

**Objectives**

- ★ Calculate the mean, median, and mode
- ★ Calculate the moving average for a data set
- ★ Calculate the weighted mean
- ★ Define a measure of central tendency
- ★ Demonstrate understanding of the concept of mean, median, and mode
- ★ Demonstrate understanding of the concept of sample means
- ★ Determine the most appropriate measure of center
- ★ Estimate the average of a frequency distribution
- ★ Find a data value using the mean

Sep 26-8:38 AM

**Measures of Location**

Frequency distributions, bar charts, pie charts, and histograms can be informative visual tools for examining the big picture when analyzing data. But there is a lack of exactness in the language that we use to describe these graphs. Suppose we say that one data set is more compact than another. This only leads to the question, How much more compact is it? Graphical analysis is ill-equipped to answer that question precisely.

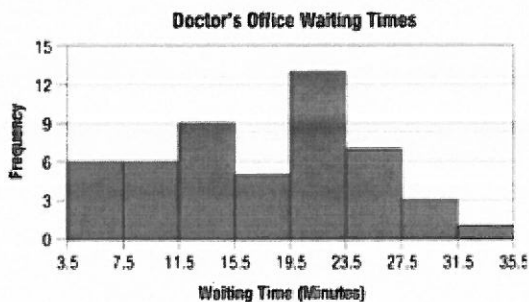


Figure 4.1

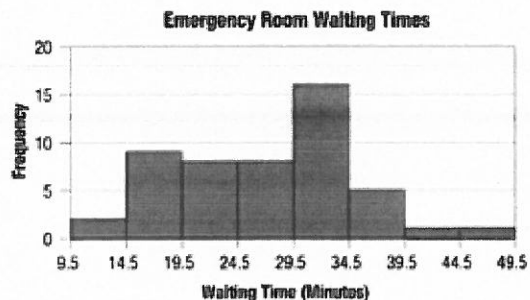
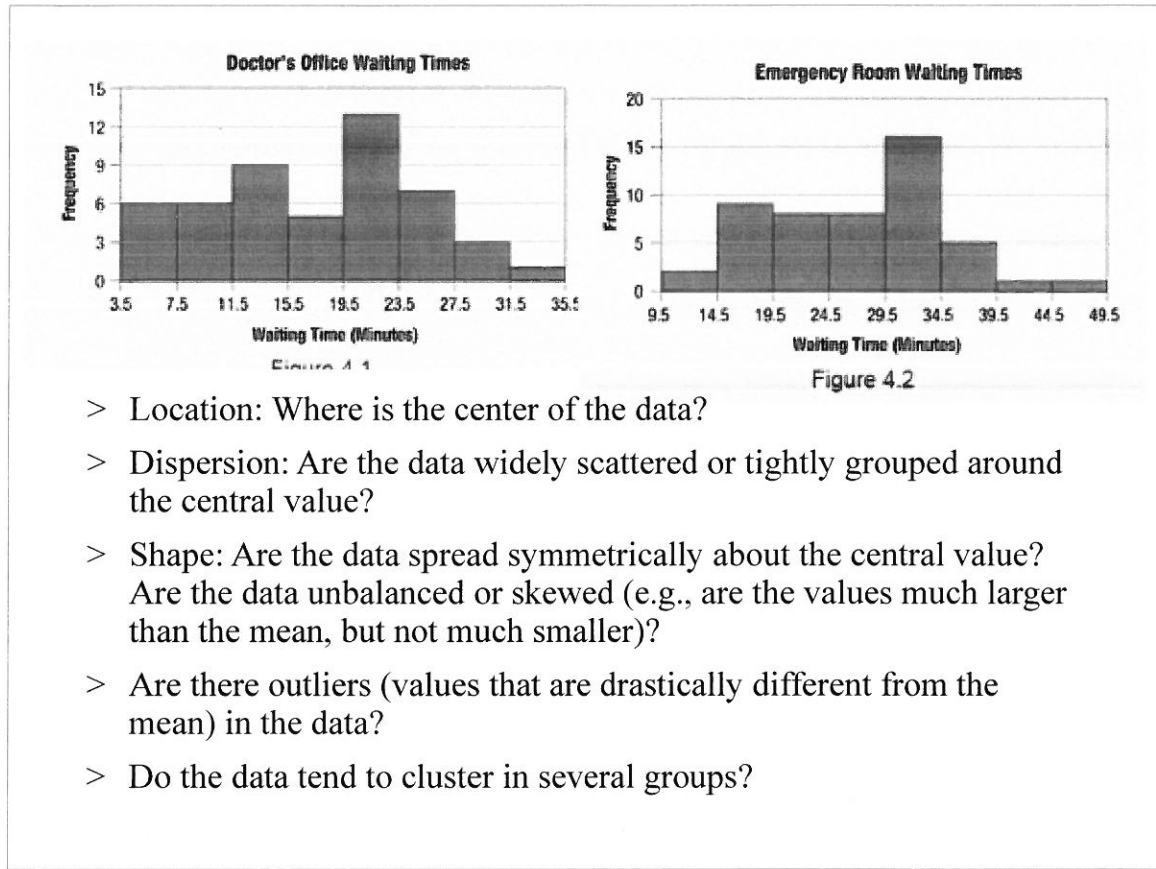


Figure 4.2

Sep 26-8:39 AM



Sep 26-8:40 AM

Different types of measurements will be developed for each of the attributes listed. For example, the mean (average) and the median measure the location or central tendency of a data set. Both of these statistical measures are trying to tell us something about the location of the middle of the data, but they use different ideas for defining that notion of middle.

From the morning paper to the evening news, general concepts are translated into specific statistical measures. Some examples are listed below.

- > The proportion of United States residents earning below the poverty level was 14.3% in 2009.

Source: U.S. Census Bureau

- > The median price of new homes sold in 2009 was \$216,700.

Source: U.S. Census Bureau

- > The average hourly manufacturing wage in the U.S. in 2010 was \$18.90.

Source: U.S. Department of Labor

Sep 26-8:41 AM

Such measures are examples of **numerical descriptive statistics**.

**Definition**

---

**Numerical descriptive statistics** are numerical summaries of data.

There is a distinction between measures that are applied to populations and measures that are applied to samples.

**Definition**

---

Measures that apply to population data are called **parameters**.

**Definition**

---

Measures that apply to sample data are called **statistics**.

Sep 26-8:41 AM

In most instances a data analyst will not know what the population parameters are, since the cost and/or feasibility of obtaining all the population data is usually prohibitive. Inferential statistics involves making conclusions regarding a population using data from a sample.

**Definition**

---

**Inferential statistics** is concerned with making conclusions about population parameters using sample statistics.

There are many formulas in this chapter that you will have to spend some time examining to appreciate. In most cases, the concepts which motivate these formulas are simple. Yet if these concepts are either ignored or forgotten, statistics becomes a meaningless assortment of symbols instead of a useful problem-solving and decision-making tool.

Sep 26-8:41 AM

The data in Table 4.1 are costs (rounded to the nearest dollar) of fill-ups at a local gasoline station in December 2010 for 400 transactions. Looking at the 400 observations without using a graphical representation would be confusing. (It seems 400 measurements would contain a great deal of information, yet the sheer volume of the data obscures comprehension.) It is the old problem of not being able to see the forest for the trees.

22	39	25	35	43	36	52	44	37	49
32	32	35	44	33	51	44	28	29	43
37	45	44	28	51	38	47	34	52	50
33	45	28	54	38	40	37	42	38	44
30	40	36	38	32	36	28	23	51	28
43	43	25	24	47	36	22	24	42	55
39	45	36	36	38	32	43	32	34	54
45	58	49	23	55	43	29	11	34	45
48	36	22	34	18	27	37	50	39	30
37	48	36	64	62	33	33	58	31	47
36	38	33	41	26	43	28	42	43	30
50	32	47	22	42	35	21	56	46	21
50	23	52	46	31	41	35	33	30	47
54	43	43	44	51	49	31	56	27	49

Sep 26-8:42 AM

Statistically speaking, the idea of location is similar to knowing the whereabouts of a person. If we think of a data set as a group of data values that cluster around some central value, then this central value provides a focal point for the data set—a location of sorts. Unfortunately, the notion of central value is a vague concept, which is as much defined by the way it is measured as by the notion itself. There are several statistical measures that can be used to define the notion of center: the arithmetic mean, weighted mean, trimmed mean, median, and mode.

Sep 26-8:42 AM

### Arithmetic Mean

The arithmetic mean is one of the more commonly used statistical measures. It appears every day in newspapers, business publications, and frequently in conversation. For example, when your instructor returns an assessment, after viewing your grade, one of the first questions asked is, What is the average? The word average is often associated with the mean.

#### Formula: Arithmetic Mean

Suppose there are  $n$  observations in a data set, consisting of the observations  $x_1, x_2, \dots, x_n$ ; then the **arithmetic mean** is defined to be

$$\frac{1}{n}(x_1 + x_2 + \dots + x_n).$$

Sep 26-8:43 AM

If we use some common mathematical notation (summation notation, represented by  $\sum$ ), the formula can be simplified to

$$\frac{\sum x_i}{n},$$

where  $x_i$  is the  $i^{\text{th}}$  data value in the data set and  $\sum$  (pronounced *sigma*) is a mathematical notation for adding values. There are two symbols that are associated with the expression given above:

$$\mu = \frac{1}{N}(x_1 + x_2 + \dots + x_N) \text{ the population mean, and}$$

$$\bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n) \text{ the sample mean.}$$

Here  $N$  refers to the size of the population and  $n$  refers to the size of the sample. Otherwise, the calculations are made in precisely the same way. The Greek letter  $\mu$ , representing the population mean, is pronounced *mu* and the symbol  $\bar{x}$  representing the sample mean, is pronounced *x-bar*.

Sep 26-8:44 AM

Population  
(Parameters)

Samples  
(Statistics)

$$\mu = \frac{1}{N}(x_1 + x_2 + \dots + x_N) \text{ the population mean} \quad \bar{x} = \frac{1}{n}(x_1 + x_2 + \dots + x_n) \text{ the sample mean.}$$

Here  $N$  refers to the size of the population and  $n$  refers to the size of the sample. Otherwise, the calculations are made in precisely the same way. The Greek letter  $\mu$ , representing the population mean, is pronounced *mu* and the symbol  $\bar{x}$  representing the sample mean, is pronounced *x-bar*.

Sep 26-8:44 AM

#### Example 4.1

Calculate the sample mean of the following sample data values: 4, 10, 7, 15.

Solution

Note that  $x_1=4$ ,  $x_2=10$ ,  $x_3=7$ ,  $x_4=15$ , and  $n=4$ .

$$\begin{aligned} \bar{x} &= \frac{1}{4}(4+10+7+15)=9 \\ &= \frac{\sum x_i}{n} = \frac{4+10+7+15}{4} = \frac{36}{4}=9 \end{aligned}$$

Sep 26-8:46 AM

The sample mean is 9. But why does adding up a group of numbers and dividing by the number of observations measure central tendency? As unlikely as it sounds, the answer is related to balancing a scale.

**Definition**

Given some point  $A$  and a data point  $x$ , then  $x - A$  represents how far  $x$  deviates from  $A$ . This difference is also called a **deviation**.

Sep 26-8:47 AM

Let's calculate the deviations from the mean for the data in Example 4.1. Examining the deviations from the mean in Table 4.2, we can see the deviations on the left side (-5 and -2) and right side (1 and 6) are in balance. In fact, the mean is considered a point of centrality because the deviations from the mean on the positive side and the negative side are equal (See Figure 4.4). The sample mean can be interpreted as a center of gravity.

Table 4.2 – Deviations from the Mean	
Data ( $x_i$ )	Deviations from the Mean ( $x_i - 9$ )
4	-5
10	1
7	-2
15	6
Total	$\sum (x_i - 9) = 0$



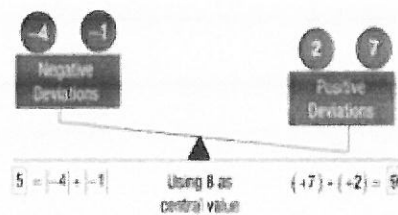
Figure 4.4

Sep 26-8:48 AM

The positive deviations (2 and 7) are not counterbalanced by the negative deviations (-4 and -1). A desirable characteristic of a central value would be to have the positive and negative deviations equal to each other in absolute value.

**Table 4.3 – Deviations from Some Other Value**

Data ( $x_i$ )	Deviations from 8 ( $x_i - 8$ )
4	-4
10	2
7	-1
15	7
Total	$\sum (x_i - 8) = 4$



Sep 26-8:48 AM

Although the arithmetic mean is frequently used, there are times when it should not be employed. Since the mean requires that the data values be added, it should only be used for quantitative data. Furthermore, if one of the data values is extremely large or small relative to others, this could be considered an **outlier**. An outlier is a data value that can have a dramatic impact on the value of the mean.

**Definition**

Statistical measures which are not affected by outliers are said to be **resistant**.

The arithmetic mean is not a **resistant measure**.

Sep 26-8:48 AM



**Weighted Mean**

The weighted mean is similar to the arithmetic mean except it allows you to give different weights (or importance) to each data value. The weighted mean gives you the flexibility to assign weights when you find it inappropriate to treat each observation the same. The weights are usually positive numbers that sum to one, with the largest weight being applied to the observation with the greatest importance. The weights can be determined in a variety of ways, such as the number of employees, market value of a company, or some other objective or subjective method. There are occasions in which it is easier to assign the weights without worrying that they will sum to one. If you are concerned about your weights summing to one, you can make your weights sum to one by dividing each weight by the sum of all the weights.

**Formula: Weighted Mean**

---

The weighted mean of a data set with values  $x_1, x_2, x_3, \dots, x_n$  is given by

$$\bar{x} = \frac{w_1x_1 + w_2x_2 + \dots + w_nx_n}{w_1 + w_2 + \dots + w_n} = \frac{\sum(w_ix_i)}{\sum w_i}$$

where  $w_i$  is the weight of observation  $x_i$ .

Sep 26-8:49 AM

**Example 4.2**

The following table consists of the November 2010 unemployment rates and civilian labor force sizes for the Mid-Atlantic states.

Table 4.4 – Unemployment Rates		
State	Size of Civilian Labor Force (Thousands)	Unemployment Rate (%)
Delaware	422.4	8.4
Maryland	2972.7	7.4
New Jersey	4490.6	9.2
New York	9658.4	8.3
Pennsylvania	6362.0	8.6
Virginia	4177.7	6.8
District of Columbia	331.6	9.8
West Virginia	777.6	9.3

Source: U.S. Department of Labor, Bureau of Labor Statistics

Using the weighted mean, calculate the average unemployment rate for the Mid-Atlantic states.

Sep 26-8:49 AM

Solution

The average unemployment rate is calculated as follows.

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

422.4	8.4
2972.7	7.4
4490.6	9.2
9658.4	8.3
6362.0	8.6
4177.7	6.8
331.6	9.8
777.6	9.3

$$\begin{aligned} \sum(w_i x_i) &= 422.4(8.4) + 2972.7(7.4) + 4490.6(9.2) + 9658.4(8.3) \\ &\quad + 6362.0(8.6) + 4177.7(6.8) + 331.6(9.8) + 777.6(9.3) \\ &= 240627.3 \end{aligned}$$

$$\begin{aligned} \sum w_i &= 422.4 + 2972.7 + 4490.6 + 9658.4 + 6362.0 + 4177.7 + 331.6 + 777.6 \\ &= 29193 \end{aligned}$$

$$\text{so, } \bar{x} = \frac{\sum(w_i x_i)}{\sum w_i} = \frac{240627.3}{29193} \approx 8.24\%$$

Thus, the average unemployment rate, calculated by the weighted mean, is 8.24%. It is appropriate to use the weighted mean to calculate the average unemployment rate since the size of the civilian labor force (the weight) is different for each state.

Sep 26-8:50 AM

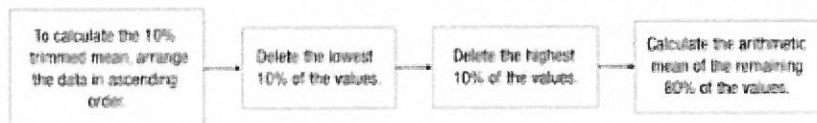
### Trimmed Mean

Since outliers can have an enormous effect on the value of the mean, the mean's usefulness as a typical measure of data is diminished if the data contain outliers.

#### Definition

The **trimmed mean** is a modification of the arithmetic mean which ignores an equal percentage of the highest and lowest data values in calculating the mean.

#### Finding the 10% Trimmed Mean



Sep 26-8:51 AM

**Example 4.3**

Consider the following data.

15, 21, 25, 31, 35, 42, 48, 51, 54, 60

Find the 10% trimmed mean.

**Solution**

Since there are 10 observations, removing the highest 10% and the lowest 10% means removing only one observation from each end of the data.

That is,

$$10\% \text{ of } 10 = 0.1 \cdot 10 = 1.$$

Note that the data are already sorted. If the mean is calculated without including the values of 15 and 60, the resulting measure is called the 10% trimmed mean.

~~15~~, 21, 25, 31, 35, 42, 48, 51, 54, ~~60~~

$$10\% \text{ trimmed mean} = \frac{21 + 25 + 31 + 35 + 42 + 48 + 51 + 54}{8} = \frac{307}{8} = 38.375.$$

If there had been 100 observations, the largest 10% and the smallest 10% (a total of 20 data values) would have been removed before the mean was calculated.

Sep 26-8:51 AM

## Measuring Figure Skating Performances

Almost every figure skating competition has some scoring controversy. The Winter Olympics of 2002 were no exception. French judge, Marie-Reine Le Gongue, said she was "pressured to vote a certain way" when she scored the Russian couple, Elena Berezhnaya and Anton Sikharulidze, over the Canadian pair, Jamie Sale and David Pelletier. In addition, very few people understood exactly how Sarah Hughes won the gold medal and how Michelle Kwan dropped to third after leading the event.

For almost a century figure skating has used a scoring method that is similar to the methodology of the trimmed mean in order to remove bias. Skaters are scored on a 0 to 6 scale. The highest and lowest scores are discarded (the data is trimmed) and the resulting score is computed.

The intent of trimming the data is to avoid bias caused by judges with nationalistic or political agendas.

Since the controversy, the International Skating Union has replaced this scoring method with a new system which, though it is different, still utilizes trimmed data to eliminate bias.

Sep 26-8:52 AM

### Median

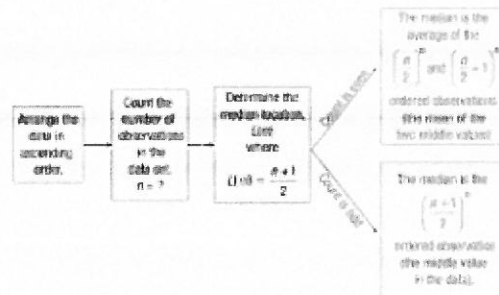
The median of a set of data provides another measure of center. It is a simple idea. To find the median, place the data in ascending order and then find the observation that has an equal number of data values on either side. That is, half of the observations are less than the median and half of the observations are greater than the median. The median is the middle value.

**Definition**

---

The median of a set of observations is the data value in the middle of an ordered array. The same number of data values is on either side of the median value.

#### Finding the Median



Sep 26-8:53 AM

#### Example 4.4

Given the following eleven observations, find the median.

2, 4, 0, 3, 0, 1, 8, 5, 1, 5, 9

Solution

First, the data set must be ordered.

0, 0, 1, 1, 2, 3, 4, 5, 5, 8, 9

The number of observations,  $n=11$ .

Next, calculate the median location,  $L(m) = \frac{n+1}{2} = \frac{11+1}{2} = 6$ . Therefore, the median is the 6<sup>th</sup> ordered observation, which is 3.

0, 0, 1, 1, 2, 3, 4, 5, 5, 8, 9  
 5 data values                      5 data values

Sep 26-8:54 AM

**Example 4.5**

Consider the following ten test scores.

65, 98, 76, 83, 94, 79, 88, 72, 90, 85

Find the median.

**Solution**

The number of observations,  $n=10$ .

The median location,  $L(n) = \frac{10+1}{2} = \frac{11}{2} = 5.5$ .

If there is an even number of observations, average the two center values in the ordered array. The median is the average of the  $\frac{10}{2} = 5^{\text{th}}$  and  $6^{\text{th}}$  ordered observations. Thus, we find the median as follows.

65, 72, 76, 79, 83, 85, 88, 90, 94, 98

4 data values                      4 data values

$$\frac{83+85}{2} = 84 \text{ (the median).}$$

Sep 26-8:55 AM

The median possesses a rather obvious notion of centrality, since it is defined as the central value in an ordered list. It is not affected by outliers and is thus a resistant measure. For example, if we replaced 98 with 200,000,000 in the data set from Example 4.5, the median would not change at all. The median does possess one limitation: it cannot be applied to nominal data. In order to calculate the median, the data must be placed in order. To accomplish this task meaningfully, the level of measurement must be at least ordinal.

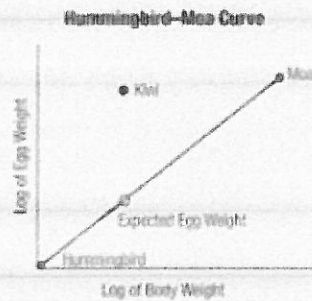
Unless the data set is skewed or contains outliers, the median and the mean usually have similar values.

Sep 26-8:56 AM

### Kiwi Eggs are Outliers!

Both plant and animal kingdoms offer spectacularly odd and beautiful sights. A kiwi bird (one of the many interesting life forms from New Zealand) lays eggs that are close to 25% of its body weight and sometimes lays two or three such eggs at a time. For most species of birds, eggs usually correspond to about 5% of the bird's body weight.

If you draw a graph relating (log) egg weight to (log) body weight you get a so-called hummingbird-moa curve (Moa is an extinct ostrich-like bird of the New Zealand area). In this curve, the kiwi bird is an outlier. Using the kiwi body weight (about 5 lb) one expects an egg weight of about 55 to 100 grams while the actual weight of kiwi eggs is about 400 to 435 grams. This egg weight matches an expected body weight of about 40 lb according to the hummingbird-moa curve. Why is this the case, and what accounts for such an anomaly?



The most reasonable explanation provided by biologists is that kiwis and moa birds are members of the same species except the kiwis have dwarfed through their evolutionary history. A sub-area of biology called "allometry" states that as body size decreases the internal organs decrease relatively slowly which supports the dwarfism hypothesis. The kiwis have lost body weight but not their internal wood-structure which still holds large eggs. Outliers are important because they force you to think about data more seriously.

Sep 26-8:56 AM

### Mode

The mode is another measure of location. It is not used as frequently as the mean or the median, and its relation to these values is not so predictable. The mode is the only measure of location that can be used for nominal data. Of the three measures of location, the mode is used the least due to the limited information it provides. Sometimes sorting the data (in ascending or descending order) makes it easier to find the mode.

#### Definition

The mode of a data set is the most frequently occurring value.

Sep 26-8:57 AM

**Example 4.7**

Find the mode of the following data set.

**0, 1, 4, 3, 9, 8, 10, 0, 1, 3, 0**

**Solution**

Since the value of 0 occurs more than any other value, it is the mode. In this instance, as a measure of location, the modal value is not a particularly appealing choice. However, the mode does possess one very favorable property—it is the only measure of location that can be applied to nominal data. Thus, for nominal measurements like color preferences, it would be perfectly reasonable to discuss the modal color.

---

Sep 26-8:58 AM

Suppose we added one more value to the data set in Example 4.7. If this value were a 1, then both 0 and 1 would be repeated three times and there would be two modes. When this occurs, the data is said to be **bimodal**. Any time data has more than two modes it is said to be **multimodal**. If all observations in a data set occur with the same frequency, then there is **no mode** for that data set.

**0, 1, 4, 3, 9, 8, 10, 0, 1, 3, 0**

Sep 26-8:59 AM

### The Relationship between the Mean, Median, and Mode

Oftentimes, the shape of the data determines how the mean, median, and mode are related. For a bell-shaped distribution, the mean, median, and mode are identical.

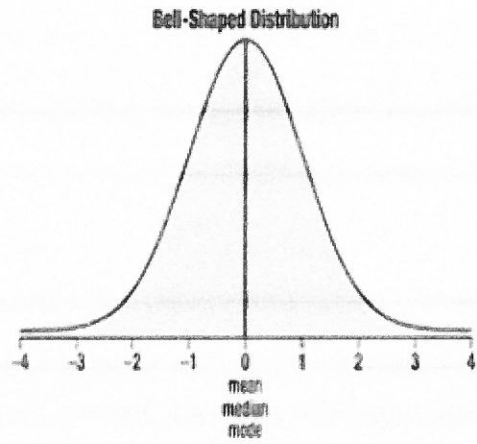


Figure 4.6

Sep 26-9:00 AM

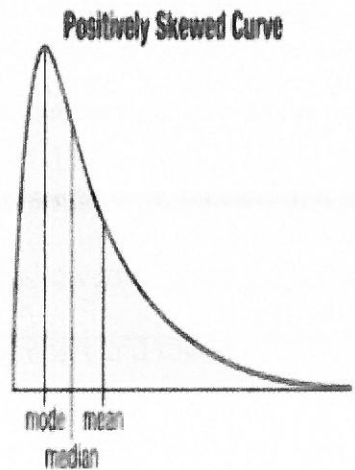


Figure 4.7

If the data are positively skewed, the median will be smaller than the mean.

Sep 26-9:00 AM



### Selecting a Measure of Location

The objective of using descriptive statistics is to provide measures that convey useful summary information about the data. When selecting a statistic to represent the central value of a data set, the first thing to consider is the type of data being analyzed.

Table 4.5 – Applicable Levels of Measurement					Table 4.6 – Sensitivity to Outliers			
Applicable Level of Measurement						Not Sensitive	Very Sensitive	
					Qualitative		Quantitative	
					Nominal	Ordinal	Interval	Ratio
Mean			✓	✓			✓	
Median		✓	✓	✓	✓			
Mode	✓	✓	✓	✓	✓			
Trimmed Mean			✓	✓		✓		

Sep 26-9:01 AM

The median is also a good measure of central tendency. It is not sensitive to outliers and can be applied to data gathered from all levels of measurement except nominal.

If the level of measurement of the data is interval or ratio and there are no outliers, the mean is a reasonable choice. If the data set appears to have any unusual values, then the trimmed mean or the median would be more appropriate.

If the data's level of measurement is nominal or ordinal (the data are qualitative), appropriate measures of center are limited. If the data are ordinal, then the median is the best choice. If the data are nominal, there is only one choice, the mode. The mode is applicable to any level of data, although it is usually not very useful for quantitative data.

Sep 26-9:04 AM

## Time Series Data and Measures of Location

We discussed two types of time series data in an earlier lesson, stationary and nonstationary. Stationary time series wobbled around some central value, so calculating a central value is perfectly reasonable, and the methods we previously discussed are applicable. A nonstationary time series is another story. Nonstationary time series possess trend. That means there is no central value for the time series. Instead, the series trends in one direction or another. Computing a central value using the methods discussed earlier would be inappropriate for such data.

Table 4.7 shows the average U.S. gas price over a 20-year period. In this nonstationary time series, the central value of the process is trending upward as shown in Figure 4.9. One way to capture this movement is with a **moving average**.

Year	Average U.S. Gas Price	2-Period Moving Average	3-Period Moving Average
1991	1.14		
1992	1.13	1.135	
1993	1.11	1.120	1.127
1994	1.11	1.110	1.117
1995	1.15	1.130	1.123
1996	1.23	1.190	1.163
1997	1.23	1.230	1.203
1998	1.06	1.145	1.173
1999	1.17	1.115	1.153

Sep 26-9:04 AM

## Definition

A **moving average** is obtained by adding consecutive observations for a number of periods and dividing the result by the number of periods included in the average.

The moving average can be used to forecast the new level of the series over time or as a descriptive method. By averaging just 2 or 3 periods at a time we can still see long-term trends, but at the same time smooth out some short-term variability in the time series. How the average is associated with a specific period is dependent on its purpose. We will assume the moving average is to be used as a method of forecasting the next level of the time series. Suppose a two-period moving average is calculated for the gas price data and is used to specify the level of the series at a given point in time. The two-period moving average for 1992 averages the values of the time series in 1991 and 1992.

$$\frac{1.14 + 1.13}{2} = 1.135$$

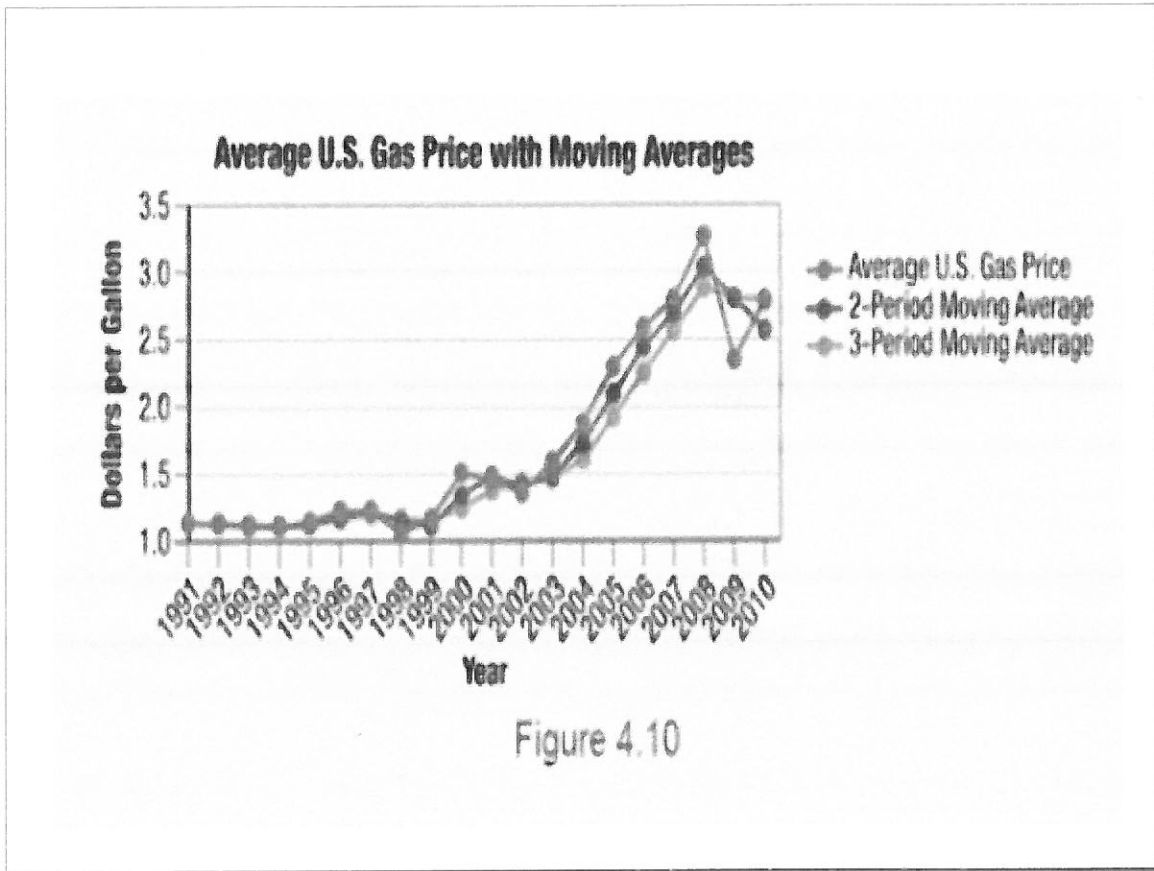
Similarly, the two-period moving average for 1993 would be the average of the time series values in 1992 and 1993.

$$\frac{1.13 + 1.11}{2} = 1.120$$

Since data are not available for 1989 or 1990, the three-period moving average for 1991 cannot be calculated. The three-period moving average associated with 1993 is the average of the time series values in 1991, 1992, and 1993.

$$\frac{1.14 + 1.13 + 1.11}{3} \approx 1.127$$

Sep 26-9:04 AM



Sep 26-9:05 AM

Consider the following data

0, 0, 9, 9, -13, 0, 9

Step 1 of 3: Determine the mean of the given data.

Sep 26-9:06 AM

Consider the following data.

0, 0, 9, 9, -13, 0, 9

Step 2 of 3 : Determine the median of the given data.

Sep 26-9:06 AM

Consider the following data.

0, 0, 9, 9, -13, 0, 9

Step 3 of 3 : Determine if the data set is unimodal, bimodal, multimodal, or has no mode. Identify the mode(s), if any exist.

Sep 26-9:07 AM