

Central Tendency and Variability

Central Tendency

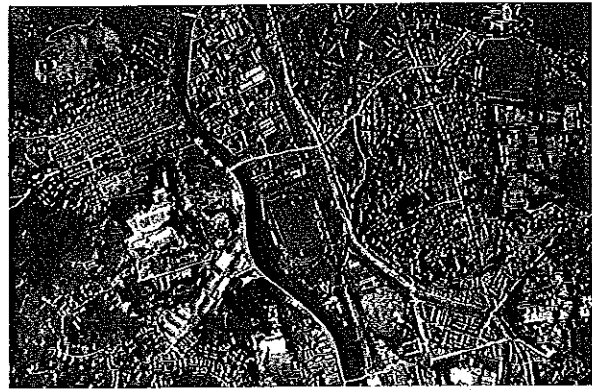
- Mean, the Arithmetic Average
- Median, the Middle Score
- Mode, the Most Common Score
- How Outliers Affect Measures of Central Tendency
- Which Measure of Central Tendency Is Best?

Measures of Variability

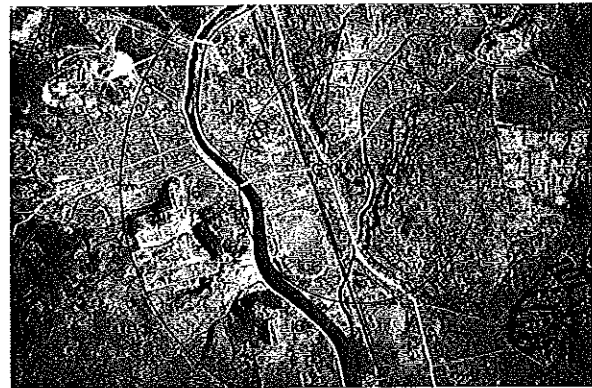
- Range
- Variance
- Standard Deviation

BEFORE YOU GO ON

- You should understand what a distribution is (Chapter 2).
- You should be able to explain histograms and frequency polygons (Chapter 2).



Nagasaki, Two Days Before the Atomic Bomb



Nagasaki, Three Days After the Atomic Bomb

On August 9, 1945, chance variability in the cloud cover diverted a B-29 bomber from Kokura, Japan, to its secondary target, the city of Nagasaki. When the atomic bomb exploded a few hundred feet above a tennis court, all of the buildings and most of the people who lived in the city of Nagasaki simply disappeared; the people in Kokura survived. Chance variability matters.

How does any nation recover from such devastation? In 1950, an American statistician named W. Edwards Deming persuaded Japan's leading engineers and businesspeople that a statistical idea could re-create their entire industrial-based economy: variability. Deming's core statistical insight was that people were happy to pay for cars, kitchen appliances, and electronics with high reliability (low variability).

The Japanese industrial leadership embraced Deming, as well as his idea that it was management's job to reduce anything that contributes to product variability (an unreliable product). In manufacturing, variability might be due to using different suppliers because they submitted the lowest bid, using worn-out machinery to save money in the short term, or making working conditions unpleasant for employees.

Deming provided practical, statistical guidelines so that Japanese businesses could lower product variability. As Japan's industrial leaders applied Deming's statistical insight, they quickly discovered that controlling variability could be translated into thousands of different manufacturing solutions. The insight transformed the reputation of Japanese companies as manufacturers of cheap junk into one of manufacturers of high-quality products. To this day, the Japanese are specific about how they transformed their devastated nation from an economic disaster to an industrial leader: W. Edwards Deming.

Deming's statistical approach to manufacturing centered on variability, and in fact, variability is one of the basic building blocks of most statistical techniques. In this chapter, we learn about three common measures of variability in a distribution: range, variance, and standard deviation. But to fully understand variability, we first have to know how to identify the middle of a distribution. So before we learn about variability, we first introduce three measures of the middle of a distribution, or the central tendency: mean, median, and mode.

Central Tendency

▶ MASTERING THE CONCEPT

4-1: Central tendency is one of the most important ways to understand a distribution of data. We can use the mean, median, or mode as an indicator of central tendency.

Central tendency refers to the descriptive statistic that best represents the center of a data set, the particular value that all the other data seem to be gathering around. It's what we mean when we refer to the "typical" score. Simply creating a visual representation of the distribution, as we did in Chapter 2, often reveals its central tendency. The central tendency is usually at (or near) the highest point in the histogram or the polygon (Figure 4-1), but the specific way that data cluster around a distribution's central tendency can be measured three different ways: mean, median, and mode.

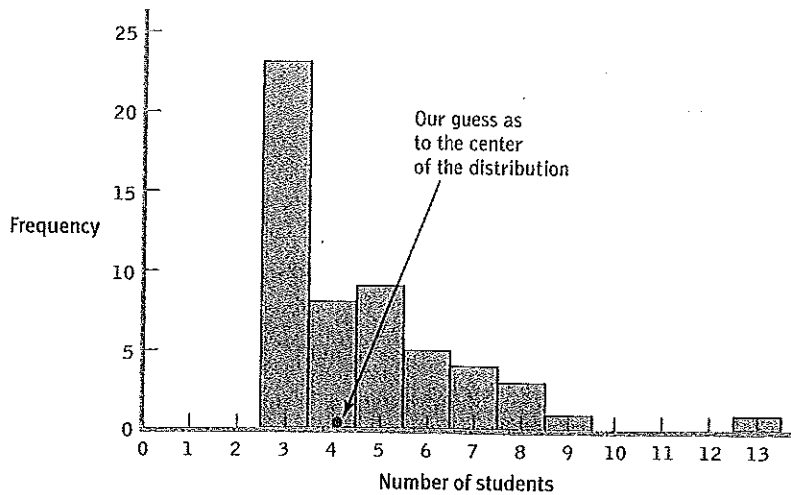


FIGURE 4-1
Estimating Central Tendency with Histograms

Histograms and frequency polygons allow us to see the likely center of our sample's distribution. The arrow points to our guess as to the center of the distribution of numbers of students mentored by chemistry professors.

Mean, the Arithmetic Average

The mean is simple to calculate and is the gateway to understanding statistical formulas. The mean is such an important concept in statistics that we provide you with four distinct ways to think about it: verbally, arithmetically, visually, and symbolically (using statistical notation).

The Mean in Plain English The most commonly reported measure of central tendency is *the mean, the arithmetic average of a group of scores*. The mean, often called the average, is used to represent the "typical" score in a distribution. This is different from the way we often use the word *average* in everyday conversation. We may refer to a person as average in a somewhat derogatory way, noting that someone is "just" average in athletic ability or a movie was "only" average. The word *average* connotes so many different shades of meaning that we need to define the mean arithmetically.

The Mean in Plain Arithmetic The mean is calculated by summing all the scores in a data set and then dividing this sum by the total number of scores. You likely have calculated means many times in your life.

For example, when we explore the numbers of students mentored by the 54 graduate advisors that we considered in Chapter 2, the mean would be calculated by first adding the number of students for each mentor, then dividing by the total number of mentors.

EXAMPLE 4.1

STEP 1. Add all of the scores together.

$$\begin{aligned}
 &3 + 3 + 3 + 4 + 5 + 9 + 5 + 3 + 3 + 5 + 6 + 3 + 4 + 8 + 6 + 3 + 3 + 3 + 4 \\
 &+ 4 + 4 + 7 + 6 + 3 + 5 + 5 + 7 + 13 + 3 + 3 + 3 + 3 + 3 + 4 + 4 + 4 + 5 \\
 &+ 6 + 7 + 6 + 7 + 8 + 8 + 3 + 3 + 3 + 3 + 5 + 3 + 3 + 5 + 3 + 5 + 3 + 3 = 250
 \end{aligned}$$

STEP 2. Divide the sum of all scores by the total number of scores.

In this case, the sum of all scores is 250, which is then divided by 54, the number of scores in this sample:

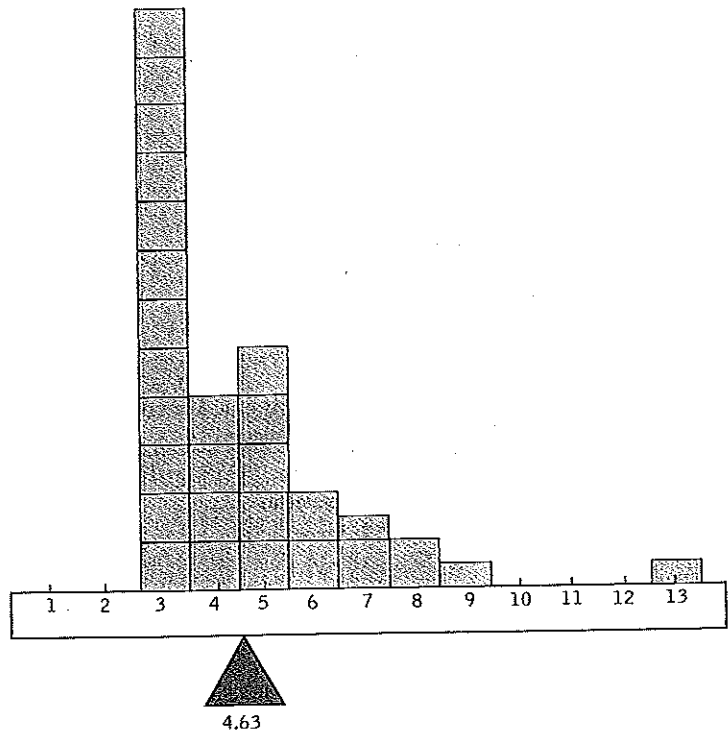
$$250/54 = 4.63$$

☐ **Central tendency** refers to the descriptive statistic that best represents the center of a data set, the particular value that all the other data seem to be gathering around.

☐ **The mean** is the arithmetic average of a group of scores. It is calculated by summing all the scores and dividing by the total number of scores.

FIGURE 4-2
The Mean as the Fulcrum of Our Data

The mean, 4.63, is the balancing point for all the scores for “numbers of students mentored” for the advisors in our sample. Mathematically, the scores always balance around the mean for any sample.



Visual Representations of the Mean Think of the mean as the visual point that perfectly balances two sides of a distribution. For example, the mean of 4.63 “students mentored” is represented visually as the point that perfectly balances that distribution, shown in the histogram in Figure 4-2.

The Mean Expressed by Symbolic Notation Symbolic notation may sound far more difficult to understand than it really is. After all, you just calculated a mean without symbolic notation and without a formula. Fortunately, we only need to understand a handful of symbols to express the ideas necessary to understand statistics.

Here are the several symbols that represent the mean. For the mean of a sample, statisticians typically use M or \bar{X} . In this book, we use M ; many other books also use M , but some use \bar{X} (pronounced “X bar”). For a population, statisticians use the Greek letter μ (pronounced “mew”) to symbolize the mean. (Although there are exceptions, Latin letters such as M tend to refer to numbers based on samples, and Greek letters such as μ tend to refer to numbers based on populations.) *The numbers based on samples are called statistics; M is a statistic. The numbers based on populations are called parameters; μ is a parameter.* Table 4-1 summarizes how these terms are used. As shown in

- A **statistic** is a number based on a sample taken from a population; statistics are usually symbolized by Latin letters.
- A **parameter** is a number based on the whole population; parameters are usually symbolized by Greek letters.

TABLE 4-1. The Mean in Symbols

The mean of a sample is an example of a statistic, whereas the mean of a population is an example of a parameter. The symbols we use depend on whether we are referring to the mean of a sample or a population.

Number	Used for	Symbol	Pronounced
Statistic	Sample	M or \bar{X}	“M” or “X bar”
Parameter	Population	μ	“mew”

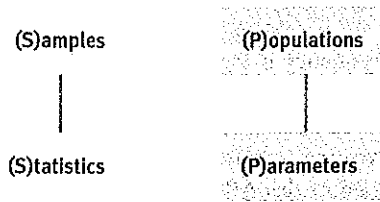


FIGURE 4-3

Try using a mnemonic trick to remember the distinction between samples and parameters. The letter *s* means that numbers based on (*s*)amples are called (*s*)tatistics. The letter *p* means that numbers based on (*p*)opulations are called (*p*)arameters.

Figure 4-3, you can remember this distinction by the first letters of these words: *statistic* and *sample* both begin with *s*, and *parameter* and *population* both begin with *p*. These symbols are part of the language of statistics and help us to communicate with other statisticians.

A formula to calculate the mean of a sample would use the symbol M on the left side of the equation. The right side would provide information on the actual calculation of the mean. A single score is typically symbolized as X . We know that we're summing all the scores, all the X 's, so the first step is to use the summation sign, Σ (pronounced "sigma"), to indicate that we're summing a list of scores. As you might guess, the full expression for summing all the scores would be ΣX . This instructs us to add up all of the X 's in the sample.

Step 1: Add up all of the scores in the sample. In statistical notation, this is ΣX .

Step 2: Divide the sum of all of the scores by the total number of scores. The total number of scores in a sample is typically represented by N . (Note that the capital letter N is typically used when we refer to the number of scores in the entire data set; if we break the sample down into smaller parts, as we'll see in later chapters, we typically use the lowercase letter n for the number of scores in each part.) The full equation would be:

$$M = \frac{\Sigma X}{N}$$

MASTERING THE FORMULA

4-1: The formula for the mean is:

$M = \frac{\Sigma X}{N}$. To calculate the mean, we add up every score, then divide by the total number of scores.

Let's look at the mean for the mentoring data that we considered earlier in this section.

EXAMPLE 4.2

STEP 1. Add up every score.

The sum of all scores is 250.

STEP 2. Divide by the total number of scores.

The total number of scores is 54. So, if we divide 250 by 54, the result is 4.63. Here's how it would look as a formula:

$$M = \frac{\Sigma X}{N} = \frac{250}{54} = 4.63$$

Statisticians tend to be as specific with their symbols as they are with their words. For example, almost all symbols are italicized, whether in the formulas to calculate statistics or in the reporting of statistics. However, the actual numerical values of the statistics are not italicized. Furthermore, whether or not a symbol is capitalized usually has meaning. Changing a symbol from uppercase to lowercase often changes what it means. When you practice calculating means, use this formula, being sure to italicize the symbols and use capital letters for M , X , and N .

Median, the Middle Score

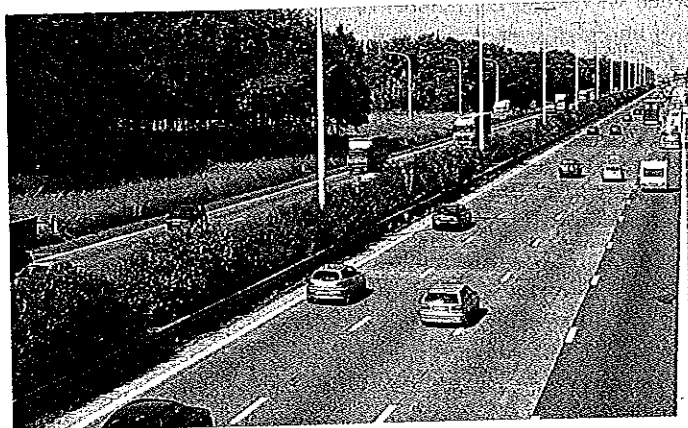
The second most common measure of central tendency is the median. *The median is the middle score of all the scores in a sample when the scores are arranged in ascending order.* We can think of the median as the 50th percentile. The median does not tend to be denoted by a symbol, although in APA style, the writing style of the American Psychological Association, it can be abbreviated as *mdn*. (Note that APA style, despite the word *psychological* in its name, is used across many of the social sciences; you are likely to use it in your courses regardless of your social science major.) To determine the median:

Step 1: Line up all the scores in ascending order.

Step 2: Find the middle score. With an odd number of scores, there will be an actual middle score. With an even number of scores, there will be no actual middle score. In this case, take the mean of the two middle scores.

Here are more specific instructions for finding the median. Keep in mind that with a distribution of only a few data points, we won't want to use the formula—just count how many numbers there are in the distribution and find the score that has the same number of scores above it and below it. Even with a distribution with many scores, the calculation is easy. All we do is divide the number of scores (N) by 2 and add $\frac{1}{2}$ —that is, 0.5. That number is the ordinal position (rank) of the median, or middle score. As illustrated below, simply count that many places over from the start of your scores and report that number.

Mean versus Median The median is the part of the roadway that divides the directions in which vehicles are permitted to drive. It can be dangerous to confuse the mean and the median, especially when you are calculating the "middle" of the roadway!



David De Looze/PhotoDisc, Green/Getty Images

EXAMPLE 4.3

Here is an example with an odd number of scores (representing numbers of students mentored for a sample of nine professors):

3, 4, 9, 4, 7, 7, 6, 3, 3

STEP 1. Line up the scores in ascending order:

3, 3, 3, 4, 4, 6, 7, 7, 9

STEP 2. Find the middle score.

First we count the total number of scores. There are 9 scores: $9/2 = 4.5$. If we add 0.5 to this result, we get 5. Therefore, the median is the 5th score. We now count across to the 5th score. The median is 4.

Here is an example with an even number of scores. Using the same data as in the example with the odd number of scores, we add one more data point: 13.

EXAMPLE 4.4**STEP 1. Line up the scores in ascending order.**

3, 3, 3, 4, 4, 6, 7, 7, 9, 13

STEP 2. Find the middle score.

First we count the total number of scores. There are 10 scores. We then divide the number of scores by two: $10/2 = 5$. If we add 0.5 to this result, we get 5.5; therefore, the median is the average of the 5th and 6th scores. The 5th and 6th scores are 4 and 6. The median is their mean, the mean of 4 and 6 = 5.

Mode, the Most Common Score

The *mode* is perhaps the easiest of the three measures of central tendency to calculate. *The mode is the most common score of all the scores in a sample.* It is readily picked out on a frequency table, histogram, or frequency polygon. Like the median, the mode does not tend to be represented by a symbol. It does not even have an APA abbreviation. When reporting modes, we use the word itself (e.g., the mode is . . .).

Determine the mode for the data for 54 graduate advisors from earlier in this chapter. Remember that each score represents the number of that graduate advisor's former students who went on to top jobs. The mode can be found either by searching the list of numbers for the most common score or by constructing a frequency table:

EXAMPLE 4.5

3, 3, 3, 4, 5, 9, 5, 3, 3, 5, 6, 3, 4, 8, 6, 3, 3, 3, 4, 4, 4, 7, 6, 3, 5, 5, 7, 13, 3, 3, 3, 3, 3, 4,
4, 4, 5, 6, 7, 6, 7, 8, 8, 3, 3, 3, 5, 3, 3, 5, 3, 5, 3, 3

Mode: _____

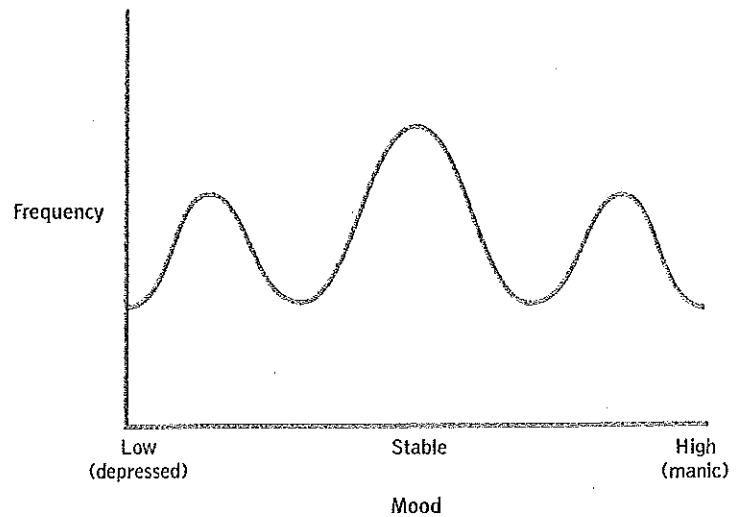
Did you get 3? If you didn't, you might have made a common mistake. The mode is the score that occurs most frequently, not the frequency of that score. So, in the data set above, the score 3 occurs 23 times. The mode is 3, *not* 23.

The mode in this example is particularly easy to determine because there is one most common score. Sometimes a data set has no specific mode. This is especially true when the scores are reported to several decimal places (and no number occurs twice). When there is no specific mode, we sometimes report the most common interval as the mode. Other data sets have more than one specific mode, where two or more different scores are the most common. When there is more than one mode, we report both, or all, of the most common scores. *When a distribution of scores has one mode, we refer to it as unimodal. When a distribution has two modes, we call it*

- The median is the middle score of all the scores in a sample when the scores are arranged in ascending order. If there is no single middle score, the median is the mean of the two middle scores.
- The mode is the most common score of all the scores in a sample.
- A unimodal distribution has one mode, or most common score.

FIGURE 4-4
Bipolar Disorder and
the Modal Mood

Because people with bipolar disorder, especially those who are not receiving treatment, have three different mood states in their lives, it might be hard to determine a true center for their daily mood scores. The distribution might be multimodal, with one mode for depressive days, one for stable days, and one for manic days.



bimodal. When a distribution has more than two modes, we call it *multimodal*. A histogram describing bipolar disease, for example, might be multimodal, as illustrated in Figure 4-4.

As demonstrated in the example above, the mode can be used with scale data; however, it is more commonly used with nominal data. For example, Cancer Research UK (2003) reported that lung cancer was the most common cause of cancer death in the United Kingdom (22%). No other type of cancer came close. Colorectal cancer accounted for 10% of cancer deaths, breast cancer for 8%, and all other types for 7% or less. In this data set, the modal type of cancer death is lung cancer.

How Outliers Affect Measures of Central Tendency

The mean usually appears in journal articles and media reports. However, we use the median and mode when the data are skewed (lopsided). One common reason for skewed data is a statistical outlier. *An outlier is an extreme score that is either very high or very low in comparison with the rest of the scores in the sample.* To demonstrate the effect of outliers on the mean, as well as the median's resistance to the effect of outliers, let's use the statistical archives of America's national pastime, baseball.

Some baseball players have made a career out of their ability to steal bases. But one major league player eclipsed all others in terms of the total number of stolen bases: Rickey Henderson. Reported below are five top base-stealers in major league history and the number of bases stolen:

- A **bimodal** distribution has two modes, or most common scores.
- A **multimodal** distribution has more than two modes, or most common scores.
- An **outlier** is an extreme score that is either very high or very low in comparison with the rest of the scores in a sample.

Rickey Henderson	1406
Lou Brock	938
Billy Hamilton	912
Ty Cobb	892
Tim Lincecum	808

To get a sense of the lifetime achievement of the best base-stealers, we might want to calculate the mean for these five players, using the formula to get a little more practice with the symbols of statistics.

$$M = \frac{\Sigma X}{N} = \frac{\Sigma(1406 + 938 + 912 + 892 + 808)}{5} = \frac{4956}{5} = 991.2$$

As often happens when there is an outlier, this mean is not the same as any of the scores in the sample. The mean of 991.2 is not typical for any of these five baseball players. An important feature of the mean, however, is that it is the point at which all the other scores would balance. Figure 4-5 demonstrates this using the analogy of a balance beam with the numbers of stolen bases from 808 to 1406 marked on it. Weights are placed to represent each of the scores in our sample. The seesaw is perfectly balanced if we put its fulcrum at the mean of 991.2.

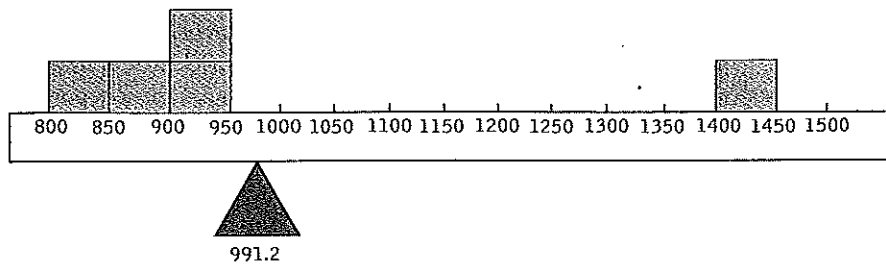


FIGURE 4-5
Outliers and the Mean

When there is an outlier, sometimes the mean is not representative of any one actual score. With the base-stealing data, the mean of 991.2 is above the lowest four scores and well below the highest. Rickey Henderson's score pulls the mean higher, even among the very best base-stealers ever.

When we look at the stolen base data, we notice that Rickey Henderson's score is very different from the others. Four of the scores are between 808 and 938, not a very wide range. But Rickey Henderson stole 1406 bases. When there is an outlier, like Rickey Henderson, it is important to consider what his score would do to the mean, especially if we have a small number of observations.

When we eliminate Rickey Henderson's score, the data are now 808, 892, 912, and 938, and the mean is now:

$$M = \frac{\Sigma X}{N} = \frac{\Sigma(808 + 892 + 912 + 938)}{4} = 887.5$$

The mean of these scores, 887.5, is a good deal lower than the mean of the scores that included Rickey Henderson's very high number of stolen bases. We see from Figure 4-6 that this mean, like the previous mean, marks the point at which all other scores are perfectly balanced around it. However, this mean is a little more representative of the scores—887.5 does seem to be a typical score for these four players.

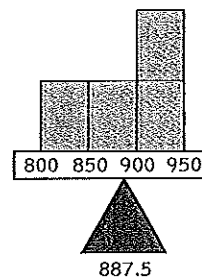


FIGURE 4-6
The Mean Without the Outlier

When the outlier—Rickey Henderson—is omitted from the base-stealing data, the mean is now more representative of the actual scores in the sample.

Which Measure of Central Tendency Is Best?

Different measures of central tendency can lead to very different conclusions. When a decision needs to be made about which measure to use, the choice is usually between the mean and the median. Typically, the mean is the measure of choice. However,

EXAMPLE 4.7



Mario Mignani/Getty Images

Celebrity Outliers Reports of the cost of a typical Manhattan apartment depend on whether the mean or the median is reported. Film star Gwyneth Paltrow and her husband, Coldplay lead singer Chris Martin, sold their Manhattan apartment in the spring of 2005 for around \$7 million. Such a sale would be an outlier and would boost the mean; however, it would not affect the median. Of course, either way, the typical Manhattan apartment is not within the budget of the typical college graduate!

Manhattan, the *New York Times* provided a model of responsible journalism by demonstrating that there is a story behind how central tendency is used to communicate real estate prices. Before the U.S. real estate bubble burst, William Neuman (2005) reported on record-high Manhattan housing prices of \$750,000 (median) and \$1,276,202 (mean). The mean was inflated by a few sales in the millions, outliers that would not affect the median. For example, the film star Gwyneth Paltrow and her husband, Chris Martin of the rock band Coldplay, sold their Manhattan apartment right around that time for about \$7 million. This expensive price certainly would have inflated the mean, but it would not have affected the median.

whenever the distribution is skewed by an outlier (or when the distribution of observations itself is skewed), the median is used to measure central tendency.

The mode is generally used in three situations: (1) when one particular score dominates a distribution; (2) to describe bimodal or multimodal distributions; and (3) when the data are nominal. When you are uncertain as to which measure is the best indicator of central tendency, report all three.

Central tendency communicates an enormous amount of information with a single number, so it is not surprising that measures of central tendency are among the most widely reported of descriptive statistics. Unfortunately, many people use them incorrectly. One particular statistical “lie” or trick that is used on consumers more than any other is reporting the mean instead of the median. To avoid being tricked when you see a report of central tendency, first notice whether it is reporting an average (mean) or a median. Second, if it is reporting a mean, think about whether that distribution is likely to be skewed by one extremely high number (as in the base-stealing example).

Here is another example in which the mean and median would lead to quite different conclusions: In an article on housing prices in

▶ MASTERING THE CONCEPT

4-2: The mean is the most common indicator of central tendency, but it is not always the best. When there is an outlier, it is usually better to use the median.

CHECK YOUR LEARNING

Reviewing the Concepts

- > The central tendency of a distribution is the one number that best describes what is typical in that distribution (often its high point).
- > The three measures of central tendency are the mean (arithmetic average), the median (middle score), and the mode (most frequently occurring score).

- > The mean is the most commonly used measure of central tendency, but the median is preferred when the distribution is skewed (lopsided) or there is an outlier. If you are unsure of which measure to report, then report all three.
- > The symbols used in statistics have very specific meanings; changing a symbol even slightly can change its meaning significantly.

Clarifying the Concepts	4-1 What is the difference between statistics and parameters?
	4-2 Does an outlier have the greatest effect on the mean, median, or mode?
<hr/>	
Calculating the Statistics	4-3 Calculate the mean, median, and mode of the following sets of numbers.
	a. 10, 8, 22, 5, 6, 1, 19, 8, 13, 12, 8
	b. 122.5, 123.8, 121.2, 125.8, 120.2, 123.8, 120.5, 119.8, 126.3, 123.6
	c. 0.100, 0.866, 0.781, 0.555, 0.222, 0.245, 0.234

Applying the Concepts	4-4 Let's examine fictional data for 20 seniors in college. Each score represents the number of nights a student spends socializing in one week: 1, 0, 1, 2, 5, 3, 2, 3, 1, 3, 1, 7, 2, 3, 2, 2, 0, 4, 6
	a. Using the formula, calculate the mean of these scores.
	b. If the researcher reported the mean of these scores to the university as an estimate for the whole university population, what symbol would be used for the mean? Why?
	c. If the researcher were interested only in the scores of these 20 students, what symbol would be used for the mean? Why?
	d. What is the median of these scores?
	e. What is the mode of these scores?
	f. Are the median and mean similar to or different from each other? What does this tell you about the distribution of scores?

Solutions to these Check Your Learning questions can be found in Appendix D.

Measures of Variability

After World War II, people often poked fun at the poor quality of Japanese transistor radios and other products. But it took just three years to transform Japanese manufacturing into an industry that made high-quality products through low variability. In statistics, *variability* is a numerical way of describing how much spread there is in a distribution. The measures of variability we learn about next provide new ways to describe the distribution of our data. One way to numerically describe the variability of a distribution is by computing its *range*. A second and more common way to describe variability is by computing *variance* and its square root, known as *standard deviation*.

Range

The range is the easiest measure of variability to calculate. *The range is a measure of variability calculated by subtracting the lowest score (the minimum) from the highest score (the maximum).* Maximum and minimum are sometimes substituted in this formula to

◀ MASTERING THE CONCEPT

4-3: After central tendency, variability is the second most common concept to help us understand the shape of a distribution.

Common indicators of variability are range, variance, and standard deviation.

- **Variability** is a numerical way of describing how much spread there is in a distribution.
- The **range** is a measure of variability calculated by subtracting the lowest score (the minimum) from the highest score (the maximum).

describe the highest and lowest scores, and some statistical computer programs abbreviate these as *max* and *min*. The range is represented in formula as:

$$\text{range} = X_{\text{highest}} - X_{\text{lowest}}$$

MASTERING THE FORMULA

4-2: The formula for the range is: $\text{range} = X_{\text{highest}} - X_{\text{lowest}}$. We simply subtract the lowest score from the highest score to calculate the range.

Here are the scores for the chemistry professors we discussed earlier in the chapter. Each score represents the number of a professor's students who went on to obtain top academic jobs.

3, 3, 3, 4, 5, 9, 5, 3, 3, 5, 6, 3, 4, 8, 6, 3, 3, 3, 4, 4, 4, 7, 6, 3, 5, 5, 7, 13, 3, 3, 3, 3, 4, 4, 4, 5, 6, 7, 6, 7, 8, 8, 3, 3, 3, 5, 3, 3, 5, 3, 3, 3

EXAMPLE 4.8

We can determine the highest and lowest scores either by reading through the data or, more easily, by glancing at the frequency table for these data.

STEP 1. Determine the highest score. In this case, the highest score is 13.

STEP 2. Determine the lowest score. In this case, the lowest score is 3.

STEP 3. Calculate the range. Subtract the lowest score from the highest score:

$$\text{range} = X_{\text{highest}} - X_{\text{lowest}} = 13 - 3 = 10.$$

The range can be a useful initial measure of variability, but what we learn from the range is limited. It is affected by our highest and lowest scores only. It does not take any other data points into account. The other scores could all be very close to the highest score or all huddled near the center. They could also be spread out evenly or have some other unexpected pattern. We can't know based only on the range.

Variance

Variance is the average of the squared deviations from the mean. It is a concept that we'll soon learn to calculate. Basically, however, variance refers to variability. When something varies, it must vary from, or be different from, some standard. That standard is the mean. So, when we compute variance, the number we arrive at is a number that describes the degree to which a distribution varies with respect to the mean. A small number indicates a small amount of spread or deviation around the mean, and a bigger number indicates a great deal of spread or deviation around the mean.

EXAMPLE 4.9

Students who seek therapy at university counseling centers often do not attend many sessions. For example, in one study, the median number of therapy sessions was 3 and the mean was 4.6 (Hatchett, 2003). Let's examine the spread of fictional scores for a sample of five students: 1, 2, 4, 4, and 10 numbers of therapy sessions, with a mean of 4.2 sessions. Next we find out how far each score deviates from the mean by subtract-

ing the mean from every score. As you might expect, we label the column that lists our scores with an X . Here, our second column includes the results we get when we subtract the mean from each score, or $X - M$. We call each of these a *deviation from the mean* (or just a *deviation*)—the amount that a score in a sample differs from the mean of the sample.

X	$X - M$
1	-3.2
2	-2.2
4	-0.2
4	-0.2
10	5.8

But we can't just take the mean of the deviations. If we do (and if you try this, don't forget the signs—negative and positive), we get 0. In fact, every time we do this with any data set, the mean is 0. Are you surprised? Remember, the mean is the point at which all scores are perfectly balanced. Mathematically, the scores *have* to balance out. Yet we know that there *is* variability among these scores. The number representing the amount of variability is certainly not 0!

When we ask students for ways to eliminate the negative signs, two suggestions typically come up: (1) take the absolute value of the deviations, thus making them all positive, or (2) square all the scores, again making them all positive. It turns out that the latter, squaring all the deviations, is how statisticians solve this problem. Once we square our deviations, we can take their average and get a measure of variability.

To recap:

STEP 1. Subtract the mean from every score.

We call these deviations from the mean, or deviations.

STEP 2. Square every deviation from the mean.

We call these squared deviations.

STEP 3. Sum all of the squared deviations.

This is often called the sum of squared deviations, or the sum of squares for short.

STEP 4. Divide the sum of squares by the total number in the sample (N).

That is, we take the average of the squared deviations.

This number represents the mathematical definition of variance—the average of the squared deviations from the mean.

Let's calculate variance for our therapy session data. We add a third column to contain the squares of each of the deviations, then add all of these numbers up to compute the *sum of squares* (symbolized as SS), the sum of each score's squared deviation from the mean. In this case, the sum of the squared deviations is 48.80, so the average squared deviation is $48.80/5 = 9.76$. Thus, the variance equals 9.76.

- Variance is the average of the squared deviations from the mean.
- A deviation from the mean is the amount that a score in a sample differs from the mean of the sample; also called *deviation*.
- The sum of squares, symbolized as SS , is the sum of the squared deviations from the mean for each score.

X	$X - M$	$(X - M)^2$
1	-3.2	10.24
2	-2.2	4.84
4	-0.2	0.04
4	-0.2	0.04
10	5.8	33.64

MASTERING THE FORMULA

4-3: The formula for variance is: $SD^2 = \frac{\Sigma(X - M)^2}{N}$. To calculate variance, subtract the mean (M) from every score (X) to calculate deviations from the mean; then square these deviations, sum them, and divide by the sample size (N). By summing the squared deviations and dividing by sample size, we are taking their mean.

Now let's put this in equation form, which will make it look more complicated than it is but will continue to acclimate us to symbolic notation. We need a few new symbols at this point, because variance has several different symbols when it's calculated from a sample, including SD^2 , s^2 , and MS . SD^2 and s^2 come from the words *standard deviation squared*. MS comes from the words *mean square* (referring to the average of the squared deviations). We'll use SD^2 at this point, but we will alert you when we switch to other symbols for variance later. When variance is calculated from a population, it typically has just one symbol, σ^2 (pronounced "sigma squared"), and is a parameter. (Remember, Latin letters are used for statistics, which are calculated from samples, and Greek letters are used for parameters, which are calculated from or hypothesized for populations.)

We already know all the symbols needed to calculate variance: X to indicate the individual scores, M to indicate the mean, and N to indicate the sample size.

$$SD^2 = \frac{\Sigma(X - M)^2}{N}$$

As you can see, variance is really just a mean—the mean of squared deviations.

Standard Deviation

EXAMPLE 4.10

Variance is useful, but not as useful as standard deviation. *Standard deviation is the square root of the average of the squared deviations from the mean, and is the typical amount that each score varies, or deviates, from the mean.* Standard deviation is perhaps better known as the square root of variance. The problem with variance—and the reason that we need standard deviation—is that it's not very easy to understand at a glance. Remember, the numbers of therapy sessions for the five students were 1, 2, 4, 4, and 10, with a mean of 4.2. The typical score does not vary from the mean by 9.76. The variance is based on *squared* deviations, not deviations, so it is too big. When we ask our students how to solve this problem, they invariably say "unsquare it," and that's just what we do. We take the square root of variance to come up with a much more useful number, the standard deviation. The square root of 9.76 is 3.12. Now we have a number that "makes sense" to us. We can now say that the typical number of therapy sessions for students in this sample is 4.2 and the typical amount a student varies from that is 3.12.

As you read journal articles, you often will see the mean and standard deviation reported as: ($M = 4.2$, $SD = 3.12$). A glance at our original data (1, 2, 4, 4, 10) tells us that these numbers make sense: 4.2 does seem to be approximately in the center; scores do seem to vary from 4.2 roughly by 3.12. The score of 10 is a bit of an outlier, but not so much that the mean and standard deviation are not somewhat representative of the typical score and typical deviation.

We didn't actually need a formula to get the standard deviation. We just took the square root of the variance. Perhaps you guessed the symbols for standard deviation

- Standard deviation is the square root of the average of the squared deviations from the mean, and is the typical amount that each score varies, or deviates, from the mean.

by just taking the square root of those for variance. With a sample, standard deviation is either SD or s . With a population, standard deviation is σ . Table 4-2 presents this information concisely. We can write the formula showing how standard deviation is calculated from variance:

$$SD = \sqrt{SD^2}$$

We also can write the formula showing how standard deviation is calculated from the original X 's, M , and N :

$$SD = \sqrt{\frac{\sum(X - M)^2}{N}}$$

MASTERING THE FORMULA

4-4: The most basic formula for standard deviation is: $SD = \sqrt{SD^2}$. We simply take the square root of variance.

MASTERING THE FORMULA

4-5: The full formula for standard deviation is: $SD = \sqrt{\frac{\sum(X - M)^2}{N}}$. To calculate standard deviation, subtract the mean from every score to calculate deviations from the mean. Then, square the deviations from the mean. Sum the squared deviations, then divide by the sample size. Finally, take the square root of the mean of the squared deviations.

TABLE 4-2. Variance and Standard Deviation in Symbols

The variance or standard deviation of a sample is an example of a statistic, whereas the variance or standard deviation of a population is an example of a parameter. The symbols we use depend on whether we are referring to the spread of a sample or a population.

Number	Used for ...	Standard Deviation Symbol	Pronounced	Variance Symbol	Pronounced
Statistic	Sample	SD or s	As written	SD^2 , s^2 , or MS	Letters as written; if superscript 2, then followed by "squared" (e.g., "ess squared")
Parameter	Population	σ	"Sigma"	σ^2	"Sigma squared"

CHECK YOUR LEARNING

Reviewing the Concepts

- > The simplest way to measure variability is the range, which is calculated by subtracting the lowest score from the highest score.
- > Variance and standard deviation both measure the degree to which scores in a distribution vary from the mean. The standard deviation is simply the square root of the variance, it represents the typical deviation of a score from the mean.

Clarifying the Concepts

- 4-5 In your own words, what is variability?
- 4-6 Distinguish the range from the standard deviation. What does each tell us about the distribution?

Calculating the Statistics

- 4-7 Calculate the range, variance, and standard deviation for the following data sets (the same ones from the section on central tendency).
 - a. 10, 8, 22, 5, 6, 1, 19, 8, 13, 12, 8
 - b. 122.5, 123.8, 121.2, 125.8, 120.2, 123.8, 120.5, 119.8, 126.3, 123.6
 - c. 0.100, 0.866, 0.781, 0.555, 0.222, 0.245, 0.234

Applying the Concepts

- 4-8 Final exam week is approaching and students are not eating as well as usual. Four students were asked how many calories of junk food they had consumed between

continued on next page

noon and 10:00 P.M. on the day before an exam. The estimated numbers of nutritionless calories, calculated with the help of a nutritional software program, were 450, 670, 1130, and 1460.

- Using the formula, calculate the range for these scores.
- What information can't you glean from the range?
- Using the formula, calculate variance for these scores.
- Using the formula, calculate standard deviation for these scores.
- If a researcher were interested only in these four students, what symbols would she use for variance and standard deviation, respectively?
- If another researcher hoped to generalize from these four students to all students at the university, what symbols would he use for variance and standard deviation?

Solutions to these Check Your Learning questions can be found in Appendix D.



REVIEW OF CONCEPTS

Central Tendency

Three measures of *central tendency* are commonly used in research. The *mean* is the arithmetic average of the data. The *median* is the midpoint of the data set; 50% of scores fall on either side of the median. The *mode* is the most common score in the data set. When there's one mode, the distribution is *unimodal*; when there are two modes, it's *bimodal*; and when there are three or more modes, it's *multimodal*. The mean is highly influenced by *outliers*, whereas the median and mode are resistant to outliers. It is important to consider whether outliers are present in our data set when deciding which measure of central tendency to use. Usually, however, the mean is the preferred measure. When the mean and other measures are used to describe samples, they're called *statistics* and symbolized by Latin letters; when they're used to describe populations, they're called *parameters* and symbolized by Greek letters.

Measures of Variability

The *range* is the simplest measure of *variability* to calculate. It is often used when the preferred measure of central tendency is the median. It is calculated by subtracting the minimum score in our data set from the maximum score. Variance and standard deviation are much more common measures of variability. They are used when the preferred measure of central tendency is the mean. *Variance* is the average of the squared deviations. It is calculated by subtracting the mean from every score to get *deviations from the mean*, then squaring each of the deviations and taking their mean. (The *sum of squares* of the deviations is used in many inferential statistics.) *Standard deviation* is the square root of variance. It is the typical amount that a score deviates from the mean.

SPSS®

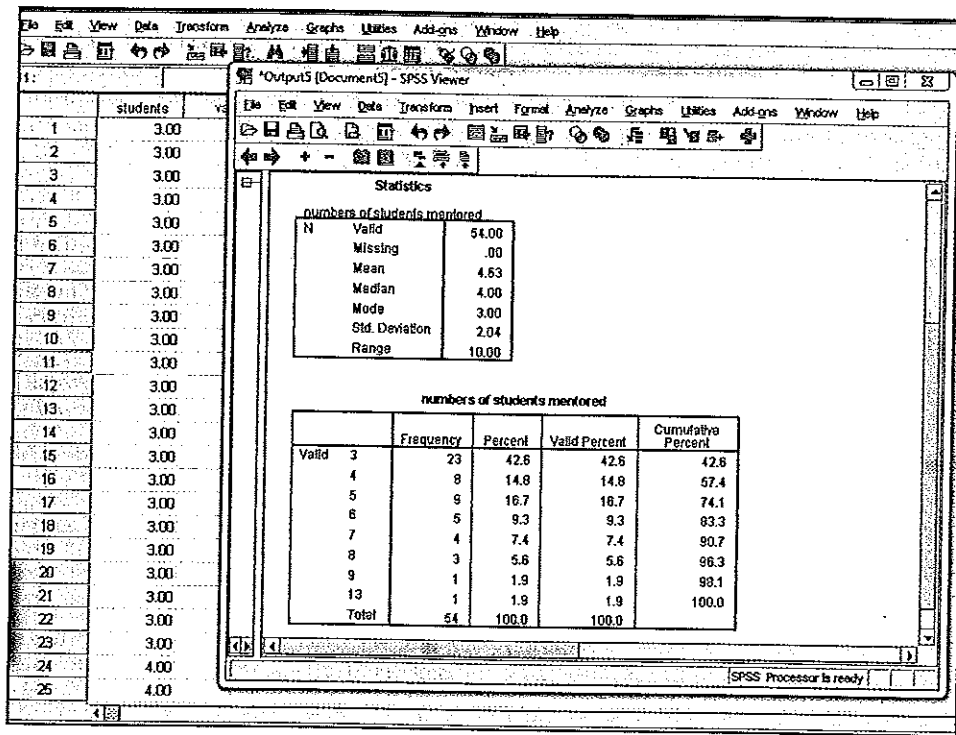
The left-hand column in "Data View" is prenumbered, beginning with 1. Each column to the right of that number contains data for a particular variable; each row below that number represents a unique individual. Notice the choices at the top of the "Data View" screen. Enter the data for the numbers of

students mentored by each professor that were used several times in this chapter, including in example 4.5.

To get a numerical description of a variable, select **Analyze** → **Descriptive Statistic** → **Frequencies**. Then, select the variable "students" by highlighting it and then clicking the **ar-**

row to move it from the left side to the right side. Then click: Statistics → Mean, Median, Mode, Std. deviation, Range → Continue → OK.

Your output will look like that in the screen shot shown here.



How It Works

4.1 CALCULATING THE MEAN

Here are data for the numbers of nights out socializing in a week for 20 students.

1, 2, 7, 6, 1, 2, 6, 5, 4, 4, 0, 3, 2, 2, 3, 4, 3, 5, 4, 4

How can we calculate the mean? First, we add up all of the scores:

$$1 + 2 + 7 + 6 + 1 + 2 + 6 + 5 + 4 + 4 + 0 + 3 + 2 + 2 + 3 + 4 + 3 + 5 + 4 + 4 = 68$$

Then we divide by 20, the number of scores:

$$68/20 = 3.4$$

With the formula $M = \frac{\sum X}{N}$, we'd calculate:

$$M = \frac{(1+2+7+6+1+2+6+5+4+4+0+3+2+2+3+4+3+5+4+4)}{20}$$

4.2 CALCULATING THE MEDIAN

Using the data for "nights out socializing," how can we calculate the median? The median is simply the middle score, or the average of the two middle scores. For these data, we first arrange the data from lowest score to highest score:

0 1 1 2 2 2 2 3 3 3 4 4 4 4 4 5 5 6 6 7

We divide the number of scores, 20, by 2 and add 1/2 to get 10.5. The mean of the 10th and 11th scores in the ordered list is the median— $(3 + 4)/2 = 3.5$. The median is 3.5.

4.3 CALCULATING THE MODE

How can we calculate the mode for the “nights out socializing” data? The mode is the most common score. We can determine this for these data by looking at the frequency distribution. Five people have a score of 4. The mode is 4.

4.4 CALCULATING VARIANCE

How can we calculate variance for the “nights out socializing” data? To calculate variance for these data, we first subtract the mean, 3.4, from every score to calculate a deviation. We then square the deviations. These calculations are shown in the table below.

X	$(X - M)$	$(X - M)^2$
1	-2.4	5.76
2	-1.4	1.96
7	3.6	12.96
6	2.6	6.76
1	-2.4	5.76
2	-1.4	1.96
6	2.6	6.76
5	1.6	2.56
4	0.6	0.36
4	0.6	0.36
0	-3.4	11.56
3	-0.4	0.16
2	-1.4	1.96
2	-1.4	1.96
3	-0.4	0.16
4	0.6	0.36
3	-0.4	0.16
5	1.6	2.56
4	0.6	0.36
4	0.6	0.36

We then add all of the scores in the third column to get the sum of squared deviations, or the sum of squares. This sum is 64.8.

We can use the formula to complete our calculations:

$$SD^2 = \frac{\sum(X - M)^2}{N} = \frac{64.8}{20} = 3.24$$

The variance is 3.24.

4.5 CALCULATING STANDARD DEVIATION

How can we calculate standard deviation for the “nights out socializing” data? The standard deviation is the typical amount that the scores in a sample vary, or deviate, from the mean, and is the square root of the variance. For these data, we can calculate standard deviation directly from the variance we calculated above using this formula:

$$SD = \sqrt{SD^2} = \sqrt{3.24} = 1.80$$

The standard deviation is 1.80.

Exercises

Clarifying the Concepts

- 4.1 Define the three measures of central tendency: mean, median, and mode.
- 4.2 The mean can be assessed visually and arithmetically. Describe each method.
- 4.3 Explain how the mean mathematically balances the distribution.
- 4.4 Explain what is meant by unimodal, bimodal, and multimodal distributions.
- 4.5 In what situations is the mode typically used?

- 4.6 What is an outlier?
- 4.7 Are the mean and median affected by outliers?
- 4.8 Define the symbols used in the equation for variance:

$$SD^2 = \frac{\sum(X - M)^2}{N}$$

- 4.9 Why is the standard deviation typically reported rather than the variance?
- 4.10 Find the incorrectly used symbol or symbols in each of the following statements or formulas. For each statement or formula, (i) state which symbol(s) is/are used incorrectly, (ii) explain why the symbol(s) in the original statement is/are incorrect, and (iii) state what symbol(s) *should* be used.

- a. The mean and standard deviation of the sample of reaction times were calculated ($m = 54.2$, $SD^2 = 9.87$).
- b. The mean of the sample of high school student GPAs was $\mu = 3.08$.
- c. Range = $X_{highest} - X_{lowest}$

Calculating the Statistics

- 4.11 Calculate the mean, median, and mode for the following data: 15, 34, 32, 46, 22, 36, 34, 28, 52, 28.
- 4.12 Calculate the mean, median, and mode for the following salaries: \$44,751, \$52,000, \$41,500, \$38,862, \$51,380, \$61,774.
- 4.13 Add another data point, 112, to the data presented in Exercise 4.11. Calculate the mean, median, and mode again. How does this new data point affect your calculations?
- 4.14 Add another salary, \$97,582, to the data presented in Exercise 4.12. Calculate the mean, median, and mode again. How does this new salary affect your calculations?
- 4.15 Calculate the range, variance, and standard deviation for the data in Exercise 4.11.
- 4.16 Calculate the range, variance, and standard deviation for the salaries in Exercise 4.12.
- 4.17 How does the range change when you include the outlier salary, \$97,582, with the data from Exercise 4.12?
- 4.18 Here are the *U.S. News & World Report* data again on alumni giving at the top 70 national universities.

48 61 45 39 46 37 38 34 33 47
 29 38 38 34 29 29 36 48 27 25
 15 25 14 26 33 16 33 32 25 34
 26 32 11 15 25 9 25 40 12 20
 32 10 24 9 16 21 12 14 18 20
 18 25 18 20 23 9 16 17 19 15
 14 18 16 17 20 24 25 11 16 13

- a. Calculate the mean of these data, showing that you know how to use the symbols and formula.
- b. Determine the median of these data.
- 4.19 Describe the variability in the data presented in Exercise 4.18 by computing the range.
- 4.20 The Mount Washington Observatory (MWO) in New Hampshire claims to have the world's worst weather. Below are some data on the weather extremes recorded at the MWO. Calculate the mean and median normal daily minimum temperature across the year.

	Normal Daily Maximum (°F)	Normal Daily Minimum (°F)	Record Low in °F (Year)	Peak Wind Gust in Miles per Hour (Year)
January	14.0	-3.7	-47 (1934)	173 (1985)
February	14.8	-1.7	-46 (1943)	166 (1972)
March	21.3	5.9	-38 (1950)	180 (1942)
April	29.4	16.4	-20 (1995)	231 (1934)
May	41.6	29.5	-2 (1966)	164 (1945)
June	50.3	38.5	8 (1945)	136 (1949)
July	54.1	43.3	24 (2001)	154 (1996)
August	53.0	42.1	20 (1986)	142 (1954)
September	46.1	34.6	9 (1992)	174 (1979)
October	36.4	24.0	-5 (1939)	161 (1943)
November	27.6	13.6	-20 (1958)	163 (1983)
December	18.5	1.7	-46 (1933)	178 (1980)

- 4.21 Calculate the mean, median, and mode for the record low temperatures recorded on top of Mount Washington presented in Exercise 4.20.
- 4.22 Calculate the mean, median, and mode for the peak wind gust data presented in Exercise 4.20.
- 4.23 When no mode appears in the raw data, we can compute a mode by breaking the data into intervals. How might you do this for the peak wind gust data presented in Exercise 4.20?
- 4.24 Calculate the range, variance, and standard deviation for normal daily minimum temperature across the years presented in Exercise 4.20.
- 4.25 Calculate the range, variance, and standard deviation for the record low temperatures recorded on top of Mount Washington presented in Exercise 4.20.
- 4.26 Calculate the range, variance, and standard deviation for the peak wind gust data presented in Exercise 4.20.

Applying the Concepts

- 4.27 For the data presented in Exercise 4.20, the "normal" daily maximum and minimum temperatures recorded

at the Mount Washington Observatory are presented for each month. These are likely to be measures of central tendency for each month over time. Explain why these “normal” temperatures might be calculated as means or medians. What would be the reasoning for using one statistic over the other?

- 4.28 Back in Exercises 4.13 and 4.14 we saw how the mean and median changed when an outlier was included in the computations. If you were reporting the typical salary at a company, how might the mean and median give different impressions to potential applicants?
- 4.29 The “normal” weather data from the Mount Washington Observatory are broken down by months. Why might you not want to average across all months in a year? How else could you summarize the year?
- 4.30 There appears to be an outlier in the data for peak wind gust recorded on top of Mount Washington (see data in Exercise 4.20). Where do you see an outlier and how would excluding this data point affect the different calculations of central tendency?
- 4.31 Here are winning percentages for 11 players for their best four-year pitching performances: 0.755, 0.721, 0.708, 0.773, 0.782, 0.747, 0.477, 0.817, 0.617, 0.650, 0.651.
- What is the mean of these scores?
 - What is the median of these scores?
 - Compare the mean and median. Does the difference between them suggest that the data are skewed very much?
- 4.32 Briefly describe a real-life situation in which the median is preferable to the mean. Give hypothetical numbers for mean and median in your explanation. Be original! (Don't use home prices or another example from the chapter.)
- 4.33 Find an advertisement for a weight-loss product either online or in the print media—the more unbelievable the claims, the better!
- What does the ad promise that this product will do for the consumer?
 - What data does it offer for its promised benefits? Does it offer any descriptive statistics or merely testimonials? If it offers descriptive statistics, what are the limitations of what they report?
 - If you were considering this product, what measures of central tendency would you most like to see? Explain your answer, noting why not all measures of central tendency would be helpful.
 - If a friend with no statistical background were considering this product, what would you tell him or her?
- 4.34 When you see an ad on TV for a body-shaping product (e.g., an abdominal muscle machine), often a

person with a wonderful success story is featured in the ad. The statement “individual results may vary” hints at what kind of data the advertisement may be presenting.

- What kind of data is being presented in these ads?
 - What statistics could be presented to help inform the public about how much “individual results might vary”?
- 4.35 The National Survey of Student Engagement asked students how often they asked questions in class or participated in classroom discussions. The options were “never,” “sometimes,” “often,” and “very often.” Here are the percentages, reported in 2005, of students who responded “very often” for the 31 institutions classified as liberal arts colleges that allowed their 2004 data to become public through the *U.S. News & World Report* Web site.

58	45	53	45	65	41	50	46	54
59	52	60	59	62	54	52	53	54
83	60	32	62	50	50	43	32	53
60	52	55	53					

- What is the range of these data?
 - The top college is Marlboro College in Vermont, and the two tied for lowest are Randolph-Macon Women's College in Virginia and Texas A&M University in Galveston. What research questions do these data suggest to you? State at least one research question generated by these data.
- 4.36 Here again are the data from the National Survey of Student Engagement for a sample of 19 national universities, as reported in 2005. These are the percentages of students who said they were assigned between 5 and 10 20-page papers.

0	5	3	3	1	10	2
2	3	1	2	4	2	1
1	1	4	3	5		

- Calculate the mean of these data using the symbols and formula.
 - Calculate the variance of these data using the symbols and formula, but also using columns to show all calculations.
 - Calculate the standard deviation using the symbols and formula.
 - In your own words, describe what the mean and standard deviation of these data tell us about these scores.
- 4.37 For each of the following situations, state whether the mean would be a statistic or a parameter. Explain your answer.

- ered in
vary”
may be
- ads?
inform
results
- asked
for par-
were
- Here
who
ified
ata to
report
- ont,
con
uni-
do
rich
of
an-
ges
and
- 4.37
- According to 1991 Canadian census data, the mean income (from employment only) of French-speaking Canadians living in Ontario was \$29,527, higher than the general population mean of \$28,838.
 - In the 2004–2005 National Basketball Association season, the 30 teams won a mean 41.00 games.
 - The General Social Survey (GSS) includes a vocabulary test in which participants are given a series of words and asked to choose the appropriate synonym from a multiple-choice list of five words (e.g., *beast* with the choices *afraid*, *words*, *large*, *animal*, and *separate*). The mean vocabulary test score was 5.98.
 - The National Survey of Student Engagement (NSSE) asked students at participating institutions how often they discussed ideas or readings with their professors outside of class. Among the 19 national universities that made their data public, the mean percentage of students who responded “very often” was 8%.
- 4.38 Consider the many possible distributions of grades on a quiz in a statistics class; imagine that the grades could range from 0 to 100. For each of the following situations, give a hypothetical mean and median (that is, make up a mean and a median that might occur with a distribution that has this shape). Explain your answer.
- Normal distribution
 - Positively skewed distribution
 - Negatively skewed distribution
- 4.39 For each of the following distributions, state whether it’s more likely to be unimodal or bimodal. Explain your answer.
- Age of patients in a hospital maternity ward
 - Depression scores on a Beck Depression Inventory
 - GRE scores of applicants to sociology graduate programs
 - The cost of an AIDS drug that is sold in developed countries in Europe, as well as in developing countries in Africa
- 4.40 Here are the numbers of wins for the 30 National Basketball Association teams in the 2004–2005 season.
- 45 43 42 33 33 54 47 44 42 30
59 45 36 18 13 52 49 44 27 26
62 50 37 34 34 59 58 51 45 18
- Determine the mean, median, and mode of these data. Use symbols and the formula when showing your calculation of the mean.
 - Using software, calculate the range and standard deviation of these data.
 - Write a one- to two-paragraph summary describing the distribution of these data. Mention center, variability, and shape. Be sure to discuss the number of modes (i.e., unimodal, bimodal, multimodal), any possible outliers, and the presence and direction of any skew.
- 4.41 The U.S. Census Bureau collects and analyzes data on numerous aspects of American life by state, including the percentage of people with high school degrees, bachelor’s degrees, and advanced degrees. If you wanted to calculate the “average” percentage of people with advanced degrees across all states, would you report a mean, median, or mode? Explain your answer clearly.
- 4.42 According to a 2007 article on the Economist.com Web site, Americans are the international leaders in TV viewing, averaging 8 hours and 11 minutes a day. Below are approximate, daily average viewing times for 12 countries based on this source:
- United States—8.2 hours
Turkey—5 hours
Italy—4.05 hours
Japan—3.75 hours
Spain—3.6 hours
Portugal—3.5 hours
Australia—3.2 hours
South Korea—3.16 hours
Canada—3.1 hours
Britain—3 hours
Denmark—3 hours
Finland—2.8 hours
- Compute the mean and median across these 12 data points.
 - How are these statistics affected by including or excluding the United States?
- 4.43 Refer to the data from Exercise 4.42.
- How do you think these daily “averages” were calculated—using means or medians?
 - Do you think TV viewing habits might vary by other personal or demographic characteristics? Could these represent confounds?
 - How might you collect samples to more specifically describe TV viewing habits as a function of other personal characteristics?
- 4.44 When the typical height or typical weight of children is plotted to create growth charts, do you think it would be appropriate to use the mean for these data? There are often outliers for height, but why might we not have to be concerned with their effect on these data?
- 4.45 Guinness World Records relies on what kind of data for its amazing claims?

Terms

central tendency (p. 72)

mean (p. 73)

statistic (p. 74)

parameter (p. 74)

median (p. 76)

mode (p. 77)

unimodal (p. 77)

bimodal (p. 78)

multimodal (p. 78)

outlier (p. 78)

variability (p. 81)

range (p. 81)

variance (p. 82)

deviation from the mean (p. 83)

sum of squares (p. 83)

standard deviation (p. 84)

Symbols

M (p. 74)

\bar{X} (p. 74)

μ (p. 74)

X (p. 75)

N (p. 75)

mdn (p. 76)

SS (p. 83)

SD^2 (p. 84)

s^2 (p. 84)

MS (p. 84)

σ^2 (p. 84)

SD (p. 85)

s (p. 85)

σ (p. 85)

Formulas

$$M = \frac{\Sigma X}{N} \quad (\text{p. 75})$$

$$\text{range} = X_{\text{highest}} - X_{\text{lowest}} \quad (\text{p. 82})$$

$$SD^2 = \frac{\Sigma(X - M)^2}{N} \quad (\text{p. 84})$$

$$SD = \sqrt{SD^2} \quad (\text{p. 85})$$

$$SD = \sqrt{\frac{\Sigma(X - M)^2}{N}} \quad (\text{p. 85})$$