to .40. Managers viewing these statistical summaries would be pleased with the quality of the detergent production process.

In this chapter, you will learn about tabular and graphical methods of descriptive statistics such as frequency distributions, bar charts, histograms, stem-and-leaf displays, crosstabulations, and others. The goal of



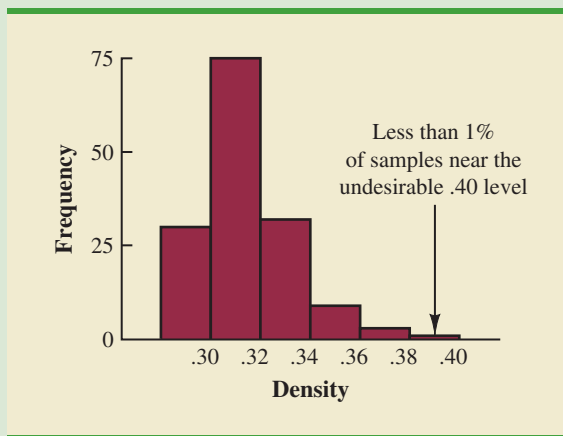
Graphical summaries help track the demand for Colgate-Palmolive products. © Victor Fisher/ Bloomberg News/Landov.

these methods is to summarize data so that the data can be easily understood and interpreted.

Frequency Distribution of Density Data

Density	Frequency
.29–.30	30
.31–.32	75
.33–.34	32
.35–.36	9
.37–.38	3
.39–.40	1
Total	150

Histogram of Density Data



*The authors are indebted to William R. Fowle, Manager of Quality Assurance, Colgate-Palmolive Company, for providing this Statistics in Practice.

As indicated in Chapter 1, data can be classified as either categorical or quantitative. **Categorical data** use labels or names to identify categories of like items. **Quantitative data** are numerical values that indicate how much or how many.

This chapter introduces tabular and graphical methods commonly used to summarize both categorical and quantitative data. Tabular and graphical summaries of data can be found in annual reports, newspaper articles, and research studies. Everyone is exposed to these types of presentations. Hence, it is important to understand how they are prepared and how they should be interpreted. We begin with tabular and graphical methods for summarizing data concerning a single variable. The last section introduces methods for summarizing data when the relationship between two variables is of interest.

Modern statistical software packages provide extensive capabilities for summarizing data and preparing graphical presentations. Minitab and Excel are two packages that are widely available. In the chapter appendixes, we show some of their capabilities.

2.1

Summarizing Categorical Data

Frequency Distribution

We begin the discussion of how tabular and graphical methods can be used to summarize categorical data with the definition of a **frequency distribution**.

FREQUENCY DISTRIBUTION

A frequency distribution is a tabular summary of data showing the number (frequency) of items in each of several nonoverlapping classes.

Let us use the following example to demonstrate the construction and interpretation of a frequency distribution for categorical data. Coke Classic, Diet Coke, Dr. Pepper, Pepsi, and Sprite are five popular soft drinks. Assume that the data in Table 2.1 show the soft drink selected in a sample of 50 soft drink purchases.

TABLE 2.1 DATA FROM A SAMPLE OF 50 SOFT DRINK PURCHASES

Coke Classic	Sprite	Pepsi
Diet Coke	Coke Classic	Coke Classic
Pepsi	Diet Coke	Coke Classic
Diet Coke	Coke Classic	Coke Classic
Coke Classic	Diet Coke	Pepsi
Coke Classic	Coke Classic	Dr. Pepper
Dr. Pepper	Sprite	Coke Classic
Diet Coke	Pepsi	Diet Coke
Pepsi	Coke Classic	Pepsi
Pepsi	Coke Classic	Pepsi
Coke Classic	Coke Classic	Pepsi
Dr. Pepper	Pepsi	Pepsi
Sprite	Coke Classic	Coke Classic
Coke Classic	Sprite	Dr. Pepper
Diet Coke	Dr. Pepper	Pepsi
Coke Classic	Pepsi	Sprite
Coke Classic	Diet Coke	



TABLE 2.2

FREQUENCY DISTRIBUTION OF SOFT DRINK PURCHASES	
Soft Drink	Frequency
Coke Classic	19
Diet Coke	8
Dr. Pepper	5
Pepsi	13
Sprite	<u>5</u>
Total	50

To develop a frequency distribution for these data, we count the number of times each soft drink appears in Table 2.1. Coke Classic appears 19 times, Diet Coke appears 8 times, Dr. Pepper appears 5 times, Pepsi appears 13 times, and Sprite appears 5 times. These counts are summarized in the frequency distribution in Table 2.2.

This frequency distribution provides a summary of how the 50 soft drink purchases are distributed across the five soft drinks. This summary offers more insight than the original data shown in Table 2.1. Viewing the frequency distribution, we see that Coke Classic is the leader, Pepsi is second, Diet Coke is third, and Sprite and Dr. Pepper are tied for fourth. The frequency distribution summarizes information about the popularity of the five soft drinks.

Relative Frequency and Percent Frequency Distributions

A frequency distribution shows the number (frequency) of items in each of several nonoverlapping classes. However, we are often interested in the proportion, or percentage, of items in each class. The *relative frequency* of a class equals the fraction or proportion of items belonging to a class. For a data set with n observations, the relative frequency of each class can be determined as follows:

RELATIVE FREQUENCY

$$\text{Relative frequency of a class} = \frac{\text{Frequency of the class}}{n} \quad (2.1)$$

The *percent frequency* of a class is the relative frequency multiplied by 100.

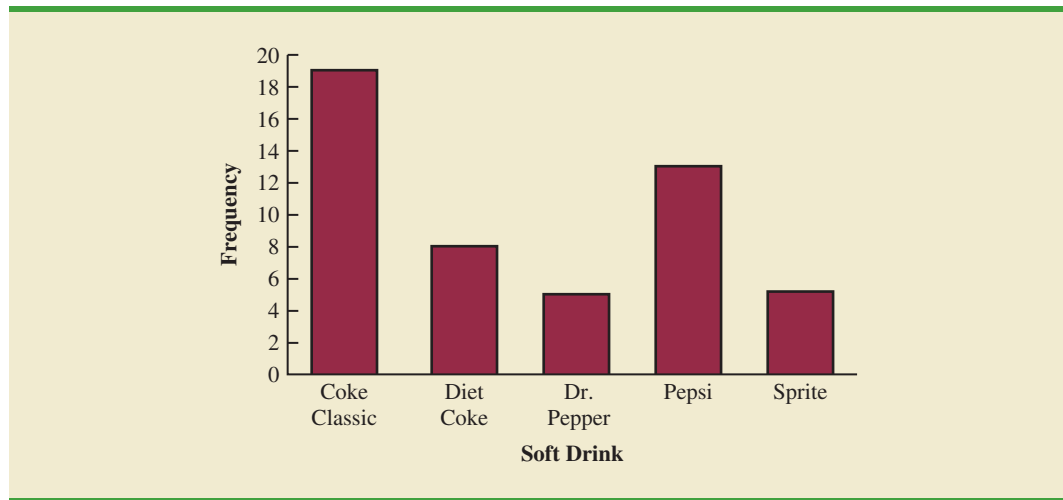
A **relative frequency distribution** gives a tabular summary of data showing the relative frequency for each class. A **percent frequency distribution** summarizes the percent frequency of the data for each class. Table 2.3 shows a relative frequency distribution and a percent frequency distribution for the soft drink data. In Table 2.3 we see that the relative frequency for Coke Classic is $19/50 = .38$, the relative frequency for Diet Coke is $8/50 = .16$, and so on. From the percent frequency distribution, we see that 38% of the purchases were Coke Classic, 16% of the purchases were Diet Coke, and so on. We can also note that $38\% + 26\% + 16\% = 80\%$ of the purchases were the top three soft drinks.

Bar Charts and Pie Charts

A **bar chart** is a graphical device for depicting categorical data summarized in a frequency, relative frequency, or percent frequency distribution. On one axis of the graph (usually the horizontal axis), we specify the labels that are used for the classes (categories). A frequency, relative frequency, or percent frequency scale can be used for the other axis of the chart

TABLE 2.3 RELATIVE FREQUENCY AND PERCENT FREQUENCY DISTRIBUTIONS OF SOFT DRINK PURCHASES

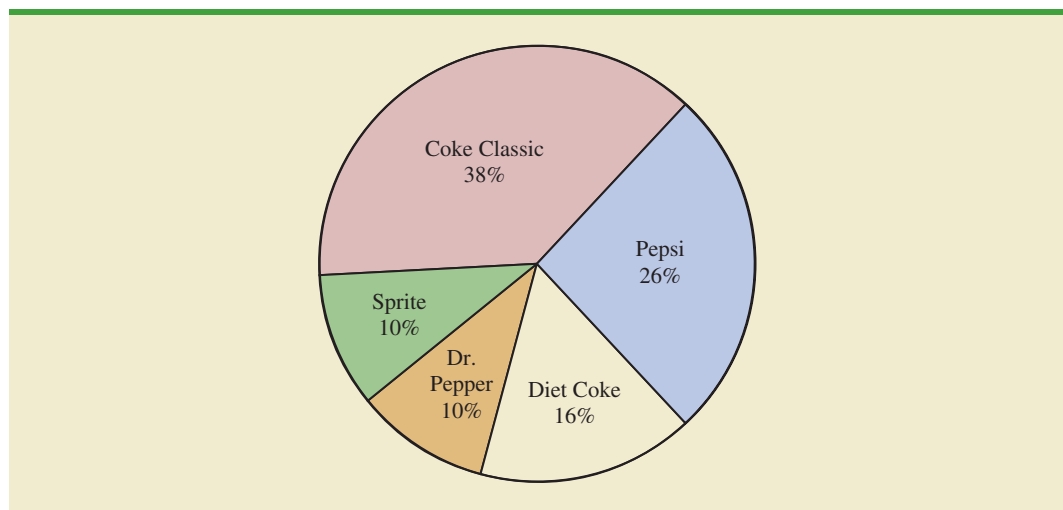
Soft Drink	Relative Frequency	Percent Frequency
Coke Classic	.38	38
Diet Coke	.16	16
Dr. Pepper	.10	10
Pepsi	.26	26
Sprite	<u>.10</u>	<u>10</u>
Total	1.00	100

FIGURE 2.1 BAR CHART OF SOFT DRINK PURCHASES

In quality control applications, bar charts are used to identify the most important causes of problems. When the bars are arranged in descending order of height from left to right with the most frequently occurring cause appearing first, the bar chart is called a pareto diagram. This diagram is named for its founder, Vilfredo Pareto, an Italian economist.

(usually the vertical axis). Then, using a bar of fixed width drawn above each class label, we extend the length of the bar until we reach the frequency, relative frequency, or percent frequency of the class. For categorical data, the bars should be separated to emphasize the fact that each class is separate. Figure 2.1 shows a bar chart of the frequency distribution for the 50 soft drink purchases. Note how the graphical presentation shows Coke Classic, Pepsi, and Diet Coke to be the most preferred brands.

The **pie chart** provides another graphical device for presenting relative frequency and percent frequency distributions for categorical data. To construct a pie chart, we first draw a circle to represent all the data. Then we use the relative frequencies to subdivide the circle into sectors, or parts, that correspond to the relative frequency for each class. For example, because a circle contains 360 degrees and Coke Classic shows a relative frequency of .38, the sector of the pie chart labeled Coke Classic consists of $.38(360) = 136.8$ degrees. The sector of the pie chart labeled Diet Coke consists of $.16(360) = 57.6$ degrees. Similar calculations for the other classes yield the pie chart in Figure 2.2. The

FIGURE 2.2 PIE CHART OF SOFT DRINK PURCHASES

numerical values shown for each sector can be frequencies, relative frequencies, or percent frequencies.

NOTES AND COMMENTS

- Often the number of classes in a frequency distribution is the same as the number of categories found in the data, as is the case for the soft drink purchase data in this section. The data involve only five soft drinks, and a separate frequency distribution class was defined for each one. Data that included all soft drinks would require many categories, most of which would have a small number of purchases. Most statisticians recommend that classes with smaller frequencies be grouped into an aggregate class called “other.” Classes with frequencies of 5% or less would most often be treated in this fashion.
- The sum of the frequencies in any frequency distribution always equals the number of observations. The sum of the relative frequencies in any relative frequency distribution always equals 1.00, and the sum of the percentages in a percent frequency distribution always equals 100.

Exercises

Methods

- The response to a question has three alternatives: A, B, and C. A sample of 120 responses provides 60 A, 24 B, and 36 C. Show the frequency and relative frequency distributions.
- A partial relative frequency distribution is given.

Class	Relative Frequency
A	.22
B	.18
C	.40
D	

- What is the relative frequency of class D?
 - The total sample size is 200. What is the frequency of class D?
 - Show the frequency distribution.
 - Show the percent frequency distribution.
- A questionnaire provides 58 Yes, 42 No, and 20 no-opinion answers.
 - In the construction of a pie chart, how many degrees would be in the section of the pie showing the Yes answers?
 - How many degrees would be in the section of the pie showing the No answers?
 - Construct a pie chart.
 - Construct a bar chart.

SELF test

Applications

- The top four prime-time television shows were *Law & Order*, *CSI*, *Without a Trace*, and *Desperate Housewives* (Nielsen Media Research, January 1, 2007). Data indicating the preferred shows for a sample of 50 viewers follow.

DH	CSI	DH	CSI	L&O
Trace	CSI	L&O	Trace	CSI
CSI	DH	Trace	CSI	DH
L&O	L&O	L&O	CSI	DH
CSI	DH	DH	L&O	CSI
DH	Trace	CSI	Trace	DH
DH	CSI	CSI	L&O	CSI
L&O	CSI	Trace	Trace	DH
L&O	CSI	CSI	CSI	DH
CSI	DH	Trace	Trace	L&O

- Are these data categorical or quantitative?
 - Provide frequency and percent frequency distributions.
 - Construct a bar chart and a pie chart.
 - On the basis of the sample, which television show has the largest viewing audience? Which one is second?
5. In alphabetical order, the six most common last names in the United States are Brown, Davis, Johnson, Jones, Smith, and Williams (*The World Almanac*, 2006). Assume that a sample of 50 individuals with one of these last names provided the following data.

WEB file
Names

Brown	Williams	Williams	Williams	Brown
Smith	Jones	Smith	Johnson	Smith
Davis	Smith	Brown	Williams	Johnson
Johnson	Smith	Smith	Johnson	Brown
Williams	Davis	Johnson	Williams	Johnson
Williams	Johnson	Jones	Smith	Brown
Johnson	Smith	Smith	Brown	Jones
Jones	Jones	Smith	Smith	Davis
Davis	Jones	Williams	Davis	Smith
Jones	Johnson	Brown	Johnson	Davis

Summarize the data by constructing the following:

- Relative and percent frequency distributions
 - A bar chart
 - A pie chart
 - Based on these data, what are the three most common last names?
6. The Nielsen Media Research television rating measures the percentage of television owners who are watching a particular television program. The highest-rated television program in television history was the *M*A*S*H Last Episode Special* shown on February 28, 1983. A 60.2 rating indicated that 60.2% of all television owners were watching this program. Nielsen Media Research provided the list of the 50 top-rated single shows in television history (*The New York Times Almanac*, 2006). The following data show the television network that produced each of these 50 top-rated shows.

WEB file
Networks

ABC	ABC	ABC	NBC	CBS
ABC	CBS	ABC	ABC	NBC
NBC	NBC	CBS	ABC	NBC
CBS	ABC	CBS	NBC	ABC
CBS	NBC	NBC	CBS	NBC
CBS	CBS	CBS	NBC	NBC
FOX	CBS	CBS	ABC	NBC
ABC	ABC	CBS	NBC	NBC
NBC	CBS	NBC	CBS	CBS
ABC	CBS	ABC	NBC	ABC

- Construct a frequency distribution, percent frequency distribution, and bar chart for the data.

SELF test

- b. Which network or networks have done the best in terms of presenting top-rated television shows? Compare the performance of ABC, CBS, and NBC.
7. Leverock's Waterfront Steakhouse in Maderia Beach, Florida, uses a questionnaire to ask customers how they rate the server, food quality, cocktails, prices, and atmosphere at the restaurant. Each characteristic is rated on a scale of outstanding (O), very good (V), good (G), average (A), and poor (P). Use descriptive statistics to summarize the following data collected on food quality. What is your feeling about the food quality ratings at the restaurant?

G	O	V	G	A	O	V	O	V	G	O	V	A
V	O	P	V	O	G	A	O	O	O	G	O	V
V	A	G	O	V	P	V	O	O	G	O	O	V
O	G	A	O	V	O	O	G	V	A	G		

8. Data for a sample of 55 members of the Baseball Hall of Fame in Cooperstown, New York, are shown here. Each observation indicates the primary position played by the Hall of Famers: pitcher (P), catcher (H), 1st base (1), 2nd base (2), 3rd base (3), shortstop (S), left field (L), center field (C), and right field (R).

L	P	C	H	2	P	R	1	S	S	1	L	P	R	P
P	P	P	R	C	S	L	R	P	C	C	P	P	R	P
2	3	P	H	L	P	1	C	P	P	P	S	1	L	R
R	1	2	H	S	3	H	2	L	P					

- a. Use frequency and relative frequency distributions to summarize the data.
- b. What position provides the most Hall of Famers?
- c. What position provides the fewest Hall of Famers?
- d. What outfield position (L, C, or R) provides the most Hall of Famers?
- e. Compare infielders (1, 2, 3, and S) to outfielders (L, C, and R).
9. The Pew Research Center's Social & Demographic Trends project found that 46% of U.S. adults would rather live in a different type of community than the one where they are living now (Pew Research Center, January 29, 2009). The national survey of 2260 adults asked: "Where do you live now?" and "What do you consider to be the ideal community?" Response options were City (C), Suburb (S), Small Town (T), or Rural (R). A representative portion of this survey for a sample of 100 respondents is as follows.

Where do you live now?

S	T	R	C	R	R	T	C	S	T	C	S	C	S	T
S	S	C	S	S	T	T	C	C	S	T	C	S	T	C
T	R	S	S	T	C	S	C	T	C	T	C	T	C	R
C	C	R	T	C	S	S	T	S	C	C	C	R	S	C
S	S	C	C	S	C	R	T	T	T	C	R	T	C	R
C	T	R	R	C	T	C	C	R	T	T	R	S	R	T
T	S	S	S	S	S	C	C	R	T					

What do you consider to be the ideal community?

S	C	R	R	R	S	T	S	S	T	T	S	C	S	T
C	C	R	T	R	S	T	T	S	S	C	C	T	T	S
S	R	C	S	C	C	S	C	R	C	T	S	R	R	R
C	T	S	T	T	T	R	R	S	C	C	R	R	S	S
S	T	C	T	T	C	R	T	T	T	C	T	T	R	R
C	S	R	T	C	T	C	C	T	T	T	R	C	R	T
T	C	S	S	C	S	T	S	S	R					

- a. Provide a percent frequency distribution for each question.
- b. Construct a bar chart for each question.
- c. Where are most adults living now?
- d. Where do most adults consider the ideal community?

WEB file
LivingArea



- e. What changes in living areas would you expect to see if people moved from where they currently live to their ideal community?
10. The *Financial Times*/Harris Poll is a monthly online poll of adults from six countries in Europe and the United States. The poll conducted in January 2008 included 1015 adults. One of the questions asked was, “How would you rate the Federal Bank in handling the credit problems in the financial markets?” Possible responses were Excellent, Good, Fair, Bad, and Terrible (Harris Interactive website, January 2008). The 1015 responses for this question can be found in the data file named FedBank.
- Construct a frequency distribution.
 - Construct a percent frequency distribution.
 - Construct a bar chart for the percent frequency distribution.
 - Comment on how adults in the United States think the Federal Bank is handling the credit problems in the financial markets.
 - In Spain, 1114 adults were asked, “How would you rate the European Central Bank in handling the credit problems in the financial markets?” The percent frequency distribution obtained follows:

Rating	Percent Frequency
Excellent	0
Good	4
Fair	46
Bad	40
Terrible	10

Compare the results obtained in Spain with the results obtained in the United States.

2.2

Summarizing Quantitative Data

Frequency Distribution

TABLE 2.4

YEAR-END AUDIT TIMES (IN DAYS)			
12	14	19	18
15	15	18	17
20	27	22	23
22	21	33	28
14	18	16	13

As defined in Section 2.1, a frequency distribution is a tabular summary of data showing the number (frequency) of items in each of several nonoverlapping classes. This definition holds for quantitative as well as qualitative data. However, with quantitative data we must be more careful in defining the nonoverlapping classes to be used in the frequency distribution.

For example, consider the quantitative data in Table 2.4. These data show the time in days required to complete year-end audits for a sample of 20 clients of Sanderson and Clifford, a small public accounting firm. The three steps necessary to define the classes for a frequency distribution with quantitative data are:

- Determine the number of nonoverlapping classes.
- Determine the width of each class.
- Determine the class limits.



Let us demonstrate these steps by developing a frequency distribution for the audit time data in Table 2.4.

Number of classes Classes are formed by specifying ranges that will be used to group the data. As a general guideline, we recommend using between 5 and 20 classes. For a small number of data items, as few as five or six classes may be used to summarize the data. For a larger number of data items, a larger number of classes is usually required. The goal is to use enough classes to show the variation in the data, but not so many classes that some contain only a few data items. Because the number of data items in Table 2.4 is relatively small ($n = 20$), we chose to develop a frequency distribution with five classes.

Making the classes the same width reduces the chance of inappropriate interpretations by the user.

Width of the classes The second step in constructing a frequency distribution for quantitative data is to choose a width for the classes. As a general guideline, we recommend that the width be the same for each class. Thus the choices of the number of classes and the width of classes are not independent decisions. A larger number of classes means a smaller class width, and vice versa. To determine an approximate class width, we begin by identifying the largest and smallest data values. Then, with the desired number of classes specified, we can use the following expression to determine the approximate class width.

$$\text{Approximate class width} = \frac{\text{Largest data value} - \text{Smallest data value}}{\text{Number of classes}} \quad (2.2)$$

The approximate class width given by equation (2.2) can be rounded to a more convenient value based on the preference of the person developing the frequency distribution. For example, an approximate class width of 9.28 might be rounded to 10 simply because 10 is a more convenient class width to use in presenting a frequency distribution.

For the data involving the year-end audit times, the largest data value is 33 and the smallest data value is 12. Because we decided to summarize the data with five classes, using equation (2.2) provides an approximate class width of $(33 - 12)/5 = 4.2$. We therefore decided to round up and use a class width of five days in the frequency distribution.

In practice, the number of classes and the appropriate class width are determined by trial and error. Once a possible number of classes is chosen, equation (2.2) is used to find the approximate class width. The process can be repeated for a different number of classes. Ultimately, the analyst uses judgment to determine the combination of the number of classes and class width that provides the best frequency distribution for summarizing the data.

For the audit time data in Table 2.4, after deciding to use five classes, each with a width of five days, the next task is to specify the class limits for each of the classes.

Class limits Class limits must be chosen so that each data item belongs to one and only one class. The *lower class limit* identifies the smallest possible data value assigned to the class. The *upper class limit* identifies the largest possible data value assigned to the class. In developing frequency distributions for qualitative data, we did not need to specify class limits because each data item naturally fell into a separate class. But with quantitative data, such as the audit times in Table 2.4, class limits are necessary to determine where each data value belongs.

Using the audit time data in Table 2.4, we selected 10 days as the lower class limit and 14 days as the upper class limit for the first class. This class is denoted 10–14 in Table 2.5. The smallest data value, 12, is included in the 10–14 class. We then selected 15 days as the lower class limit and 19 days as the upper class limit of the next class. We continued defining the lower and upper class limits to obtain a total of five classes: 10–14, 15–19, 20–24, 25–29, and 30–34. The largest data value, 33, is included in the 30–34 class. The difference between the lower class limits of adjacent classes is the class width. Using the first two lower class limits of 10 and 15, we see that the class width is $15 - 10 = 5$.

With the number of classes, class width, and class limits determined, a frequency distribution can be obtained by counting the number of data values belonging to each class. For example, the data in Table 2.4 show that four values—12, 14, 14, and 13—belong to the 10–14 class. Thus, the frequency for the 10–14 class is 4. Continuing this counting process for the 15–19, 20–24, 25–29, and 30–34 classes provides the frequency distribution in Table 2.5. Using this frequency distribution, we can observe the following:

1. The most frequently occurring audit times are in the class of 15–19 days. Eight of the 20 audit times belong to this class.
2. Only one audit required 30 or more days.

Other conclusions are possible, depending on the interests of the person viewing the frequency distribution. The value of a frequency distribution is that it provides insights about the data that are not easily obtained by viewing the data in their original unorganized form.

No single frequency distribution is best for a data set. Different people may construct different, but equally acceptable, frequency distributions. The goal is to reveal the natural grouping and variation in the data.

TABLE 2.5
FREQUENCY
DISTRIBUTION
FOR THE AUDIT
TIME DATA

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

TABLE 2.6 RELATIVE FREQUENCY AND PERCENT FREQUENCY DISTRIBUTIONS FOR THE AUDIT TIME DATA

Audit Time (days)	Relative Frequency	Percent Frequency
10–14	.20	20
15–19	.40	40
20–24	.25	25
25–29	.10	10
30–34	.05	5
Total	1.00	100

Class midpoint In some applications, we want to know the midpoints of the classes in a frequency distribution for quantitative data. The **class midpoint** is the value halfway between the lower and upper class limits. For the audit time data, the five class midpoints are 12, 17, 22, 27, and 32.

Relative Frequency and Percent Frequency Distributions

We define the relative frequency and percent frequency distributions for quantitative data in the same manner as for qualitative data. First, recall that the relative frequency is the proportion of the observations belonging to a class. With n observations,

$$\text{Relative frequency of class} = \frac{\text{Frequency of the class}}{n}$$

The percent frequency of a class is the relative frequency multiplied by 100.

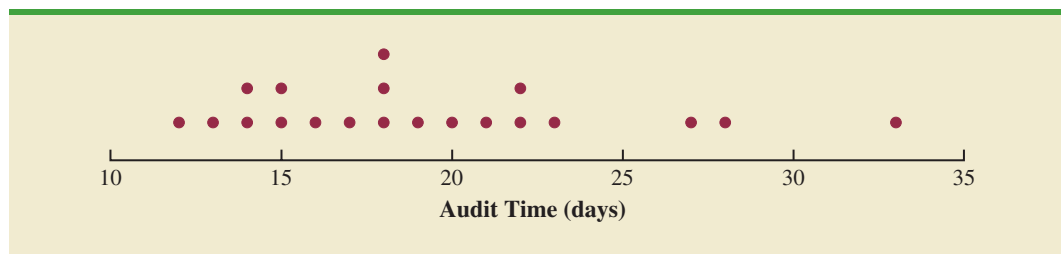
Based on the class frequencies in Table 2.5 and with $n = 20$, Table 2.6 shows the relative frequency distribution and percent frequency distribution for the audit time data. Note that .40 of the audits, or 40%, required from 15 to 19 days. Only .05 of the audits, or 5%, required 30 or more days. Again, additional interpretations and insights can be obtained by using Table 2.6.

Dot Plot

One of the simplest graphical summaries of data is a **dot plot**. A horizontal axis shows the range for the data. Each data value is represented by a dot placed above the axis. Figure 2.3 is the dot plot for the audit time data in Table 2.4. The three dots located above 18 on the horizontal axis indicate that an audit time of 18 days occurred three times. Dot plots show the details of the data and are useful for comparing the distribution of the data for two or more variables.

Histogram

A common graphical presentation of quantitative data is a **histogram**. This graphical summary can be prepared for data previously summarized in either a frequency, relative frequency, or percent frequency distribution. A histogram is constructed by placing the

FIGURE 2.3 DOT PLOT FOR THE AUDIT TIME DATA

variable of interest on the horizontal axis and the frequency, relative frequency, or percent frequency on the vertical axis. The frequency, relative frequency, or percent frequency of each class is shown by drawing a rectangle whose base is determined by the class limits on the horizontal axis and whose height is the corresponding frequency, relative frequency, or percent frequency.

Figure 2.4 is a histogram for the audit time data. Note that the class with the greatest frequency is shown by the rectangle appearing above the class of 15–19 days. The height of the rectangle shows that the frequency of this class is 8. A histogram for the relative or percent frequency distribution of these data would look the same as the histogram in Figure 2.4 with the exception that the vertical axis would be labeled with relative or percent frequency values.

As Figure 2.4 shows, the adjacent rectangles of a histogram touch one another. Unlike a bar graph, a histogram contains no natural separation between the rectangles of adjacent classes. This format is the usual convention for histograms. Because the classes for the audit time data are stated as 10–14, 15–19, 20–24, 25–29, and 30–34, one-unit spaces of 14 to 15, 19 to 20, 24 to 25, and 29 to 30 would seem to be needed between the classes. These spaces are eliminated when constructing a histogram. Eliminating the spaces between classes in a histogram for the audit time data helps show that all values between the lower limit of the first class and the upper limit of the last class are possible.

One of the most important uses of a histogram is to provide information about the shape, or form, of a distribution. Figure 2.5 contains four histograms constructed from relative frequency distributions. Panel A shows the histogram for a set of data moderately skewed to the left. A histogram is said to be skewed to the left if its tail extends farther to the left. This histogram is typical for exam scores, with no scores above 100%, most of the scores above 70%, and only a few really low scores. Panel B shows the histogram for a set of data moderately skewed to the right. A histogram is said to be skewed to the right if its tail extends farther to the right. An example of this type of histogram would be for data such as housing prices; a few expensive houses create the skewness in the right tail.

Panel C shows a symmetric histogram. In a symmetric histogram, the left tail mirrors the shape of the right tail. Histograms for data found in applications are never perfectly symmetric, but the histogram for many applications may be roughly symmetric. Data for SAT scores, heights and weights of people, and so on lead to histograms that are roughly symmetric. Panel D shows a histogram highly skewed to the right. This histogram was constructed from data on the amount of customer purchases over one day at a women's apparel store. Data from applications in business and economics often lead to histograms that

FIGURE 2.4 HISTOGRAM FOR THE AUDIT TIME DATA

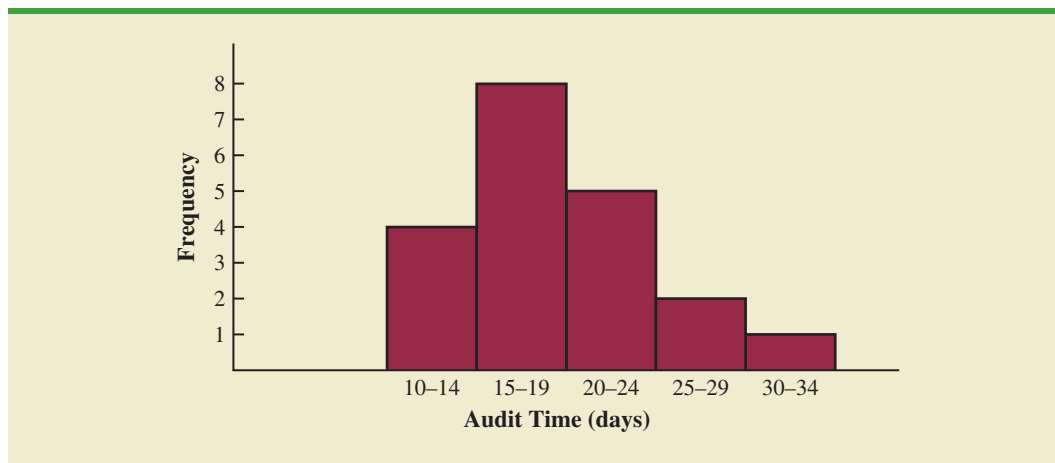
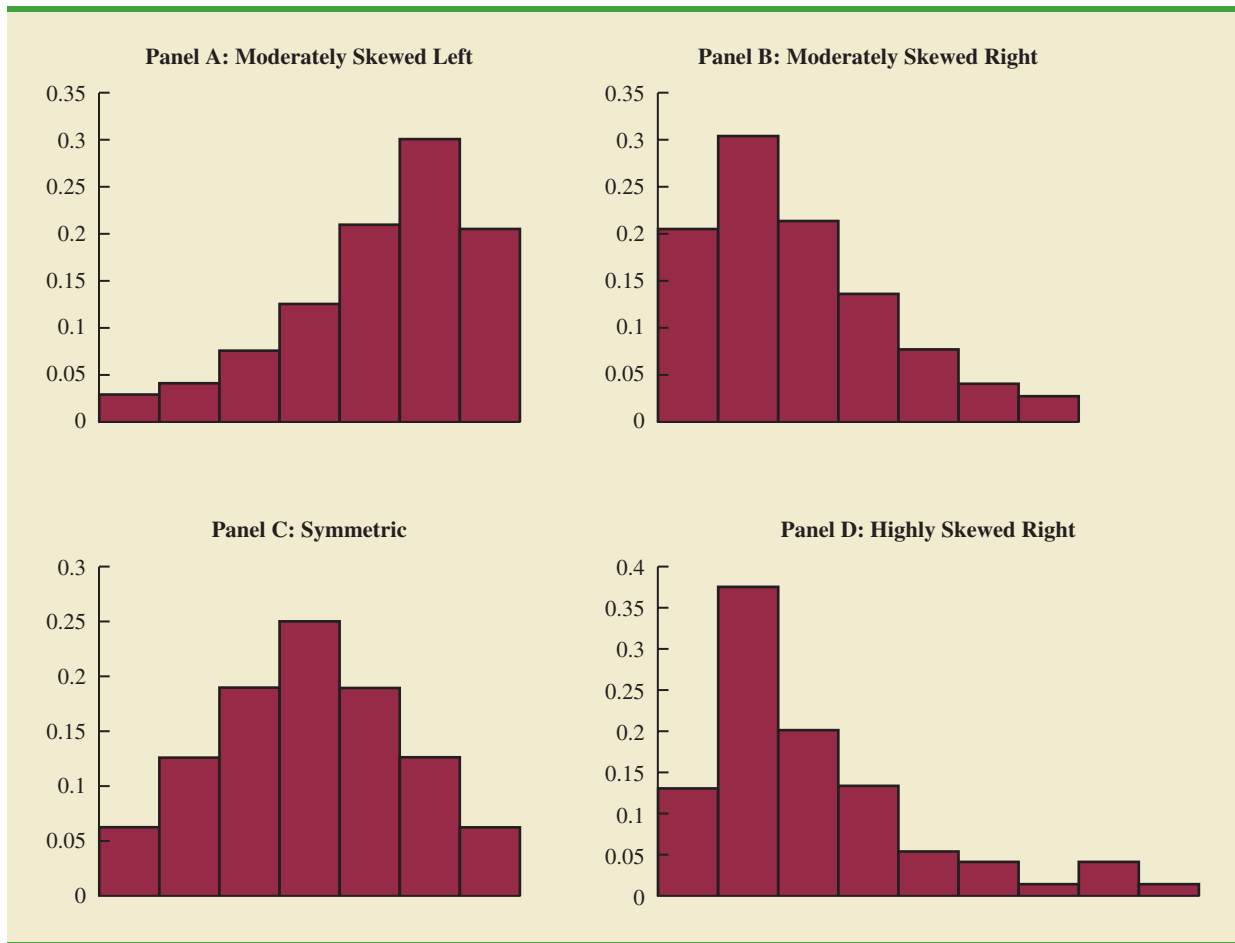


FIGURE 2.5 HISTOGRAMS SHOWING DIFFERING LEVELS OF SKEWNESS


are skewed to the right. For instance, data on housing prices, salaries, purchase amounts, and so on often result in histograms skewed to the right.

Cumulative Distributions

A variation of the frequency distribution that provides another tabular summary of quantitative data is the **cumulative frequency distribution**. The cumulative frequency distribution uses the number of classes, class widths, and class limits developed for the frequency distribution. However, rather than showing the frequency of each class, the cumulative frequency distribution shows the number of data items with values *less than or equal to the upper class limit* of each class. The first two columns of Table 2.7 provide the cumulative frequency distribution for the audit time data.

To understand how the cumulative frequencies are determined, consider the class with the description “less than or equal to 24.” The cumulative frequency for this class is simply the sum of the frequencies for all classes with data values less than or equal to 24. For the frequency distribution in Table 2.5, the sum of the frequencies for classes 10–14, 15–19, and 20–24 indicates that $4 + 8 + 5 = 17$ data values are less than or equal to 24. Hence, the cumulative frequency for this class is 17. In addition, the cumulative frequency distribution in Table 2.7 shows that four audits were completed in 14 days or less and 19 audits were completed in 29 days or less.

TABLE 2.7 CUMULATIVE FREQUENCY, CUMULATIVE RELATIVE FREQUENCY, AND CUMULATIVE PERCENT FREQUENCY DISTRIBUTIONS FOR THE AUDIT TIME DATA

Audit Time (days)	Cumulative Frequency	Cumulative Relative Frequency	Cumulative Percent Frequency
Less than or equal to 14	4	.20	20
Less than or equal to 19	12	.60	60
Less than or equal to 24	17	.85	85
Less than or equal to 29	19	.95	95
Less than or equal to 34	20	1.00	100

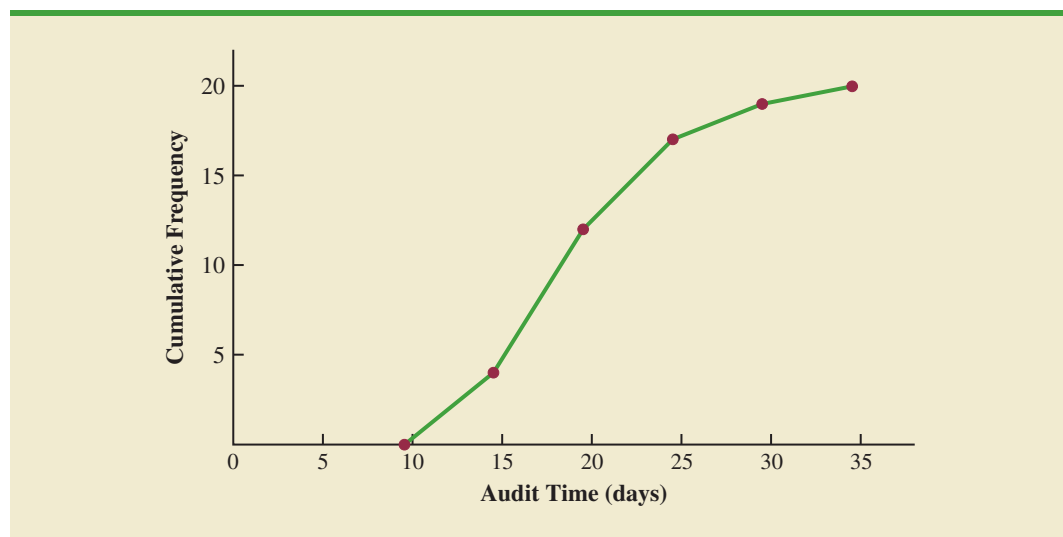
As a final point, we note that a **cumulative relative frequency distribution** shows the proportion of data items, and a **cumulative percent frequency distribution** shows the percentage of data items with values less than or equal to the upper limit of each class. The cumulative relative frequency distribution can be computed either by summing the relative frequencies in the relative frequency distribution or by dividing the cumulative frequencies by the total number of items. Using the latter approach, we found the cumulative relative frequencies in column 3 of Table 2.7 by dividing the cumulative frequencies in column 2 by the total number of items ($n = 20$). The cumulative percent frequencies were again computed by multiplying the relative frequencies by 100. The cumulative relative and percent frequency distributions show that .85 of the audits, or 85%, were completed in 24 days or less, .95 of the audits, or 95%, were completed in 29 days or less, and so on.

Ogive

A graph of a cumulative distribution, called an **ogive**, shows data values on the horizontal axis and either the cumulative frequencies, the cumulative relative frequencies, or the cumulative percent frequencies on the vertical axis. Figure 2.6 illustrates an ogive for the cumulative frequencies of the audit time data in Table 2.7.

The ogive is constructed by plotting a point corresponding to the cumulative frequency of each class. Because the classes for the audit time data are 10–14, 15–19, 20–24, and so

FIGURE 2.6 OGIVE FOR THE AUDIT TIME DATA



on, one-unit gaps appear from 14 to 15, 19 to 20, and so on. These gaps are eliminated by plotting points halfway between the class limits. Thus, 14.5 is used for the 10–14 class, 19.5 is used for the 15–19 class, and so on. The “less than or equal to 14” class with a cumulative frequency of 4 is shown on the ogive in Figure 2.6 by the point located at 14.5 on the horizontal axis and 4 on the vertical axis. The “less than or equal to 19” class with a cumulative frequency of 12 is shown by the point located at 19.5 on the horizontal axis and 12 on the vertical axis. Note that one additional point is plotted at the left end of the ogive. This point starts the ogive by showing that no data values fall below the 10–14 class. It is plotted at 9.5 on the horizontal axis and 0 on the vertical axis. The plotted points are connected by straight lines to complete the ogive.

NOTES AND COMMENTS

1. A bar chart and a histogram are essentially the same thing; both are graphical presentations of the data in a frequency distribution. A histogram is just a bar chart with no separation between bars. For some discrete quantitative data, a separation between bars is also appropriate. Consider, for example, the number of classes in which a college student is enrolled. The data may only assume integer values. Intermediate values such as 1.5, 2.73, and so on are not possible. With continuous quantitative data, however, such as the audit times in Table 2.4, a separation between bars is not appropriate.
2. The appropriate values for the class limits with quantitative data depend on the level of accuracy of the data. For instance, with the audit time data of Table 2.4 the limits used were integer values. If the data were rounded to the nearest tenth of a day (e.g., 12.3, 14.4, and so on), then the limits would be stated in tenths of days. For instance, the first class would be 10.0–14.9. If the data were recorded to the nearest hundredth of a day (e.g., 12.34, 14.45, and so on), the limits would be stated in hundredths of days. For instance, the first class would be 10.00–14.99.
3. An *open-end* class requires only a lower class limit or an upper class limit. For example, in the audit time data of Table 2.4, suppose two of the audits had taken 58 and 65 days. Rather than continue with the classes of width 5 with classes 35–39, 40–44, 45–49, and so on, we could simplify the frequency distribution to show an open-end class of “35 or more.” This class would have a frequency of 2. Most often the open-end class appears at the upper end of the distribution. Sometimes an open-end class appears at the lower end of the distribution, and occasionally such classes appear at both ends.
4. The last entry in a cumulative frequency distribution always equals the total number of observations. The last entry in a cumulative relative frequency distribution always equals 1.00 and the last entry in a cumulative percent frequency distribution always equals 100.

Exercises

Methods

11. Consider the following data.

14	21	23	21	16
19	22	25	16	16
24	24	25	19	16
19	18	19	21	12
16	17	18	23	25
20	23	16	20	19
24	26	15	22	24
20	22	24	22	20

- a. Develop a frequency distribution using classes of 12–14, 15–17, 18–20, 21–23, and 24–26.
- b. Develop a relative frequency distribution and a percent frequency distribution using the classes in part (a).

WEB file
Frequency

SELF test

12. Consider the following frequency distribution.

Class	Frequency
10–19	10
20–29	14
30–39	17
40–49	7
50–59	2

Construct a cumulative frequency distribution and a cumulative relative frequency distribution.

13. Construct a histogram and an ogive for the data in exercise 12.
14. Consider the following data.

8.9 10.2 11.5 7.8 10.0 12.2 13.5 14.1 10.0 12.2
 6.8 9.5 11.5 11.2 14.9 7.5 10.0 6.0 15.8 11.5

- Construct a dot plot.
- Construct a frequency distribution.
- Construct a percent frequency distribution.

Applications**SELF test**

15. A doctor's office staff studied the waiting times for patients who arrive at the office with a request for emergency service. The following data with waiting times in minutes were collected over a one-month period.

2 5 10 12 4 4 5 17 11 8 9 8 12 21 6 8 7 13 18 3

Use classes of 0–4, 5–9, and so on in the following:

- Show the frequency distribution.
 - Show the relative frequency distribution.
 - Show the cumulative frequency distribution.
 - Show the cumulative relative frequency distribution.
 - What proportion of patients needing emergency service wait 9 minutes or less?
16. A shortage of candidates has required school districts to pay higher salaries and offer extras to attract and retain school district superintendents. The following data show the annual base salary (\$1000s) for superintendents in 20 districts in the greater Rochester, New York, area (*The Rochester Democrat and Chronicle*, February 10, 2008).

187	184	174	185
175	172	202	197
165	208	215	164
162	172	182	156
172	175	170	183

Use classes of 150–159, 160–169, and so on in the following.

- Show the frequency distribution.
 - Show the percent frequency distribution.
 - Show the cumulative percent frequency distribution.
 - Develop a histogram for the annual base salary.
 - Do the data appear to be skewed? Explain.
 - What percentage of the superintendents make more than \$200,000?
17. The Dow Jones Industrial Average (DJIA) underwent one of its infrequent reshufflings of companies when General Motors and Citigroup were replaced by Cisco Systems and Travelers (*The Wall Street Journal*, June 8, 2009). At the time, the prices per share for the 30 companies in the DJIA were as follows:

WEB file
DJIAPrices

Company	\$/Share	Company	\$/Share
3M	61	IBM	107
Alcoa	11	Intel	16
American Express	25	J.P. Morgan Chase	35
AT&T	24	Johnson & Johnson	56
Bank of America	12	Kraft Foods	27
Boeing	52	McDonald's	59
Caterpillar	38	Merck	26
Chevron	69	Microsoft	22
Cisco Systems	20	Pfizer	14
Coca-Cola	49	Procter & Gamble	53
DuPont	27	Travelers	43
ExxonMobil	72	United Technologies	56
General Electric	14	Verizon	29
Hewlett-Packard	37	Wal-Mart Stores	51
Home Depot	24	Walt Disney	25

- What is the highest price per share? What is the lowest price per share?
- Using a class width of 10, develop a frequency distribution for the data.
- Prepare a histogram. Interpret the histogram, including a discussion of the general shape of the histogram, the midprice range, and the most frequent price range.
- Use the *The Wall Street Journal* or another newspaper to find the current price per share for these companies. Prepare a histogram of the data and discuss any changes since June 2009. What company has had the largest increase in the price per share? What company has had the largest decrease in the price per share?

18. NRF/BIG research provided results of a consumer holiday spending survey (*USA Today*, December 20, 2005). The following data provide the dollar amount of holiday spending for a sample of 25 consumers.

1200	850	740	590	340
450	890	260	610	350
1780	180	850	2050	770
800	1090	510	520	220
1450	280	1120	200	350

- What is the lowest holiday spending? The highest?
 - Use a class width of \$250 to prepare a frequency distribution and a percent frequency distribution for the data.
 - Prepare a histogram and comment on the shape of the distribution.
 - What observations can you make about holiday spending?
19. Sorting through unsolicited e-mail and spam affects the productivity of office workers. An InsightExpress survey monitored office workers to determine the unproductive time per day devoted to unsolicited e-mail and spam (*USA Today*, November 13, 2003). The following data show a sample of time in minutes devoted to this task.

2	4	8	4
8	1	2	32
12	1	5	7
5	5	3	4
24	19	4	14

Summarize the data by constructing the following:

- A frequency distribution (classes 1–5, 6–10, 11–15, 16–20, and so on)
- A relative frequency distribution
- A cumulative frequency distribution
- A cumulative relative frequency distribution
- An ogive
- What percentage of office workers spend 5 minutes or less on unsolicited e-mail and spam? What percentage of office workers spend more than 10 minutes a day on this task?

20. The *Golf Digest 50* lists the 50 professional golfers with the highest total annual income. Total income is the sum of both on-course and off-course earnings. Tiger Woods ranked first with a total annual income of \$122 million. However, almost \$100 million of this total was from off-course activities such as product endorsements and personal appearances. The 10 professional golfers with the highest *off-course* income are shown in the following table (Golf Digest website, February 2008).

Name	Off-Course Income (\$1000s)
Tiger Woods	99,800
Phil Mickelson	40,200
Arnold Palmer	29,500
Vijay Singh	25,250
Ernie Els	24,500
Greg Norman	24,000
Jack Nicklaus	20,750
Sergio Garcia	14,500
Michelle Wie	12,500
Jim Furyk	11,000

The off-course income of all 50 professional golfers in the *Golf Digest 50* can be found on the website that accompanies the text. The income data are in \$1000s. Use classes of 0–4999, 5000–9999, 10,000–14,999, and so on to answer the following questions. Include an open-ended class of 50,000 or more as the largest income class.

- a. Construct a frequency distribution and percent frequency distribution of the annual off-course income of the 50 professional golfers.
- b. Construct a histogram for these data.
- c. Comment on the shape of the distribution of off-course income.
- d. What is the most frequent off-course income class for the 50 professional golfers? Using your tabular and graphical summaries, what additional observations can you make about the off-course income of these 50 professional golfers?
21. The *Nielsen Home Technology Report* provided information about home technology and its usage. The following data are the hours of personal computer usage during one week for a sample of 50 persons.

4.1	1.5	10.4	5.9	3.4	5.7	1.6	6.1	3.0	3.7
3.1	4.8	2.0	14.8	5.4	4.2	3.9	4.1	11.1	3.5
4.1	4.1	8.8	5.6	4.3	3.3	7.1	10.3	6.2	7.6
10.8	2.8	9.5	12.9	12.1	0.7	4.0	9.2	4.4	5.7
7.2	6.1	5.7	5.9	4.7	3.9	3.7	3.1	6.1	3.1

Summarize the data by constructing the following:

- a. A frequency distribution (use a class width of three hours)
- b. A relative frequency distribution
- c. A histogram
- d. An ogive
- e. Comment on what the data indicate about personal computer usage at home.

WEB file
OffCourse

WEB file
Computer

2.3

Exploratory Data Analysis: The Stem-and-Leaf Display

The techniques of **exploratory data analysis** consist of simple arithmetic and easy-to-draw graphs that can be used to summarize data quickly. One technique—referred to as a **stem-and-leaf display**—can be used to show both the rank order and shape of a data set simultaneously.

TABLE 2.8 NUMBER OF QUESTIONS ANSWERED CORRECTLY ON AN APTITUDE TEST

112	72	69	97	107
73	92	76	86	73
126	128	118	127	124
82	104	132	134	83
92	108	96	100	92
115	76	91	102	81
95	141	81	80	106
84	119	113	98	75
68	98	115	106	95
100	85	94	106	119

To illustrate the use of a stem-and-leaf display, consider the data in Table 2.8. These data result from a 150-question aptitude test given to 50 individuals recently interviewed for a position at Haskens Manufacturing. The data indicate the number of questions answered correctly.

To develop a stem-and-leaf display, we first arrange the leading digits of each data value to the left of a vertical line. To the right of the vertical line, we record the last digit for each data value. Based on the top row of data in Table 2.8 (112, 72, 69, 97, and 107), the first five entries in constructing a stem-and-leaf display would be as follows:

6	9
7	2
8	
9	7
10	7
11	2
12	
13	
14	

For example, the data value 112 shows the leading digits 11 to the left of the line and the last digit 2 to the right of the line. Similarly, the data value 72 shows the leading digit 7 to the left of the line and last digit 2 to the right of the line. Continuing to place the last digit of each data value on the line corresponding to its leading digit(s) provides the following:

6	9 8
7	2 3 6 3 6 5
8	6 2 3 1 1 0 4 5
9	7 2 2 6 2 1 5 8 8 5 4
10	7 4 8 0 2 6 6 0 6
11	2 8 5 9 3 5 9
12	6 8 7 4
13	2 4
14	1

With this organization of the data, sorting the digits on each line into rank order is simple. Doing so provides the stem-and-leaf display shown here.

6	8 9
7	2 3 3 5 6 6
8	0 1 1 2 3 4 5 6
9	1 2 2 2 4 5 5 6 7 8 8
10	0 0 2 4 6 6 6 7 8
11	2 3 5 5 8 9 9
12	4 6 7 8
13	2 4
14	1

The numbers to the left of the vertical line (6, 7, 8, 9, 10, 11, 12, 13, and 14) form the *stem*, and each digit to the right of the vertical line is a *leaf*. For example, consider the first row with a stem value of 6 and leaves of 8 and 9.

$$6 \mid 8 \ 9$$

This row indicates that two data values have a first digit of six. The leaves show that the data values are 68 and 69. Similarly, the second row

$$7 \mid 2 \ 3 \ 3 \ 5 \ 6 \ 6$$

indicates that six data values have a first digit of seven. The leaves show that the data values are 72, 73, 73, 75, 76, and 76.

To focus on the shape indicated by the stem-and-leaf display, let us use a rectangle to contain the leaves of each stem. Doing so, we obtain the following:

6	8 9
7	2 3 3 5 6 6
8	0 1 1 2 3 4 5 6
9	1 2 2 2 4 5 5 6 7 8 8
10	0 0 2 4 6 6 6 7 8
11	2 3 5 5 8 9 9
12	4 6 7 8
13	2 4
14	1

Rotating this page counterclockwise onto its side provides a picture of the data that is similar to a histogram with classes of 60–69, 70–79, 80–89, and so on.

Although the stem-and-leaf display may appear to offer the same information as a histogram, it has two primary advantages.

1. The stem-and-leaf display is easier to construct by hand.
2. Within a class interval, the stem-and-leaf display provides more information than the histogram because the stem-and-leaf shows the actual data.

Just as a frequency distribution or histogram has no absolute number of classes, neither does a stem-and-leaf display have an absolute number of rows or stems. If we believe that our original stem-and-leaf display condensed the data too much, we can easily stretch the display by using two or more stems for each leading digit. For example, to use two stems for each leading digit,

In a stretched stem-and-leaf display, whenever a stem value is stated twice, the first value corresponds to leaf values of 0–4, and the second value corresponds to leaf values of 5–9.

we would place all data values ending in 0, 1, 2, 3, and 4 in one row and all values ending in 5, 6, 7, 8, and 9 in a second row. The following stretched stem-and-leaf display illustrates this approach.

6	8 9
7	2 3 3
7	5 6 6
8	0 1 1 2 3 4
8	5 6
9	1 2 2 2 4
9	5 5 6 7 8 8
10	0 0 2 4
10	6 6 6 7 8
11	2 3
11	5 5 8 9 9
12	4
12	6 7 8
13	2 4
13	
14	1

Note that values 72, 73, and 73 have leaves in the 0–4 range and are shown with the first stem value of 7. The values 75, 76, and 76 have leaves in the 5–9 range and are shown with the second stem value of 7. This stretched stem-and-leaf display is similar to a frequency distribution with intervals of 65–69, 70–74, 75–79, and so on.

The preceding example showed a stem-and-leaf display for data with as many as three digits. Stem-and-leaf displays for data with more than three digits are possible. For example, consider the following data on the number of hamburgers sold by a fast-food restaurant for each of 15 weeks.

1565	1852	1644	1766	1888	1912	2044	1812
1790	1679	2008	1852	1967	1954	1733	

A stem-and-leaf display of these data follows.

Leaf unit = 10	
15	6
16	4 7
17	3 6 9
18	1 5 5 8
19	1 5 6
20	0 4

A single digit is used to define each leaf in a stem-and-leaf display. The leaf unit indicates how to multiply the stem-and-leaf numbers in order to approximate the original data. Leaf units may be 100, 10, 1, 0.1, and so on.

Note that a single digit is used to define each leaf and that only the first three digits of each data value have been used to construct the display. At the top of the display we have specified Leaf unit = 10. To illustrate how to interpret the values in the display, consider the first stem, 15, and its associated leaf, 6. Combining these numbers, we obtain 156. To reconstruct an approximation of the original data value, we must multiply this number by 10, the value of the leaf unit. Thus, $156 \times 10 = 1560$ is an approximation of the original data value used to construct the stem-and-leaf display. Although it is not possible to reconstruct the exact data value from this stem-and-leaf display, the convention of using a single digit for each leaf enables stem-and-leaf displays to be constructed for data having a large number of digits. For stem-and-leaf displays where the leaf unit is not shown, the leaf unit is assumed to equal 1.

Exercises

Methods

22. Construct a stem-and-leaf display for the following data.

70	72	75	64	58	83	80	82
76	75	68	65	57	78	85	72

23. Construct a stem-and-leaf display for the following data.

11.3	9.6	10.4	7.5	8.3	10.5	10.0
9.3	8.1	7.7	7.5	8.4	6.3	8.8

24. Construct a stem-and-leaf display for the following data. Use a leaf unit of 10.

1161	1206	1478	1300	1604	1725	1361	1422
1221	1378	1623	1426	1557	1730	1706	1689

SELF test

Applications

25. A psychologist developed a new test of adult intelligence. The test was administered to 20 individuals, and the following data were obtained.

114	99	131	124	117	102	106	127	119	115
98	104	144	151	132	106	125	122	118	118

Construct a stem-and-leaf display for the data.

26. The American Association of Individual Investors conducts an annual survey of discount brokers. The following prices charged are from a sample of 24 discount brokers (*AII Journal*, January 2003). The two types of trades are a broker-assisted trade of 100 shares at \$50 per share and an online trade of 500 shares at \$50 per share.

SELF test

Broker	Broker-Assisted 100 Shares at \$50/Share	Online 500 Shares at \$50/Share	Broker	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

WEB file
Broker

- a. Round the trading prices to the nearest dollar and develop a stem-and-leaf display for 100 shares at \$50 per share. Comment on what you learned about broker-assisted trading prices.
 - b. Round the trading prices to the nearest dollar and develop a stretched stem-and-leaf display for 500 shares online at \$50 per share. Comment on what you learned about online trading prices.
27. Most major ski resorts offer family programs that provide ski and snowboarding instruction for children. The typical classes provide four to six hours on the snow with a certified instructor. The daily rate for a group lesson at 15 ski resorts follows (*The Wall Street Journal*, January 20, 2006).

Resort	Location	Daily Rate	Resort	Location	Daily Rate
Beaver Creek	Colorado	\$137	Okemo	Vermont	\$ 86
Deer Valley	Utah	115	Park City	Utah	145
Diamond Peak	California	95	Butternut	Massachusetts	75
Heavenly	California	145	Steamboat	Colorado	98
Hunter	New York	79	Stowe	Vermont	104
Mammoth	California	111	Sugar Bowl	California	100
Mount Sunapee	New Hampshire	96	Whistler-Blackcomb	British Columbia	104
Mount Bachelor	Oregon	83			

- Develop a stem-and-leaf display for the data.
 - Interpret the stem-and-leaf display in terms of what it tells you about the daily rate for these ski and snowboarding instruction programs.
28. The 2004 Naples, Florida, minimarathon (13.1 miles) had 1228 registrants (*Naples Daily News*, January 17, 2004). Competition was held in six age groups. The following data show the ages for a sample of 40 individuals who participated in the marathon.



49	33	40	37	56
44	46	57	55	32
50	52	43	64	40
46	24	30	37	43
31	43	50	36	61
27	44	35	31	43
52	43	66	31	50
72	26	59	21	47

- Show a stretched stem-and-leaf display.
- What age group had the largest number of runners?
- What age occurred most frequently?
- A *Naples Daily News* feature article emphasized the number of runners who were “20-something.” What percentage of the runners were in the 20-something age group? What do you suppose was the focus of the article?

2.4

Crosstabulations and Scatter Diagrams

Crosstabulations and scatter diagrams are used to summarize data in a way that reveals the relationship between two variables.

Thus far in this chapter, we have focused on tabular and graphical methods used to summarize the data for *one variable at a time*. Often a manager or decision maker requires tabular and graphical methods that will assist in the understanding of the *relationship between two variables*. Crosstabulation and scatter diagrams are two such methods.

Crosstabulation

A **crosstabulation** is a tabular summary of data for two variables. Let us illustrate the use of a crosstabulation by considering the following application based on data from Zagat’s Restaurant Review. The quality rating and the meal price data were collected for a sample of 300 restaurants located in the Los Angeles area. Table 2.9 shows the data for the first 10 restaurants. Data on a restaurant’s quality rating and typical meal price are reported. Quality rating is a categorical variable with rating categories of good, very good, and excellent. Meal price is a quantitative variable that ranges from \$10 to \$49.

A crosstabulation of the data for this application is shown in Table 2.10. The left and top margin labels define the classes for the two variables. In the left margin, the row labels (good, very good, and excellent) correspond to the three classes of the quality rating variable. In the top margin, the column labels (\$10–19, \$20–29, \$30–39, and \$40–49) correspond to

TABLE 2.9 QUALITY RATING AND MEAL PRICE FOR 300 LOS ANGELES RESTAURANTS

WEB file
Restaurant

Restaurant	Quality Rating	Meal Price (\$)
1	Good	18
2	Very Good	22
3	Good	28
4	Excellent	38
5	Very Good	33
6	Good	28
7	Very Good	19
8	Very Good	11
9	Very Good	23
10	Good	13
.	.	.
.	.	.
.	.	.

the four classes of the meal price variable. Each restaurant in the sample provides a quality rating and a meal price. Thus, each restaurant in the sample is associated with a cell appearing in one of the rows and one of the columns of the crosstabulation. For example, restaurant 5 is identified as having a very good quality rating and a meal price of \$33. This restaurant belongs to the cell in row 2 and column 3 of Table 2.10. In constructing a crosstabulation, we simply count the number of restaurants that belong to each of the cells in the crosstabulation table.

In reviewing Table 2.10, we see that the greatest number of restaurants in the sample (64) have a very good rating and a meal price in the \$20–29 range. Only two restaurants have an excellent rating and a meal price in the \$10–19 range. Similar interpretations of the other frequencies can be made. In addition, note that the right and bottom margins of the crosstabulation provide the frequency distributions for quality rating and meal price separately. From the frequency distribution in the right margin, we see that data on quality ratings show 84 good restaurants, 150 very good restaurants, and 66 excellent restaurants. Similarly, the bottom margin shows the frequency distribution for the meal price variable.

Dividing the totals in the right margin of the crosstabulation by the total for that column provides a relative and percent frequency distribution for the quality rating variable.

Quality Rating	Relative Frequency	Percent Frequency
Good	.28	28
Very Good	.50	50
Excellent	.22	22
Total	1.00	100

TABLE 2.10 CROSSTABULATION OF QUALITY RATING AND MEAL PRICE FOR 300 LOS ANGELES RESTAURANTS

Quality Rating	Meal Price				Total
	\$10–19	\$20–29	\$30–39	\$40–49	
Good	42	40	2	0	84
Very Good	34	64	46	6	150
Excellent	2	14	28	22	66
Total	78	118	76	28	300

From the percent frequency distribution we see that 28% of the restaurants were rated good, 50% were rated very good, and 22% were rated excellent.

Dividing the totals in the bottom row of the crosstabulation by the total for that row provides a relative and percent frequency distribution for the meal price variable.

Meal Price	Relative Frequency	Percent Frequency
\$10–19	.26	26
\$20–29	.39	39
\$30–39	.25	25
\$40–49	.09	9
Total	1.00	100

Note that the sum of the values in each column does not add exactly to the column total, because the values being summed are rounded. From the percent frequency distribution we see that 26% of the meal prices are in the lowest price class (\$10–19), 39% are in the next higher class, and so on.

The frequency and relative frequency distributions constructed from the margins of a crosstabulation provide information about each of the variables individually, but they do not shed any light on the relationship between the variables. The primary value of a crosstabulation lies in the insight it offers about the relationship between the variables. A review of the crosstabulation in Table 2.10 reveals that higher meal prices are associated with the higher quality restaurants, and the lower meal prices are associated with the lower quality restaurants.

Converting the entries in a crosstabulation into row percentages or column percentages can provide more insight into the relationship between the two variables. For row percentages, the results of dividing each frequency in Table 2.10 by its corresponding row total are shown in Table 2.11. Each row of Table 2.11 is a percent frequency distribution of meal price for one of the quality rating categories. Of the restaurants with the lowest quality rating (good), we see that the greatest percentages are for the less expensive restaurants (50% have \$10–19 meal prices and 47.6% have \$20–29 meal prices). Of the restaurants with the highest quality rating (excellent), we see that the greatest percentages are for the more expensive restaurants (42.4% have \$30–39 meal prices and 33.4% have \$40–49 meal prices). Thus, we continue to see that the more expensive meals are associated with the higher quality restaurants.

Crosstabulation is widely used for examining the relationship between two variables. In practice, the final reports for many statistical studies include a large number of crosstabulation tables. In the Los Angeles restaurant survey, the crosstabulation is based on one qualitative variable (quality rating) and one quantitative variable (meal price). Crosstabulations can also be developed when both variables are qualitative and when both variables are quantitative. When quantitative variables are used, however, we must first create classes for the values of the variable. For instance, in the restaurant example we grouped the meal prices into four classes (\$10–19, \$20–29, \$30–39, and \$40–49).

TABLE 2.11 ROW PERCENTAGES FOR EACH QUALITY RATING CATEGORY

Quality Rating	Meal Price				Total
	\$10–19	\$20–29	\$30–39	\$40–49	
Good	50.0	47.6	2.4	0.0	100
Very Good	22.7	42.7	30.6	4.0	100
Excellent	3.0	21.2	42.4	33.4	100

Simpson's Paradox

The data in two or more crosstabulations are often combined or aggregated to produce a summary crosstabulation showing how two variables are related. In such cases, we must be careful in drawing a conclusion because a conclusion based upon aggregate data can be reversed if we look at the unaggregated data. The reversal of conclusions based on aggregate and unaggregated data is called **Simpson's paradox**. To provide an illustration of Simpson's paradox we consider an example involving the analysis of verdicts for two judges in two different courts.

Judges Ron Luckett and Dennis Kendall presided over cases in Common Pleas Court and Municipal Court during the past three years. Some of the verdicts they rendered were appealed. In most of these cases the appeals court upheld the original verdicts, but in some cases those verdicts were reversed. For each judge a crosstabulation was developed based upon two variables: Verdict (upheld or reversed) and Type of Court (Common Pleas and Municipal). Suppose that the two crosstabulations were then combined by aggregating the type of court data. The resulting aggregated crosstabulation contains two variables: Verdict (upheld or reversed) and Judge (Luckett or Kendall). This crosstabulation shows the number of appeals in which the verdict was upheld and the number in which the verdict was reversed for both judges. The following crosstabulation shows these results along with the column percentages in parentheses next to each value.

Verdict	Judge		Total
	Luckett	Kendall	
Upheld	129 (86%)	110 (88%)	239
Reversed	21 (14%)	15 (12%)	36
Total (%)	150 (100%)	125 (100%)	275

A review of the column percentages shows that 86% of the verdicts were upheld for Judge Luckett, while 88% of the verdicts were upheld for Judge Kendall. From this aggregated crosstabulation, we conclude that Judge Kendall is doing the better job because a greater percentage of Judge Kendall's verdicts are being upheld.

The following unaggregated crosstabulations show the cases tried by Judge Luckett and Judge Kendall in each court; column percentages are shown in parentheses next to each value.

Verdict	Judge Luckett			Verdict	Judge Kendall		
	Common Pleas	Municipal Court	Total		Common Pleas	Municipal Court	Total
Upheld	29 (91%)	100 (85%)	129	Upheld	90 (90%)	20 (80%)	110
Reversed	3 (9%)	18 (15%)	21	Reversed	10 (10%)	5 (20%)	15
Total (%)	32 (100%)	118 (100%)	150	Total (%)	100 (100%)	25 (100%)	125

From the crosstabulation and column percentages for Judge Luckett, we see that the verdicts were upheld in 91% of the Common Pleas Court cases and in 85% of the Municipal Court cases. From the crosstabulation and column percentages for Judge Kendall, we see that the verdicts were upheld in 90% of the Common Pleas Court cases and in 80% of the Municipal Court cases. Thus, when we unaggregate the data, we see that Judge Luckett has a better record because a greater percentage of Judge Luckett's verdicts are being upheld in both courts. This result contradicts the conclusion we reached with the aggregated data crosstabulation that showed Judge Kendall had the better record. This reversal of conclusions based on aggregated and unaggregated data illustrates Simpson's paradox.

The original crosstabulation was obtained by aggregating the data in the separate crosstabulations for the two courts. Note that for both judges the percentage of appeals that resulted in reversals was much higher in Municipal Court than in Common Pleas Court. Because Judge Lockett tried a much higher percentage of his cases in Municipal Court, the aggregated data favored Judge Kendall. When we look at the crosstabulations for the two courts separately, however, Judge Lockett shows the better record. Thus, for the original crosstabulation, we see that the *type of court* is a hidden variable that cannot be ignored when evaluating the records of the two judges.

Because of the possibility of Simpson's paradox, realize that the conclusion or interpretation may be reversed depending upon whether you are viewing unaggregated or aggregate crosstabulation data. Before drawing a conclusion, you may want to investigate whether the aggregate or unaggregate form of the crosstabulation provides the better insight and conclusion. Especially when the crosstabulation involves aggregated data, you should investigate whether a hidden variable could affect the results such that separate or unaggregated crosstabulations provide a different and possibly better insight and conclusion.

Scatter Diagram and Trendline

A **scatter diagram** is a graphical presentation of the relationship between two quantitative variables, and a **trendline** is a line that provides an approximation of the relationship. As an illustration, consider the advertising/sales relationship for a stereo and sound equipment store in San Francisco. On 10 occasions during the past three months, the store used weekend television commercials to promote sales at its stores. The managers want to investigate whether a relationship exists between the number of commercials shown and sales at the store during the following week. Sample data for the 10 weeks with sales in hundreds of dollars are shown in Table 2.12.

Figure 2.7 shows the scatter diagram and the trendline¹ for the data in Table 2.12. The number of commercials (x) is shown on the horizontal axis and the sales (y) are shown on the vertical axis. For week 1, $x = 2$ and $y = 50$. A point with those coordinates is plotted on the scatter diagram. Similar points are plotted for the other nine weeks. Note that during two of the weeks one commercial was shown, during two of the weeks two commercials were shown, and so on.

The completed scatter diagram in Figure 2.7 indicates a positive relationship between the number of commercials and sales. Higher sales are associated with a higher number of commercials. The relationship is not perfect in that all points are not on a straight line. However, the general pattern of the points and the trendline suggest that the overall relationship is positive.

TABLE 2.12 SAMPLE DATA FOR THE STEREO AND SOUND EQUIPMENT STORE

Week	Number of Commercials x	Sales (\$100s) y
1	2	50
2	5	57
3	1	41
4	3	54
5	4	54
6	1	38
7	5	63
8	3	48
9	4	59
10	2	46



¹The equation of the trendline is $y = 36.15 + 4.95x$. The slope of the trendline is 4.95 and the y -intercept (the point where the line intersects the y -axis) is 36.15. We will discuss in detail the interpretation of the slope and y -intercept for a linear trendline in Chapter 14 when we study simple linear regression.

FIGURE 2.7 SCATTER DIAGRAM AND TRENDLINE FOR THE STEREO AND SOUND EQUIPMENT STORE

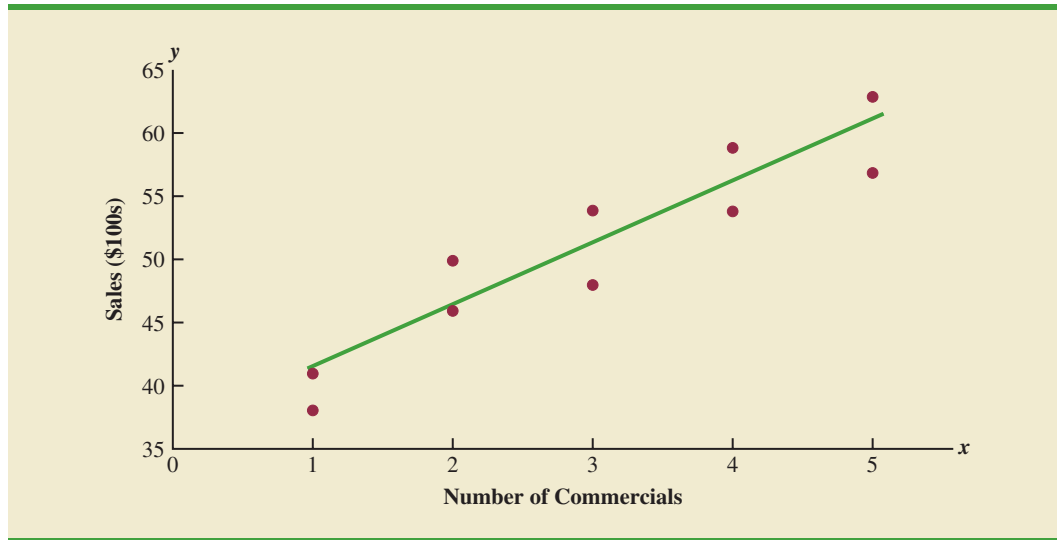
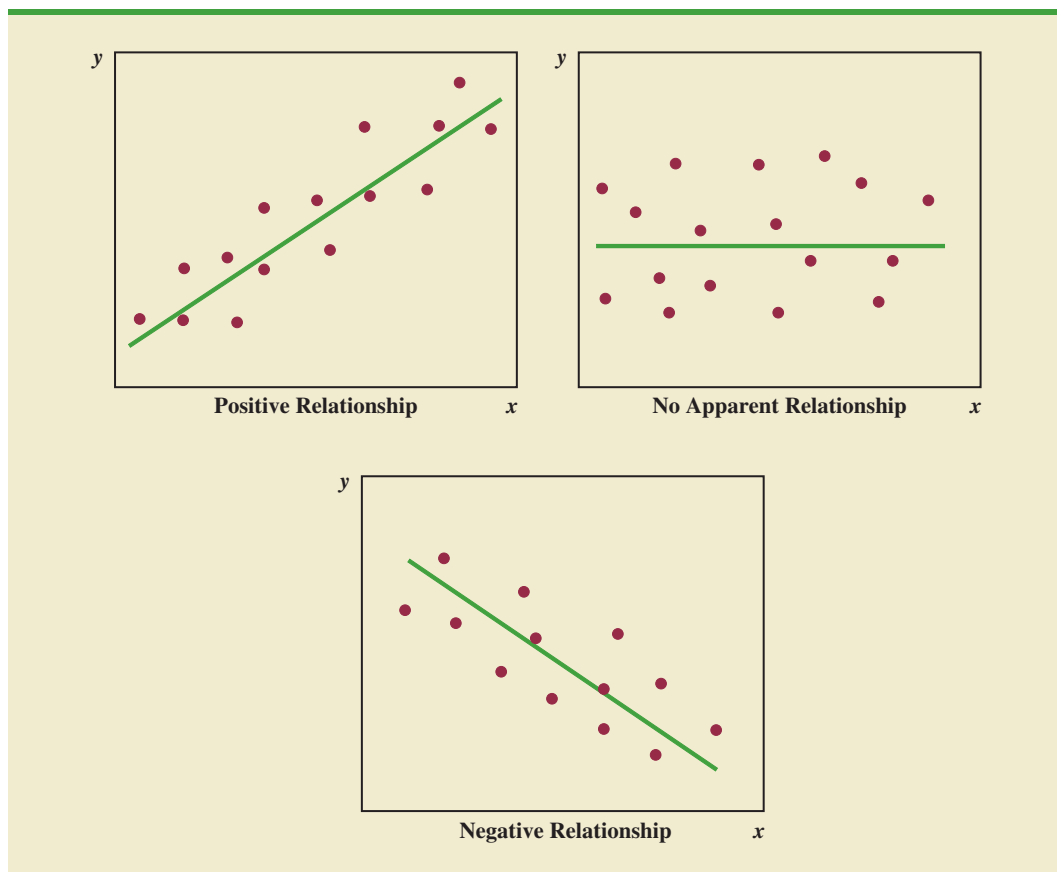


FIGURE 2.8 TYPES OF RELATIONSHIPS DEPICTED BY SCATTER DIAGRAMS



Some general scatter diagram patterns and the types of relationships they suggest are shown in Figure 2.8. The top left panel depicts a positive relationship similar to the one for the number of commercials and sales example. In the top right panel, the scatter diagram shows no apparent relationship between the variables. The bottom panel depicts a negative relationship where y tends to decrease as x increases.

Exercises

Methods

SELF test

29. The following data are for 30 observations involving two qualitative variables, x and y . The categories for x are A, B, and C; the categories for y are 1 and 2.

WEB file

Crosstab

Observation	x	y	Observation	x	y
1	A	1	16	B	2
2	B	1	17	C	1
3	B	1	18	B	1
4	C	2	19	C	1
5	B	1	20	B	1
6	C	2	21	C	2
7	B	1	22	B	1
8	C	2	23	C	2
9	A	1	24	A	1
10	B	1	25	B	1
11	A	1	26	C	2
12	B	1	27	C	2
13	C	2	28	A	1
14	C	2	29	B	1
15	C	2	30	B	2

- Develop a crosstabulation for the data, with x as the row variable and y as the column variable.
 - Compute the row percentages.
 - Compute the column percentages.
 - What is the relationship, if any, between x and y ?
30. The following 20 observations are for two quantitative variables, x and y .

SELF test

WEB file

Scatter

Observation	x	y	Observation	x	y
1	-22	22	11	-37	48
2	-33	49	12	34	-29
3	2	8	13	9	-18
4	29	-16	14	-33	31
5	-13	10	15	20	-16
6	21	-28	16	-3	14
7	-13	27	17	-15	18
8	-23	35	18	12	17
9	14	-5	19	-20	-11
10	3	-3	20	-7	-22

- Develop a scatter diagram for the relationship between x and y .
- What is the relationship, if any, between x and y ?

Applications

31. The following crosstabulation shows household income by educational level of the head of household (*Statistical Abstract of the United States: 2008*).

Educational Level	Household Income (\$1000s)					Total
	Under 25	25.0–49.9	50.0–74.9	75.0–99.9	100 or more	
Not H.S. graduate	4207	3459	1389	539	367	9961
H.S. graduate	4917	6850	5027	2637	2668	22099
Some college	2807	5258	4678	3250	4074	20067
Bachelor's degree	885	2094	2848	2581	5379	13787
Beyond bach. deg.	290	829	1274	1241	4188	7822
Total	13106	18490	15216	10248	16676	73736

- Compute the row percentages and identify the percent frequency distributions of income for households in which the head is a high school graduate and in which the head holds a bachelor's degree.
 - What percentage of households headed by high school graduates earn \$75,000 or more? What percentage of households headed by bachelor's degree recipients earn \$75,000 or more?
 - Construct percent frequency histograms of income for households headed by persons with a high school degree and for those headed by persons with a bachelor's degree. Is any relationship evident between household income and educational level?
32. Refer again to the crosstabulation of household income by educational level shown in exercise 31.
- Compute column percentages and identify the percent frequency distributions displayed. What percentage of the heads of households did not graduate from high school?
 - What percentage of the households earning \$100,000 or more were headed by a person having schooling beyond a bachelor's degree? What percentage of the households headed by a person with schooling beyond a bachelor's degree earned over \$100,000? Why are these two percentages different?
 - Compare the percent frequency distributions for those households earning "Under 25," "100 or more," and for "Total." Comment on the relationship between household income and educational level of the head of household.
33. Recently, management at Oak Tree Golf Course received a few complaints about the condition of the greens. Several players complained that the greens are too fast. Rather than react to the comments of just a few, the Golf Association conducted a survey of 100 male and 100 female golfers. The survey results are summarized here.

Male Golfers			Female Golfers		
Handicap	Greens Condition		Handicap	Greens Condition	
	Too Fast	Fine		Too Fast	Fine
Under 15	10	40	Under 15	1	9
15 or more	25	25	15 or more	39	51

- Combine these two crosstabulations into one with Male and Female as the row labels and Too Fast and Fine as the column labels. Which group shows the highest percentage saying that the greens are too fast?

- b. Refer to the initial crosstabulations. For those players with low handicaps (better players), which group (male or female) shows the highest percentage saying the greens are too fast?
 - c. Refer to the initial crosstabulations. For those players with higher handicaps, which group (male or female) shows the highest percentage saying the greens are too fast?
 - d. What conclusions can you draw about the preferences of men and women concerning the speed of the greens? Are the conclusions you draw from part (a) as compared with parts (b) and (c) consistent? Explain any apparent inconsistencies.
34. Table 2.13 shows a data set containing information for 45 mutual funds that are part of the *Morningstar Funds500* for 2008. The data set includes the following five variables:
- Fund Type: The type of fund, labeled DE (Domestic Equity), IE (International Equity), and FI (Fixed Income)
 - Net Asset Value (\$): The closing price per share
 - 5-Year Average Return (%): The average annual return for the fund over the past 5 years
 - Expense Ratio (%): The percentage of assets deducted each fiscal year for fund expenses
 - Morningstar Rank: The risk adjusted star rating for each fund; Morningstar ranks go from a low of 1-Star to a high of 5-Stars
- a. Prepare a crosstabulation of the data on Fund Type (rows) and the average annual return over the past 5 years (columns). Use classes of 0–9.99, 10–19.99, 20–29.99, 30–39.99, 40–49.99, and 50–59.99 for the 5-Year Average Return (%).
 - b. Prepare a frequency distribution for the data on Fund Type.
 - c. Prepare a frequency distribution for the data on 5-Year Average Return (%).
 - d. How has the crosstabulation helped in preparing the frequency distributions in parts (b) and (c)?
 - e. What conclusions can you draw about the fund type and the average return over the past 5 years?
35. Refer to the data in Table 2.13.
- a. Prepare a crosstabulation of the data on Fund Type (rows) and the expense ratio (columns). Use classes of .25–.49, .50–.74, .75–.99, 1.00–1.24, and 1.25–1.49 for Expense Ratio (%).
 - b. Prepare a percent frequency distribution for Expense Ratio (%).
 - c. What conclusions can you draw about fund type and the expense ratio?
36. Refer to the data in Table 2.13.
- a. Prepare a scatter diagram with 5-Year Average Return (%) on the horizontal axis and Net Asset Value (\$) on the vertical axis.
 - b. Comment on the relationship, if any, between the variables.
37. The U.S. Department of Energy’s Fuel Economy Guide provides fuel efficiency data for cars and trucks (Fuel Economy website, February 22, 2008). A portion of the data for 311 compact, midsize, and large cars is shown in Table 2.14. The data set contains the following variables:
- Size: Compact, Midsize, and Large
 - Displacement: Engine size in liters
 - Cylinders: Number of cylinders in the engine
 - Drive: Front wheel (F), rear wheel (R), and four wheel (4)
 - Fuel Type: Premium (P) or regular (R) fuel
 - City MPG: Fuel efficiency rating for city driving in terms of miles per gallon
 - Hwy MPG: Fuel efficiency rating for highway driving in terms of miles per gallon

The complete data set is contained in the file named FuelData08.

- Prepare a crosstabulation of the data on Size (rows) and Hwy MPG (columns). Use classes of 15–19, 20–24, 25–29, 30–34, and 35–39 for Hwy MPG.
- Comment on the relationship between Size and Hwy MPG.

TABLE 2.13 FINANCIAL DATA FOR A SAMPLE OF 45 MUTUAL FUNDS

Fund Name	Fund Type	Net Asset Value (\$)	5-Year Average Return (%)	Expense Ratio (%)	Morningstar Rank
Amer Cent Inc & Growth Inv	DE	28.88	12.39	0.67	2-Star
American Century Intl. Disc	IE	14.37	30.53	1.41	3-Star
American Century Tax-Free Bond	FI	10.73	3.34	0.49	4-Star
American Century Ultra	DE	24.94	10.88	0.99	3-Star
Ariel	DE	46.39	11.32	1.03	2-Star
Artisan Intl Val	IE	25.52	24.95	1.23	3-Star
Artisan Small Cap	DE	16.92	15.67	1.18	3-Star
Baron Asset	DE	50.67	16.77	1.31	5-Star
Brandywine	DE	36.58	18.14	1.08	4-Star
Brown Cap Small	DE	35.73	15.85	1.20	4-Star
Buffalo Mid Cap	DE	15.29	17.25	1.02	3-Star
Delafield	DE	24.32	17.77	1.32	4-Star
DFA U.S. Micro Cap	DE	13.47	17.23	0.53	3-Star
Dodge & Cox Income	FI	12.51	4.31	0.44	4-Star
Fairholme	DE	31.86	18.23	1.00	5-Star
Fidelity Contrafund	DE	73.11	17.99	0.89	5-Star
Fidelity Municipal Income	FI	12.58	4.41	0.45	5-Star
Fidelity Overseas	IE	48.39	23.46	0.90	4-Star
Fidelity Sel Electronics	DE	45.60	13.50	0.89	3-Star
Fidelity Sh-Term Bond	FI	8.60	2.76	0.45	3-Star
Fidelity	DE	39.85	14.40	0.56	4-Star
FPA New Income	FI	10.95	4.63	0.62	3-Star
Gabelli Asset AAA	DE	49.81	16.70	1.36	4-Star
Greenspring	DE	23.59	12.46	1.07	3-Star
Janus	DE	32.26	12.81	0.90	3-Star
Janus Worldwide	IE	54.83	12.31	0.86	2-Star
Kalmar Gr Val Sm Cp	DE	15.30	15.31	1.32	3-Star
Managers Freemont Bond	FI	10.56	5.14	0.60	5-Star
Marsico 21st Century	DE	17.44	15.16	1.31	5-Star
Mathews Pacific Tiger	IE	27.86	32.70	1.16	3-Star
Meridan Value	DE	31.92	15.33	1.08	4-Star
Oakmark I	DE	40.37	9.51	1.05	2-Star
PIMCO Emerg Mkts Bd D	FI	10.68	13.57	1.25	3-Star
RS Value A	DE	26.27	23.68	1.36	4-Star
T. Rowe Price Latin Am.	IE	53.89	51.10	1.24	4-Star
T. Rowe Price Mid Val	DE	22.46	16.91	0.80	4-Star
Templeton Growth A	IE	24.07	15.91	1.01	3-Star
Thornburg Value A	DE	37.53	15.46	1.27	4-Star
USAA Income	FI	12.10	4.31	0.62	3-Star
Vanguard Equity-Inc	DE	24.42	13.41	0.29	4-Star
Vanguard Global Equity	IE	23.71	21.77	0.64	5-Star
Vanguard GNMA	FI	10.37	4.25	0.21	5-Star
Vanguard Sht-Tm TE	FI	15.68	2.37	0.16	3-Star
Vanguard Sm Cp Idx	DE	32.58	17.01	0.23	3-Star
Wasatch Sm Cp Growth	DE	35.41	13.98	1.19	4-Star

TABLE 2.14 FUEL EFFICIENCY DATA FOR 311 CARS

Car	Size	Displacement	Cylinders	Drive	Fuel Type	City MPG	Hwy MPG
1	Compact	3.1	6	4	P	15	25
2	Compact	3.1	6	4	P	17	25
3	Compact	3.0	6	4	P	17	25
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
161	Midsize	2.4	4	F	R	22	30
162	Midsize	2.0	4	F	P	19	29
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
•	•	•	•	•	•	•	•
310	Large	3.0	6	F	R	17	25
311	Large	3.0	6	F	R	18	25

WEB file

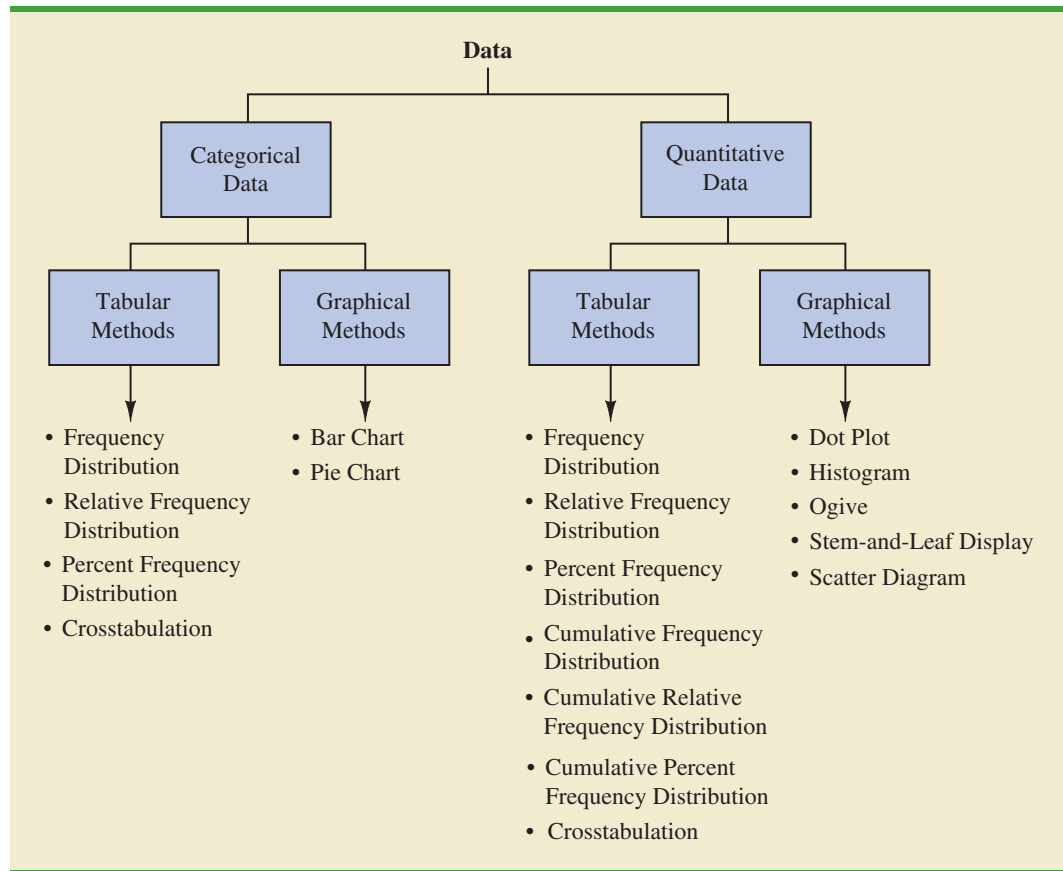
FuelData08

- c. Prepare a crosstabulation of the data on Drive (rows) and City MPG (columns). Use classes of 5–9, 10–14, 15–19, 20–24, 25–29, 30–34, and 35–39 for City MPG.
 - d. Comment on the relationship between Drive and City MPG.
 - e. Prepare a crosstabulation of the data on Fuel Type (rows) and City MPG (columns). Use classes of 5–9, 10–14, 15–19, 20–24, 25–29, 30–34, and 35–39 for City MPG.
 - f. Comment on the relationship between Fuel Type and City MPG.
38. Refer to exercise 37 and the data in the file named FuelData08.
- a. Prepare a crosstabulation of the data on Displacement (rows) and Hwy MPG (columns). Use classes of 1.0–2.9, 3.0–4.9, and 5.0–6.9 for Displacement. Use classes of 15–19, 20–24, 25–29, 30–34, and 35–39 for Hwy MPG.
 - b. Comment on the relationship, if any, between Displacement and Hwy MPG.
 - c. Develop a scatter diagram of the data on Displacement and Hwy MPG. Use the vertical axis for Hwy MPG.
 - d. What does the scatter diagram developed in part (c) indicate about the relationship, if any, between Displacement and Hwy MPG?
 - e. In investigating the relationship between Displacement and Hwy MPG you developed a tabular summary of the data (crosstabulation) and a graphical summary (scatter diagram). In this case which approach do you prefer? Explain.

Summary

A set of data, even if modest in size, is often difficult to interpret directly in the form in which it is gathered. Tabular and graphical methods provide procedures for organizing and summarizing data so that patterns are revealed and the data are more easily interpreted. Frequency distributions, relative frequency distributions, percent frequency distributions, bar charts, and pie charts were presented as tabular and graphical procedures for summarizing qualitative data. Frequency distributions, relative frequency distributions, percent frequency distributions, histograms, cumulative frequency distributions, cumulative relative frequency distributions, cumulative percent frequency distributions, and ogives were presented as ways of summarizing quantitative data. A stem-and-leaf display provides an exploratory data analysis technique that can be used to summarize quantitative data. Crosstabulation was presented as a tabular method for summarizing data for two variables. The scatter diagram was introduced as a graphical method for showing the relationship between two quantitative variables. Figure 2.9 shows the tabular and graphical methods presented in this chapter.

FIGURE 2.9 TABULAR AND GRAPHICAL METHODS FOR SUMMARIZING DATA



With large data sets, computer software packages are essential in constructing tabular and graphical summaries of data. In the chapter appendixes, we show how Minitab, Excel, and StatTools can be used for this purpose.

Glossary

Categorical data Labels or names used to identify categories of like items.

Quantitative data Numerical values that indicate how much or how many.

Frequency distribution A tabular summary of data showing the number (frequency) of data values in each of several nonoverlapping classes.

Relative frequency distribution A tabular summary of data showing the fraction or proportion of data values in each of several nonoverlapping classes.

Percent frequency distribution A tabular summary of data showing the percentage of data values in each of several nonoverlapping classes.

Bar chart A graphical device for depicting qualitative data that have been summarized in a frequency, relative frequency, or percent frequency distribution.

Pie chart A graphical device for presenting data summaries based on subdivision of a circle into sectors that correspond to the relative frequency for each class.

Class midpoint The value halfway between the lower and upper class limits.

Dot plot A graphical device that summarizes data by the number of dots above each data value on the horizontal axis.

Histogram A graphical presentation of a frequency distribution, relative frequency distribution, or percent frequency distribution of quantitative data constructed by placing the class intervals on the horizontal axis and the frequencies, relative frequencies, or percent frequencies on the vertical axis.

Cumulative frequency distribution A tabular summary of quantitative data showing the number of data values that are less than or equal to the upper class limit of each class.

Cumulative relative frequency distribution A tabular summary of quantitative data showing the fraction or proportion of data values that are less than or equal to the upper class limit of each class.

Cumulative percent frequency distribution A tabular summary of quantitative data showing the percentage of data values that are less than or equal to the upper class limit of each class.

Ogive A graph of a cumulative distribution.

Exploratory data analysis Methods that use simple arithmetic and easy-to-draw graphs to summarize data quickly.

Stem-and-leaf display An exploratory data analysis technique that simultaneously rank orders quantitative data and provides insight about the shape of the distribution.

Crosstabulation A tabular summary of data for two variables. The classes for one variable are represented by the rows; the classes for the other variable are represented by the columns.

Simpson's paradox Conclusions drawn from two or more separate crosstabulations that can be reversed when the data are aggregated into a single crosstabulation.

Scatter diagram A graphical presentation of the relationship between two quantitative variables. One variable is shown on the horizontal axis and the other variable is shown on the vertical axis.

Trendline A line that provides an approximation of the relationship between two variables.

Key Formulas

Relative Frequency

$$\frac{\text{Frequency of the class}}{n} \quad (2.1)$$

Approximate Class Width

$$\frac{\text{Largest data value} - \text{Smallest data value}}{\text{Number of classes}} \quad (2.2)$$

Supplementary Exercises

39. The Higher Education Research Institute at UCLA provides statistics on the most popular majors among incoming college freshmen. The five most popular majors are Arts and Humanities (A), Business Administration (B), Engineering (E), Professional (P), and Social Science (S) (*The New York Times Almanac*, 2006). A broad range of other (O) majors, including biological science, physical science, computer science, and education, are grouped together. The majors selected for a sample of 64 college freshmen follow.

S	P	P	O	B	E	O	E	P	O	O	B	O	O	O	A
O	E	E	B	S	O	B	O	A	O	E	O	E	O	B	P
B	A	S	O	E	A	B	O	S	S	O	O	E	B	O	B
A	E	B	E	A	A	P	O	O	E	O	B	B	O	P	B

- Show a frequency distribution and percent frequency distribution.
- Show a bar chart.



Major



- c. What percentage of freshmen select one of the five most popular majors?
 - d. What is the most popular major for incoming freshmen? What percentage of freshmen select this major?
40. General Motors had a 23% share of the automobile industry with sales coming from eight divisions: Buick, Cadillac, Chevrolet, GMC, Hummer, Pontiac, Saab, and Saturn (*Forbes*, December 22, 2008). The data set GMSales shows the sales for a sample of 200 General Motors vehicles. The division for the vehicle is provided for each sale.
- a. Show the frequency distribution and the percent frequency distribution of sales by division for General Motors.
 - b. Show a bar chart of the percent frequency distribution.
 - c. Which General Motors division was the company leader in sales? What was the percentage of sales for this division? Was this General Motors' most important division? Explain.
 - d. Due to the ongoing recession, high gasoline prices, and the decline in automobile sales, General Motors was facing bankruptcy in 2009. A government "bail-out" loan and a restructuring of the company were anticipated. Expectations were that General Motors could not continue to operate all eight divisions. Based on the percentage of sales, which of the eight divisions looked to be the best candidates for General Motors to discontinue? Which divisions looked to be the least likely candidates for General Motors to discontinue?
41. Dividend yield is the annual dividend paid by a company expressed as a percentage of the price of the stock ($\text{Dividend}/\text{Stock Price} \times 100$). The dividend yield for the Dow Jones Industrial Average companies is shown in Table 2.15 (*The Wall Street Journal*, June 8, 2009).
- a. Construct a frequency distribution and percent frequency distribution.
 - b. Construct a histogram.
 - c. Comment on the shape of the distribution.
 - d. What do the tabular and graphical summaries tell about the dividend yields among the Dow Jones Industrial Average companies?
 - e. What company has the highest dividend yield? If the stock for this company currently sells for \$20 per share and you purchase 500 shares, how much dividend income will this investment generate in one year?
42. Approximately 1.5 million high school students take the Scholastic Aptitude Test (SAT) each year and nearly 80% of the college and universities without open admissions policies use SAT scores in making admission decisions (College Board, March 2009). The current

TABLE 2.15 DIVIDEND YIELD FOR DOW JONES INDUSTRIAL AVERAGE COMPANIES

Company	Dividend Yield %	Company	Dividend Yield %
3M	3.6	IBM	2.1
Alcoa	1.3	Intel	3.4
American Express	2.9	J.P. Morgan Chase	0.5
AT&T	6.6	Johnson & Johnson	3.6
Bank of America	0.4	Kraft Foods	4.4
Boeing	3.8	McDonald's	3.4
Caterpillar	4.7	Merck	5.5
Chevron	3.9	Microsoft	2.5
Cisco Systems	0.0	Pfizer	4.2
Coca-Cola	3.3	Procter & Gamble	3.4
DuPont	5.8	Travelers	3.0
ExxonMobil	2.4	United Technologies	2.9
General Electric	9.2	Verizon	6.3
Hewlett-Packard	0.9	Wal-Mart Stores	2.2
Home Depot	3.9	Walt Disney	1.5



version of the SAT includes three parts: reading comprehension, mathematics, and writing. A perfect combined score for all three parts is 2400. A sample of SAT scores for the combined three-part SAT are as follows:



1665	1525	1355	1645	1780
1275	2135	1280	1060	1585
1650	1560	1150	1485	1990
1590	1880	1420	1755	1375
1475	1680	1440	1260	1730
1490	1560	940	1390	1175

- Show a frequency distribution and histogram. Begin with the first class starting at 800 and use a class width of 200.
 - Comment on the shape of the distribution.
 - What other observations can be made about the SAT scores based on the tabular and graphical summaries?
43. The Pittsburgh Steelers defeated the Arizona Cardinals 27 to 23 in professional football's 43rd Super Bowl. With this win, its sixth championship, the Pittsburgh Steelers became the team with the most wins in the 43-year history of the event (*Tampa Tribune*, February 2, 2009). The Super Bowl has been played in eight different states: Arizona (AZ), California (CA), Florida (FL), Georgia (GA), Louisiana (LA), Michigan (MI), Minnesota (MN), and Texas (TX). Data in the following table show the state where the Super Bowls were played and the point margin of victory for the winning team.



Super Bowl	State	Won By Points	Super Bowl	State	Won By Points	Super Bowl	State	Won By Points
1	CA	25	16	MI	5	31	LA	14
2	FL	19	17	CA	10	32	CA	7
3	FL	9	18	FL	19	33	FL	15
4	LA	16	19	CA	22	34	GA	7
5	FL	3	20	LA	36	35	FL	27
6	FL	21	21	CA	19	36	LA	3
7	CA	7	22	CA	32	37	CA	27
8	TX	17	23	FL	4	38	TX	3
9	LA	10	24	LA	45	39	FL	3
10	FL	4	25	FL	1	40	MI	11
11	CA	18	26	MN	13	41	FL	12
12	LA	17	27	CA	35	42	AZ	3
13	FL	4	28	GA	17	43	FL	4
14	CA	12	29	FL	23			
15	LA	17	30	AZ	10			

- Show a frequency distribution and bar chart for the state where the Super Bowl was played.
- What conclusions can you draw from your summary in part (a)? What percentage of Super Bowls were played in the states of Florida or California? What percentage of Super Bowls were played in northern or cold-weather states?
- Show a stretched stem-and-leaf display for the point margin of victory for the winning team. Show a histogram.
- What conclusions can you draw from your summary in part (c)? What percentage of Super Bowls have been close games with the margin of victory less than 5 points? What percentage of Super Bowls have been won by 20 or more points?
- The closest Super Bowl occurred when the New York Giants beat the Buffalo Bills. Where was this game played and what was the winning margin of victory? The biggest point margin in Super Bowl history occurred when the San Francisco 49ers beat the Denver Broncos. Where was this game played and what was the winning margin of victory?

44. Data from the U.S. Census Bureau provides the population by state in millions of people (*The World Almanac*, 2006).

WEB file
Population

State	Population	State	Population	State	Population
Alabama	4.5	Louisiana	4.5	Ohio	11.5
Alaska	0.7	Maine	1.3	Oklahoma	3.5
Arizona	5.7	Maryland	5.6	Oregon	3.6
Arkansas	2.8	Massachusetts	6.4	Pennsylvania	12.4
California	35.9	Michigan	10.1	Rhode Island	1.1
Colorado	4.6	Minnesota	5.1	South Carolina	4.2
Connecticut	3.5	Mississippi	2.9	South Dakota	0.8
Delaware	0.8	Missouri	5.8	Tennessee	5.9
Florida	17.4	Montana	0.9	Texas	22.5
Georgia	8.8	Nebraska	1.7	Utah	2.4
Hawaii	1.3	Nevada	2.3	Vermont	0.6
Idaho	1.4	New Hampshire	1.3	Virginia	7.5
Illinois	12.7	New Jersey	8.7	Washington	6.2
Indiana	6.2	New Mexico	1.9	West Virginia	1.8
Iowa	3.0	New York	19.2	Wisconsin	5.5
Kansas	2.7	North Carolina	8.5	Wyoming	0.5
Kentucky	4.1	North Dakota	0.6		

- Develop a frequency distribution, a percent frequency distribution, and a histogram. Use a class width of 2.5 million.
 - Discuss the skewness in the distribution.
 - What observations can you make about the population of the 50 states?
45. *Drug Store News* (September 2002) provided data on annual pharmacy sales for the leading pharmacy retailers in the United States. The following data are annual sales in millions.

Retailer	Sales	Retailer	Sales
Ahold USA	\$ 1700	Medicine Shoppe	\$ 1757
CVS	12700	Rite-Aid	8637
Eckerd	7739	Safeway	2150
Kmart	1863	Walgreens	11660
Kroger	3400	Wal-Mart	7250

- Show a stem-and-leaf display.
 - Identify the annual sales levels for the smallest, medium, and largest drug retailers.
 - What are the two largest drug retailers?
46. The daily high and low temperatures for 20 cities follow (*USA Today*, March 3, 2006).

WEB file
CityTemp

City	High	Low	City	High	Low
Albuquerque	66	39	Los Angeles	60	46
Atlanta	61	35	Miami	84	65
Baltimore	42	26	Minneapolis	30	11
Charlotte	60	29	New Orleans	68	50
Cincinnati	41	21	Oklahoma City	62	40
Dallas	62	47	Phoenix	77	50
Denver	60	31	Portland	54	38
Houston	70	54	St. Louis	45	27
Indianapolis	42	22	San Francisco	55	43
Las Vegas	65	43	Seattle	52	36

- a. Prepare a stem-and-leaf display of the high temperatures.
 - b. Prepare a stem-and-leaf display of the low temperatures.
 - c. Compare the two stem-and-leaf displays and make comments about the difference between the high and low temperatures.
 - d. Provide a frequency distribution for both high and low temperatures.
47. Refer to the data set for high and low temperatures for 20 cities in exercise 46.
 - a. Develop a scatter diagram to show the relationship between the two variables, high temperature and low temperature.
 - b. Comment on the relationship between high and low temperatures.
 48. One of the questions in a *Financial Times*/Harris Poll was, “How much do you favor or oppose a higher tax on higher carbon emission cars?” Possible responses were strongly favor, favor more than oppose, oppose more than favor, and strongly oppose. The following crosstabulation shows the responses obtained for 5372 adults surveyed in four countries in Europe and the United States (Harris Interactive website, February 27, 2008).

Level of Support	Country					Total
	Great Britain	Italy	Spain	Germany	United States	
Strongly favor	337	334	510	222	214	1617
Favor more than oppose	370	408	355	411	327	1871
Oppose more than favor	250	188	155	267	275	1135
Strongly oppose	130	115	89	211	204	749
Total	1087	1045	1109	1111	1020	5372

- a. Construct a percent frequency distribution for the level of support variable. Do you think the results show support for a higher tax on higher carbon emission cars?
 - b. Construct a percent frequency distribution for the country variable.
 - c. Does the level of support among adults in the European countries appear to be different than the level of support among adults in the United States? Explain.
49. Western University has only one women’s softball scholarship remaining for the coming year. The final two players that Western is considering are Allison Fealey and Emily Janson. The coaching staff has concluded that the speed and defensive skills are virtually identical for the two players, and that the final decision will be based on which player has the best batting average. Crosstabulations of each player’s batting performance in their junior and senior years of high school are as follows:

Outcome	Allison Fealey		Outcome	Emily Janson	
	Junior	Senior		Junior	Senior
Hit	15	75	Hit	70	35
No Hit	25	175	No Hit	130	85
Total At-Bats	40	250	Total At Bats	200	120

A player’s batting average is computed by dividing the number of hits a player has by the total number of at-bats. Batting averages are represented as a decimal number with three places after the decimal.

- a. Calculate the batting average for each player in her junior year. Then calculate the batting average of each player in her senior year. Using this analysis, which player should be awarded the scholarship? Explain.

- b. Combine or aggregate the data for the junior and senior years into one crosstabulation as follows:

Outcome	Player	
	Fealey	Janson
Hit		
No Hit		
Total At-Bats		

Calculate each player’s batting average for the combined two years. Using this analysis, which player should be awarded the scholarship? Explain.

- c. Are the recommendations you made in parts (a) and (b) consistent? Explain any apparent inconsistencies.
50. A survey of commercial buildings served by the Cincinnati Gas & Electric Company asked what main heating fuel was used and what year the building was constructed. A partial crosstabulation of the findings follows.

Year Constructed	Fuel Type				
	Electricity	Natural Gas	Oil	Propane	Other
1973 or before	40	183	12	5	7
1974–1979	24	26	2	2	0
1980–1986	37	38	1	0	6
1987–1991	48	70	2	0	1

- a. Complete the crosstabulation by showing the row totals and column totals.
 - b. Show the frequency distributions for year constructed and for fuel type.
 - c. Prepare a crosstabulation showing column percentages.
 - d. Prepare a crosstabulation showing row percentages.
 - e. Comment on the relationship between year constructed and fuel type.
51. Table 2.16 contains a portion of the data in the file named Fortune. Data on stockholders’ equity, market value, and profits for a sample of 50 Fortune 500 companies are shown.

TABLE 2.16 DATA FOR A SAMPLE OF 50 FORTUNE 500 COMPANIES

Company	Stockholders’ Equity (\$1000s)	Market Value (\$1000s)	Profit (\$1000s)
AGCO	982.1	372.1	60.6
AMP	2698.0	12017.6	2.0
Apple Computer	1642.0	4605.0	309.0
Baxter International	2839.0	21743.0	315.0
Bergen Brunswick	629.1	2787.5	3.1
Best Buy	557.7	10376.5	94.5
Charles Schwab	1429.0	35340.6	348.5
.	.	.	.
.	.	.	.
.	.	.	.
Walgreen	2849.0	30324.7	511.0
Westvaco	2246.4	2225.6	132.0
Whirlpool	2001.0	3729.4	325.0
Xerox	5544.0	35603.7	395.0

WEB file
Fortune

- a. Prepare a crosstabulation for the variables Stockholders' Equity and Profit. Use classes of 0–200, 200–400, . . . , 1000–1200 for Profit, and classes of 0–1200, 1200–2400, . . . , 4800–6000 for Stockholders' Equity.
 - b. Compute the row percentages for your crosstabulation in part (a).
 - c. What relationship, if any, do you notice between Profit and Stockholders' Equity?
52. Refer to the data set in Table 2.16.
- a. Prepare a crosstabulation for the variables Market Value and Profit.
 - b. Compute the row percentages for your crosstabulation in part (a).
 - c. Comment on any relationship between the variables.
53. Refer to the data set in Table 2.16.
- a. Prepare a scatter diagram to show the relationship between the variables Profit and Stockholders' Equity.
 - b. Comment on any relationship between the variables.
54. Refer to the data set in Table 2.16.
- a. Prepare a scatter diagram to show the relationship between the variables Market Value and Stockholders' Equity.
 - b. Comment on any relationship between the variables.

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 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 2.17 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.

TABLE 2.17 DATA FOR A SAMPLE OF 100 CREDIT CARD PURCHASES AT PELICAN STORES

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
.
.
96	Regular	1	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

WEB file
PelicanStores

Most of the variables shown in Table 2.17 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 tabular and graphical methods of descriptive statistics to help management develop a customer profile and to evaluate the promotional campaign. At a minimum, your report should include the following:

1. Percent frequency distribution for key variables.
2. A bar chart or pie chart showing the number of customer purchases attributable to the method of payment.
3. A crosstabulation of type of customer (regular or promotional) versus net sales. Comment on any similarities or differences present.
4. A scatter diagram to explore the relationship between net sales and customer age.

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 (\$millions), the total gross sales (\$millions), 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 2.18 shows the data for the first 10 motion pictures in this file.

Managerial Report

Use the tabular and graphical methods of descriptive statistics to learn how these variables contribute to the success of a motion picture. Include the following in your report.

TABLE 2.18 PERFORMANCE DATA FOR 10 MOTION PICTURES

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

WEB file
Movies

1. Tabular and graphical summaries for each of the four variables along with a discussion of what each summary tells us about the motion picture industry.
2. A scatter diagram to explore the relationship between Total Gross Sales and Opening Weekend Gross Sales. Discuss.
3. A scatter diagram to explore the relationship between Total Gross Sales and Number of Theaters. Discuss.
4. A scatter diagram to explore the relationship between Total Gross Sales and Number of Weeks in the Top 60. Discuss.

Appendix 2.1 Using Minitab for Tabular and Graphical Presentations

Minitab offers extensive capabilities for constructing tabular and graphical summaries of data. In this appendix we show how Minitab can be used to construct several graphical summaries and the tabular summary of a crosstabulation. The graphical methods presented include the dot plot, the histogram, the stem-and-leaf display, and the scatter diagram.

Dot Plot



We use the audit time data in Table 2.4 to demonstrate. The data are in column C1 of a Minitab worksheet. The following steps will generate a dot plot.

- Step 1.** Select the **Graph** menu and choose **Dotplot**
- Step 2.** Select **One Y, Simple** and click **OK**
- Step 3.** When the Dotplot-One Y, Simple dialog box appears:
Enter C1 in the **Graph Variables** box
Click **OK**

Histogram



We show how to construct a histogram with frequencies on the vertical axis using the audit time data in Table 2.4. The data are in column C1 of a Minitab worksheet. The following steps will generate a histogram for audit times.

- Step 1.** Select the **Graph** menu
- Step 2.** Choose **Histogram**
- Step 3.** Select **Simple** and click **OK**
- Step 4.** When the Histogram-Simple dialog box appears:
Enter C1 in the **Graph Variables** box
Click **OK**
- Step 5.** When the Histogram appears:
Position the mouse pointer over any one of the bars
Double-click
- Step 6.** When the Edit Bars dialog box appears:
Click on the **Binning** tab
Select **Cutpoint** for Interval Type
Select **Midpoint/Cutpoint positions** for Interval Definition
Enter 10:35/5 in the **Midpoint/Cutpoint positions** box*
Click **OK**

*The entry 10:35/5 indicates that 10 is the starting value for the histogram, 35 is the ending value for the histogram, and 5 is the class width.

Note that Minitab also provides the option of scaling the x -axis so that the numerical values appear at the midpoints of the histogram rectangles. If this option is desired, modify step 6 to include Select **Midpoint** for Interval Type and Enter 12:32/5 in the **Midpoint/Cutpoint positions** box. These steps provide the same histogram with the midpoints of the histogram rectangles labeled 12, 17, 22, 27, and 32.

Stem-and-Leaf Display



ApTest

We use the aptitude test data in Table 2.8 to demonstrate the construction of a stem-and-leaf display. The data are in column C1 of a Minitab worksheet. The following steps will generate the stretched stem-and-leaf display shown in Section 2.3.

- Step 1.** Select the **Graph** menu
- Step 2.** Choose **Stem-and-Leaf**
- Step 3.** When the Stem-and-Leaf dialog box appears:
Enter C1 in the **Graph Variables** box
Click **OK**

Scatter Diagram



Stereo

We use the stereo and sound equipment store data in Table 2.12 to demonstrate the construction of a scatter diagram. The weeks are numbered from 1 to 10 in column C1, the data for number of commercials are in column C2, and the data for sales are in column C3 of a Minitab worksheet. The following steps will generate the scatter diagram shown in Figure 2.7.

- Step 1.** Select the **Graph** menu
- Step 2.** Choose **Scatterplot**
- Step 3.** Select **Simple** and click **OK**
- Step 4.** When the Scatterplot-Simple dialog box appears:
Enter C3 under **Y variables** and C2 under **X variables**
Click **OK**

Crosstabulation



Restaurant

We use the data from Zagat's restaurant review, part of which is shown in Table 2.9, to demonstrate. The restaurants are numbered from 1 to 300 in column C1 of the Minitab worksheet. The quality ratings are in column C2, and the meal prices are in column C3.

Minitab can only create a crosstabulation for qualitative variables and meal price is a quantitative variable. So we need to first code the meal price data by specifying the class to which each meal price belongs. The following steps will code the meal price data to create four classes of meal price in column C4: \$10–19, \$20–29, \$30–39, and \$40–49.

- Step 1.** Select the **Data** menu
- Step 2.** Choose **Code**
- Step 3.** Choose **Numeric to Text**
- Step 4.** When the Code-Numeric to Text dialog box appears:
Enter C3 in the **Code data from columns** box
Enter C4 in the **Store coded data in columns** box
Enter 10:19 in the first **Original values** box and \$10-19 in the adjacent **New** box
Enter 20:29 in the second **Original values** box and \$20-29 in the adjacent **New** box

Enter 30:39 in the third **Original values** box and \$30-39 in the adjacent **New** box
 Enter 40:49 in the fourth **Original values** box and \$40-49 in the adjacent **New** box
 Click **OK**

For each meal price in column C3 the associated meal price category will now appear in column C4. We can now develop a crosstabulation for quality rating and the meal price categories by using the data in columns C2 and C4. The following steps will create a crosstabulation containing the same information as shown in Table 2.10.

Step 1. Select the **Stat** menu

Step 2. Choose **Tables**

Step 3. Choose **Cross Tabulation and Chi-Square**

Step 4. When the Cross Tabulation and Chi-Square dialog box appears:

Enter C2 in the **For rows** box and C4 in the **For columns** box

Select **Counts** under Display

Click **OK**

Appendix 2.2 Using Excel for Tabular and Graphical Presentations

Excel offers extensive capabilities for constructing tabular and graphical summaries of data. In this appendix, we show how Excel can be used to construct a frequency distribution, bar chart, pie chart, histogram, scatter diagram, and crosstabulation. We will demonstrate three of Excel's most powerful tools for data analysis: chart tools, PivotChart Report, and PivotTable Report.

Frequency Distribution and Bar Chart for Categorical Data

In this section we show how Excel can be used to construct a frequency distribution and a bar chart for categorical data. We illustrate each using the data on soft drink purchases in Table 2.1.

Frequency distribution We begin by showing how the COUNTIF function can be used to construct a frequency distribution for the data in Table 2.1. Refer to Figure 2.10 as we describe the steps involved. The formula worksheet (showing the functions and formulas used) is set in the background, and the value worksheet (showing the results obtained using the functions and formulas) appears in the foreground.

The label "Brand Purchased" and the data for the 50 soft drink purchases are in cells A1:A51. We also entered the labels "Soft Drink" and "Frequency" in cells C1:D1. The five soft drink names are entered into cells C2:C6. Excel's COUNTIF function can now be used to count the number of times each soft drink appears in cells A2:A51. The following steps are used.

Step 1. Select cell D2

Step 2. Enter =COUNTIF(\$A\$2:\$A\$51,C2)

Step 3. Copy cell D2 to cells D3:D6

The formula worksheet in Figure 2.10 shows the cell formulas inserted by applying these steps. The value worksheet shows the values computed by the cell formulas. This worksheet shows the same frequency distribution that we developed in Table 2.2.



FIGURE 2.10 FREQUENCY DISTRIBUTION FOR SOFT DRINK PURCHASES
CONSTRUCTED USING EXCEL'S COUNTIF FUNCTION

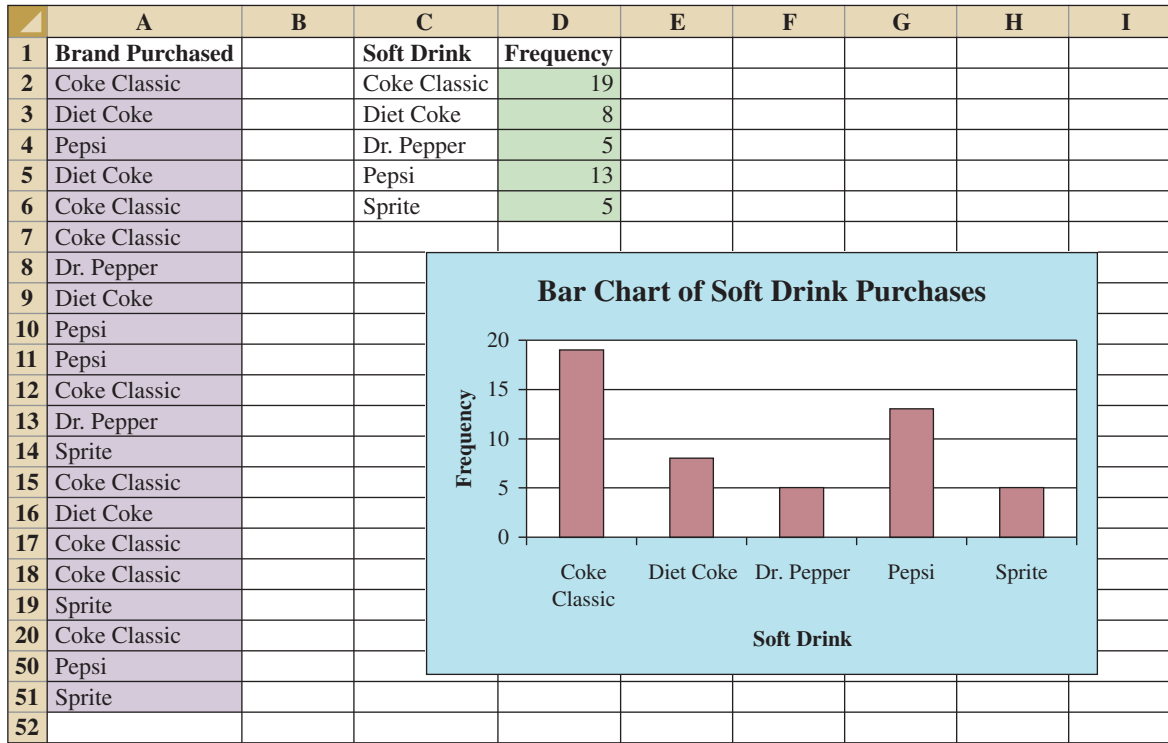
	A	B	C	D	E
1	Brand Purchased		Soft Drink	Frequency	
2	Coke Classic		Coke Classic	=COUNTIF(\$A\$2:\$A\$51,C2)	
3	Diet Coke		Diet Coke	=COUNTIF(\$A\$2:\$A\$51,C3)	
4	Pepsi		Dr. Pepper	=COUNTIF(\$A\$2:\$A\$51,C4)	
5	Diet Coke		Pepsi	=COUNTIF(\$A\$2:\$A\$51,C5)	
6	Coke Classic		Sprite	=COUNTIF(\$A\$2:\$A\$51,C6)	
7	Coke Classic				
8	Dr. Pepper				
9	Diet Coke	1	Brand Purchased	Soft Drink	Frequency
10	Pepsi	2	Coke Classic	Coke Classic	19
45	Pepsi	3	Diet Coke	Diet Coke	8
46	Pepsi	4	Pepsi	Dr. Pepper	5
47	Pepsi	5	Diet Coke	Pepsi	13
48	Coke Classic	6	Coke Classic	Sprite	5
49	Dr. Pepper	7	Coke Classic		
50	Pepsi	8	Dr. Pepper		
51	Sprite	9	Diet Coke		
52		10	Pepsi		
		45	Pepsi		
		46	Pepsi		
		47	Pepsi		
		48	Coke Classic		
		49	Dr. Pepper		
		50	Pepsi		
		51	Sprite		
		52			

Note: Rows 11–44 are hidden.



Bar chart Here we show how Excel's chart tools can be used to construct a bar chart for the soft drink data. Refer to the frequency distribution shown in the value worksheet of Figure 2.10. The bar chart that we are going to develop is an extension of this worksheet. The worksheet and the bar chart developed are shown in Figure 2.11. The steps are as follows:

- Step 1.** Select cells C2:D6
- Step 2.** Click the **Insert** tab on the Ribbon
- Step 3.** In the **Charts** group, click **Column**
- Step 4.** When the list of column chart subtypes appears:
 - Go to the **2-D Column** section
 - Click **Clustered Column** (the leftmost chart)
- Step 5.** In the **Chart Layouts** group, click the **More** button (the downward-pointing arrow with a line over it) to display all the options
- Step 6.** Choose **Layout 9**
- Step 7.** Select the **Chart Title** and replace it with **Bar Chart of Soft Drink Purchases**
- Step 8.** Select the **Horizontal (Category) Axis Title** and replace it with **Soft Drink**
- Step 9.** Select the **Vertical (Value) Axis Title** and replace it with **Frequency**
- Step 10.** Right-click the **Series 1 Legend Entry**
 - Click **Delete**
- Step 11.** Right-click the vertical axis
 - Click **Format Axis**

FIGURE 2.11 BAR CHART OF SOFT DRINK PURCHASES CONSTRUCTED USING EXCEL'S CHART TOOLS

Step 12. When the Format Axis dialog box appears:

Go to the **Axis Options** section

Select **Fixed** for **Major Unit** and enter 5.0 in the corresponding box

Click **Close**

The resulting bar chart is shown in Figure 2.11.*

Excel can produce a pie chart for the soft drink data in a similar fashion. The major difference is that in step 3 we would click **Pie** in the **Charts** group. Several style pie charts are available.

Frequency Distribution and Histogram for Quantitative Data

In a later section of this appendix we describe how to use Excel's PivotTable Report to construct a crosstabulation.

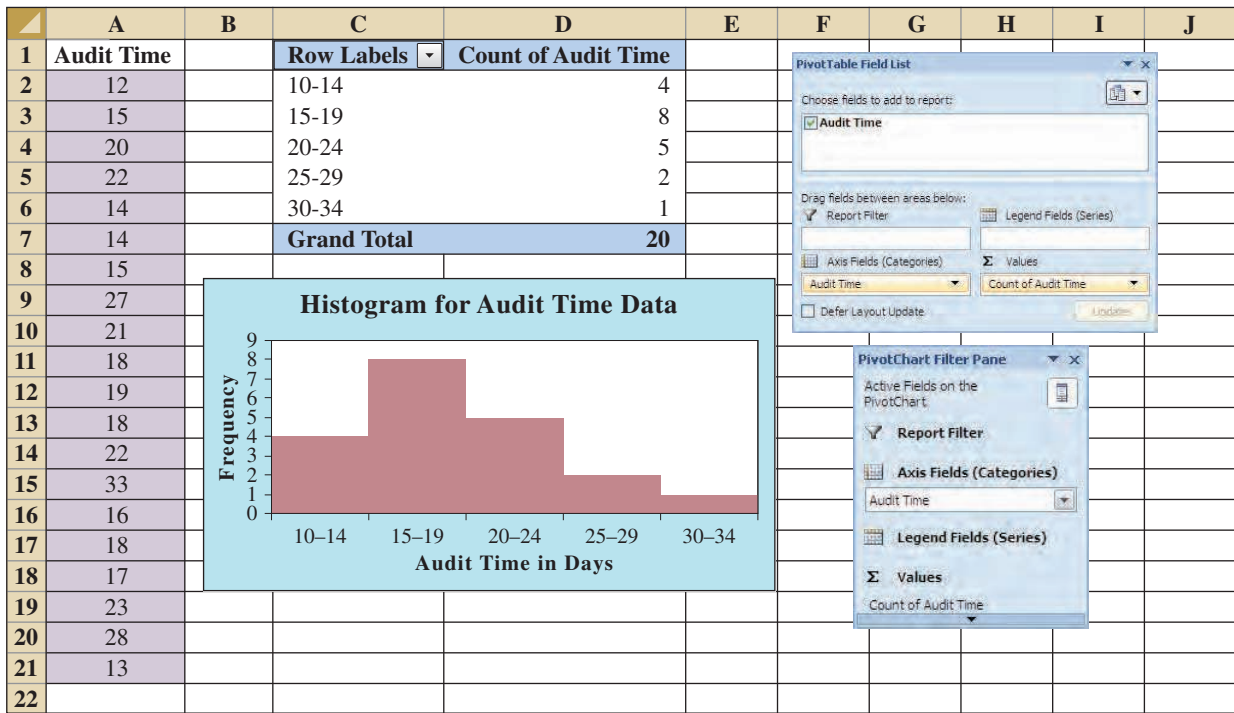
WEB file

Audit

Excel's PivotTable Report is an interactive tool that allows you to quickly summarize data in a variety of ways, including developing a frequency distribution for quantitative data. Once a frequency distribution is created using the PivotTable Report, Excel's chart tools can then be used to construct the corresponding histogram. But, using Excel's PivotChart Report, we can construct a frequency distribution and a histogram simultaneously. We will illustrate this procedure using the audit time data in Table 2.4. The label "Audit Time" and the 20 audit time values are entered into cells A1:A21 of an Excel worksheet. The following steps describe how to use Excel's PivotChart Report to construct a frequency distribution and a histogram for the audit time data. Refer to Figure 2.12 as we describe the steps involved.

*The bar chart in Figure 2.11 can be resized. Resizing an Excel chart is not difficult. First, select the chart. Sizing handles will appear on the chart border. Click on the sizing handles and drag them to resize the figure to your preference.

FIGURE 2.12 USING EXCEL'S PIVOTCHART REPORT TO CONSTRUCT A FREQUENCY DISTRIBUTION AND HISTOGRAM FOR THE AUDIT TIME DATA



Step 1. Click the **Insert** tab on the Ribbon

Step 2. In the **Tables** group, click the word **PivotTable**

Step 3. Choose **PivotChart** from the options that appear

Step 4. When the **Create PivotTable with PivotChart** dialog box appears,

Choose **Select a table or range**

Enter A1:A21 in the **Table/Range** box

Choose **Existing Worksheet** as the location for the PivotTable and PivotChart

Enter C1 in the **Location** box

Click **OK**

Step 5. In the **PivotTable Field List**, go to **Choose Fields to add to report**

Drag the **Audit Time** field to the **Axis Fields (Categories)** area

Drag the **Audit Time** field to the **Values** area

Step 6. Click **Sum of Audit Time** in the **Values** area

Step 7. Click **Value Field Settings** from the list of options that appears

Step 8. When the Value Field Settings dialog appears,

Under **Summarize value field by**, choose **Count**

Click **OK**

Step 9. Close the **PivotTable Field List**.

Step 10. Right-click cell C2 in the PivotTable report or any other cell containing an audit time

Step 11. Choose **Group** from the list of options that appears

Step 12. When the **Grouping** dialog box appears,

Enter 10 in the **Starting at** box

Enter 34 in the **Ending at** box

Enter 5 in the **By** box

Click **OK** (a PivotChart will appear)

Step 13. Click inside the resulting PivotChart

Step 14. Click the **Design** tab on the Ribbon

Step 15. In the **Chart Layouts** group, click the **More** button (the downward pointing arrow with a line over it) to display all the options

Step 16. Choose **Layout 8**

Step 17. Select the **Chart Title** and replace it with **Histogram for Audit Time Data**

Step 18. Select the **Horizontal (Category) Axis Title** and replace it with **Audit Time in Days**

Step 19. Select the **Vertical (Value) Axis Title** and replace it with **Frequency**

Figure 2.12 shows the resulting PivotTable and PivotChart. We see that the PivotTable report provides the frequency distribution for the audit time data and the PivotChart provides the corresponding histogram. If desired, we can change the labels in any cell in the frequency distribution by selecting the cell and typing in the new label.

Crosstabulation

Excel's PivotTable Report provides an excellent way to summarize the data for two or more variables simultaneously. We will illustrate the use of Excel's PivotTable Report by showing how to develop a crosstabulation of quality ratings and meal prices for the sample of 300 Los Angeles restaurants. We will use the data in the file named Restaurant; the labels "Restaurant," "Quality Rating," and "Meal Price (\$)" have been entered into cells A1:C1 of the worksheet as shown in Figure 2.13. The data for each of the restaurants in the sample have been entered into cells B2:C301.

FIGURE 2.13 EXCEL WORKSHEET CONTAINING RESTAURANT DATA



Note: Rows 12–291 are hidden.

	A	B	C	D
1	Restaurant	Quality Rating	Meal Price (\$)	
2	1	Good	18	
3	2	Very Good	22	
4	3	Good	28	
5	4	Excellent	38	
6	5	Very Good	33	
7	6	Good	28	
8	7	Very Good	19	
9	8	Very Good	11	
10	9	Very Good	23	
11	10	Good	13	
292	291	Very Good	23	
293	292	Very Good	24	
294	293	Excellent	45	
295	294	Good	14	
296	295	Good	18	
297	296	Good	17	
298	297	Good	16	
299	298	Good	15	
300	299	Very Good	38	
301	300	Very Good	31	
302				

In order to use the Pivot Table report to create a crosstabulation, we need to perform three tasks: Display the Initial PivotTable Field List and PivotTable Report; Set Up the PivotTable Field List; and Finalize the PivotTable Report. These tasks are described as follows.

Display the Initial PivotTable Field List and PivotTable Report: Three steps are needed to display the initial PivotTable Field List and PivotTable report.

Step 1. Click the **Insert** tab on the Ribbon

Step 2. In the **Tables** group, click the icon above the word PivotTable

Step 3. When the **Create PivotTable** dialog box appears,

Choose **Select a Table or Range**

Enter A1:C301 in the **Table/Range** box

Choose **New Worksheet** as the location for the PivotTable Report

Click **OK**

The resulting initial PivotTable Field List and PivotTable Report are shown in Figure 2.14.

Set Up the PivotTable Field List: Each of the three columns in Figure 2.13 (labeled Restaurant, Quality Rating, and Meal Price (\$)) is considered a field by Excel. Fields may be chosen to represent rows, columns, or values in the body of the PivotTable Report. The following steps show how to use Excel's PivotTable Field List to assign the Quality Rating field to the rows, the Meal Price (\$) field to the columns, and the Restaurant field to the body of the PivotTable report.

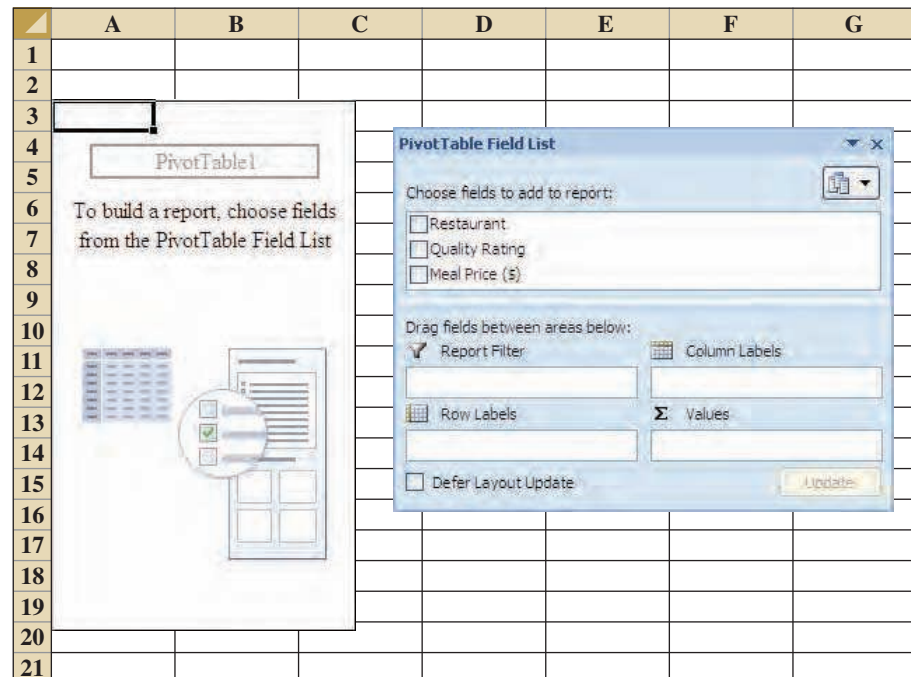
Step 1. In the **PivotTable Field List**, go to **Choose Fields to add to report**

Drag the **Quality Rating** field to the **Row Labels** area

Drag the **Meal Price (\$)** field to the **Column Labels** area

Drag the **Restaurant** field to the **Values** area

FIGURE 2.14 INITIAL PIVOTTABLE FIELD LIST AND PIVOTTABLE FIELD REPORT FOR THE RESTAURANT DATA



- Step 2.** Click on **Sum of Restaurant** in the **Values** area
- Step 3.** Click **Value Field Settings** from the list of options that appear
- Step 4.** When the Value Field Settings dialog appears,
Under **Summarize value field by**, choose **Count**
Click **OK**

Figure 2.15 shows the completed PivotTable Field List and a portion of the PivotTable worksheet as it now appears.

Finalize the PivotTable Report To complete the PivotTable Report we need to group the columns representing meal prices and place the row labels for quality rating in the proper order. The following steps accomplish this.

- Step 1.** Right-click in cell B4 or any cell containing meal prices
- Step 2.** Choose **Group** from the list of options that appears
- Step 3.** When the **Grouping** dialog box appears,
Enter 10 in the **Starting at** box
Enter 49 in the **Ending at** box
Enter 10 in the **By** box
Click **OK**
- Step 4.** Right-click on **Excellent** in cell A5
- Step 5.** Choose **Move** and click **Move “Excellent” to End**

The final PivotTable Report is shown in Figure 2.16. Note that it provides the same information as the crosstabulation shown in Table 2.10.

Scatter Diagram

We can use Excel’s chart tools to construct a scatter diagram and a trend line for the stereo and sound equipment store data presented in Table 2.12. Refer to Figures 2.17 and 2.18 as

FIGURE 2.15 COMPLETED PIVOTTABLE FIELD LIST AND A PORTION OF THE PIVOTTABLE REPORT FOR THE RESTAURANT DATA (COLUMNS H:AK ARE HIDDEN)

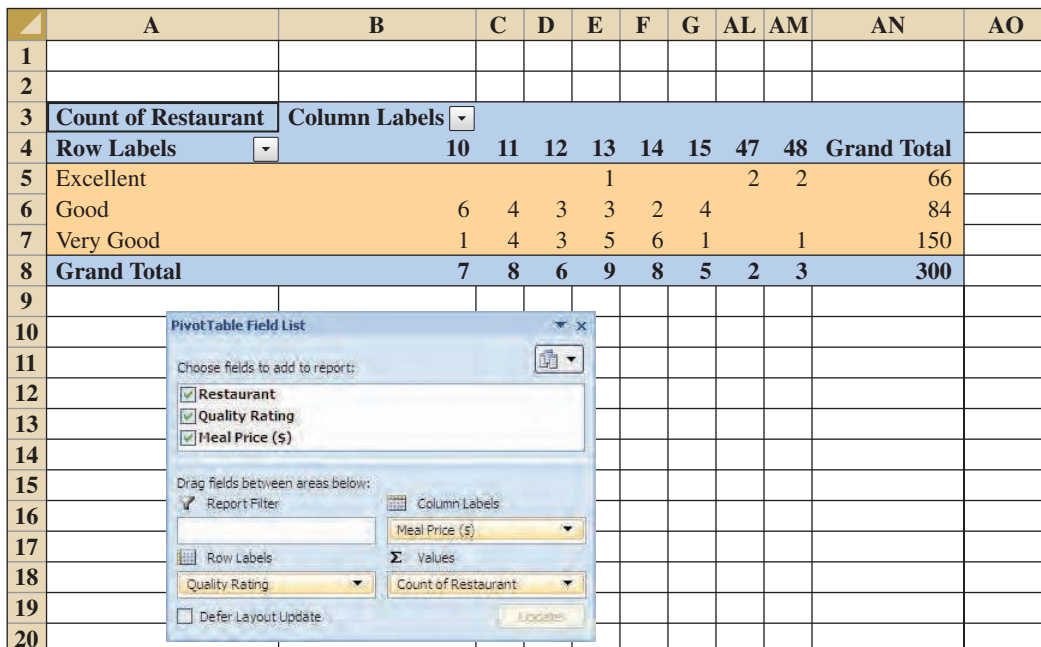


FIGURE 2.16 FINAL PIVOTTABLE REPORT FOR THE RESTAURANT DATA

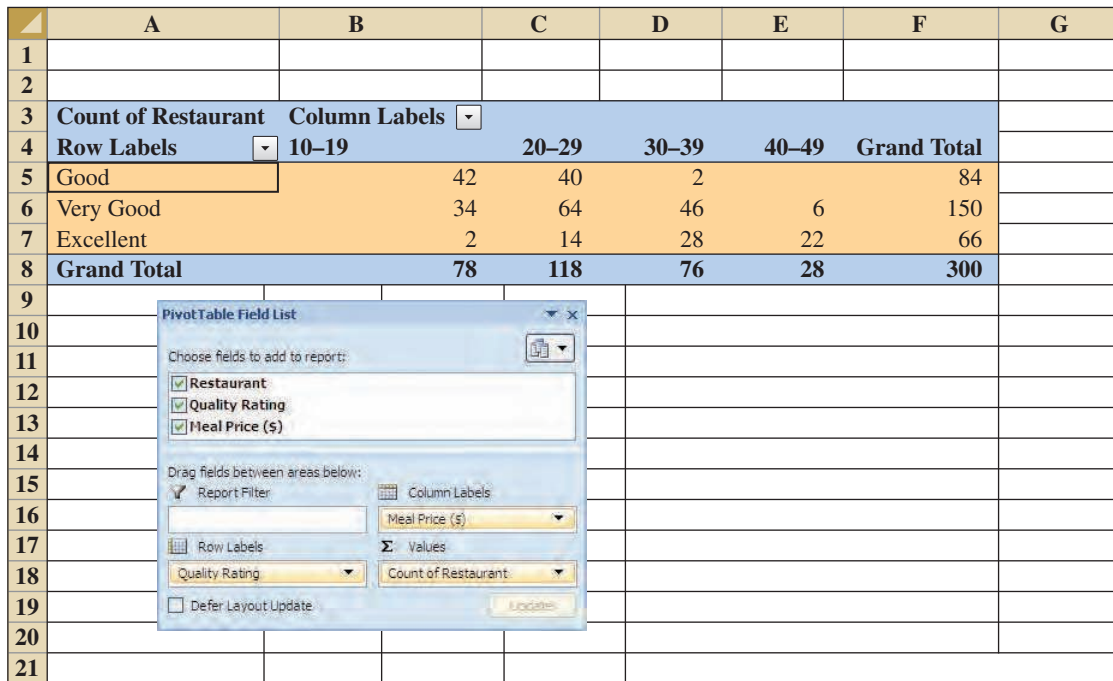


FIGURE 2.17 SCATTER DIAGRAM FOR THE STEREO AND SOUND EQUIPMENT STORE USING EXCEL'S CHART TOOLS

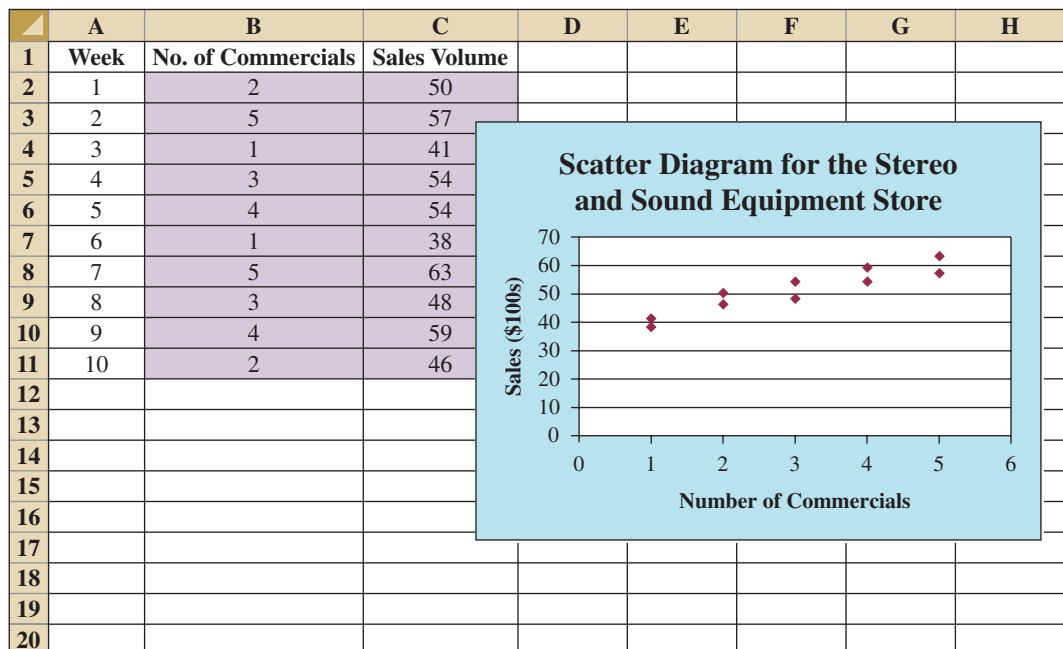
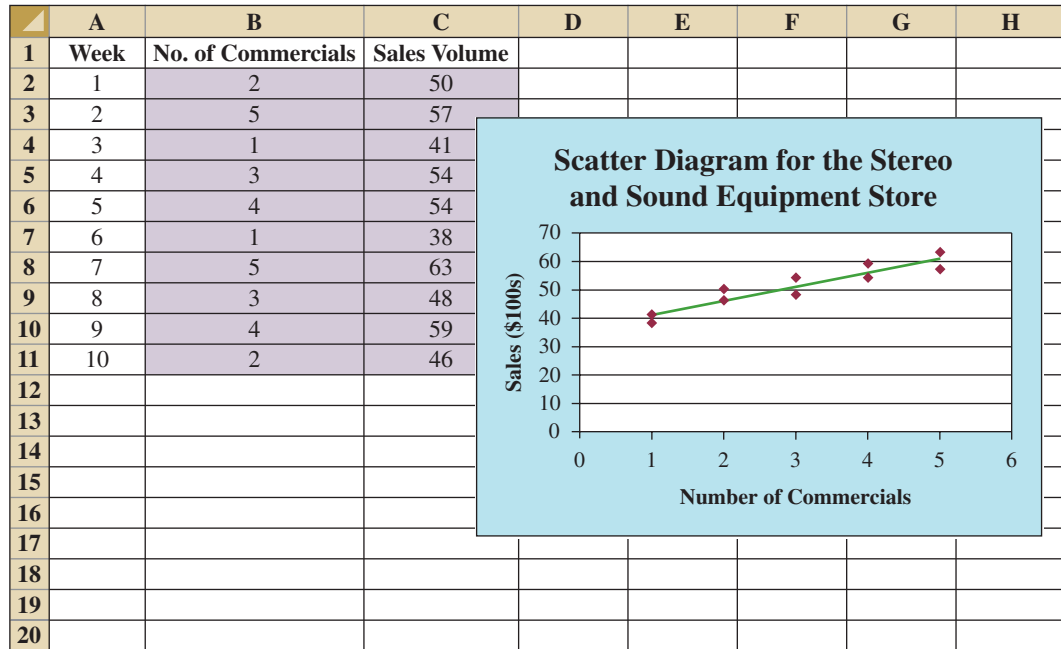


FIGURE 2.18 SCATTER DIAGRAM AND TRENDLINE FOR THE STEREO AND SOUND EQUIPMENT STORE USING EXCEL'S CHART TOOLS



we describe the steps involved. We will use the data in the file named Stereo; the labels Week, No. of Commercials, and Sales Volume have been entered into cells A1:C1 of the worksheet. The data for each of the 10 weeks are entered into cells B2:C11. The following steps describe how to use Excel's chart tools to produce a scatter diagram for the data.

- Step 1.** Select cells B2:C11
- Step 2.** Click the **Insert** tab on the Ribbon
- Step 3.** In the **Charts** group, click **Scatter**
- Step 4.** When the list of scatter diagram subtypes appears, click **Scatter with only Markers** (the chart in the upper left corner)
- Step 5.** In the **Chart Layouts** group, click **Layout 1**
- Step 6.** Select the **Chart Title** and replace it with **Scatter Diagram for the Stereo and Sound Equipment Store**
- Step 7.** Select the **Horizontal (Value) Axis Title** and replace it with **Number of Commercials**
- Step 8.** Select the **Vertical (Value) Axis Title** and replace it with **Sales (\$100s)**
- Step 9.** Right-click the **Series 1 Legend Entry** and click **Delete**

The worksheet displayed in Figure 2.17 shows the scatter diagram produced by Excel. The following steps describe how to add a trendline.

- Step 1.** Position the mouse pointer over any data point in the scatter diagram and right-click to display a list of options
- Step 2.** Choose **Add Trendline**
- Step 3.** When the **Format Trendline** dialog box appears,
 - Select **Trendline Options**
 - Choose **Linear** from the **Trend/Regression Type** list
 - Click **Close**

The worksheet displayed in Figure 2.18 shows the scatter diagram with the trendline added.

Appendix 2.3 Using StatTools for Tabular and Graphical Presentations

In this appendix we show how StatTools can be used to construct a histogram and a scatter diagram.

Histogram

We use the audit time data in Table 2.4 to illustrate. Begin by using the Data Set Manager to create a StatTools data set for these data using the procedure described in the appendix in Chapter 1. The following steps will generate a histogram.



- Step 1.** Click the **StatTools** tab on the Ribbon
- Step 2.** In the **Analyses Group**, click **Summary Graphs**
- Step 3.** Choose the **Histogram** option
- Step 4.** When the StatTools—Histogram dialog box appears,
 - In the **Variables** section, select **Audit Time**
 - In the **Options** section,
 - Enter 5 in the **Number of Bins** box
 - Enter 9.5 in the **Histogram Minimum** box
 - Enter 34.5 in the **Histogram Maximum** box
 - Choose **Categorical** in the **X-Axis** box
 - Choose **Frequency** in the **Y-Axis** box
 - Click **OK**

A histogram for the audit time data similar to the histogram shown in Figure 2.12 will appear. The only difference is the histogram developed using StatTools shows the class mid-points on the horizontal axis.

Scatter Diagram

We use the stereo and sound equipment data in Table 2.12 to demonstrate the construction of a scatter diagram. Begin by using the Data Set Manager to create a StatTools data set for these data using the procedure described in the appendix in Chapter 1. The following steps will generate a scatter diagram.



- Step 1.** Click the **StatTools** tab on the Ribbon
- Step 2.** In the **Analyses Group**, click **Summary Graphs**
- Step 3.** Choose the **Scatterplot** option
- Step 4.** When the StatTools—Scatterplot dialog box appears,
 - In the **Variables** section,
 - In the column labeled **X**, select **No. of Commercials**
 - In the column labeled **Y**, select **Sales Volume**
 - Click **OK**

A scatter diagram similar to the one shown in Figure 2.17 will appear.