

## Chapter 2: Data

What are data?

In order to determine the context of data, consider the “W’s”

- Who –
- What (and in what units) –
- When –
- Where –
- Why –
- How –

There are two major ways to treat data:

- A \_\_\_\_\_ is used to answer questions about how cases fall into categories. A categorical variable may be comprised of word labels, or it may use numbers as labels.

Examples:

- A \_\_\_\_\_ is used to answer questions about the quantity of what is being measured. A quantitative variable is comprised of numeric values.

Examples:

What is a statistic?

Are the numbers **17, 21, 44, 76** data?

Data must have \_\_\_\_\_ to be meaningful. The numbers listed above could be test scores, ages of a group of golfers, or the uniform numbers of the starting backfield on the football team. Without \_\_\_\_\_ data cannot be interpreted.

Suppose a Consumer Reports article (published in June 2005) on energy bars gave the brand name, flavor, price, number of calories, and grams of protein and fat. Identify the following:

- Who:
- What:
- When:
- Where:
- How:
- Why:
- Categorical variables:
- Quantitative variables (with units):

A report on the Boston Marathon listed each runner's gender, county, age, and time. Identify the following:

- Who:
- What:
- When:
- Where:
- How:
- Why:
- Categorical variables:
- Quantitative variables (with units):

## Chapter 2: Data

What are data?

*Data are values along with their context. Data can be numbers or labels.*

In order to determine the context of data, consider the “W’s”

- Who – *the cases (about whom the data was collected). People are referred to as **respondents**, **subjects**, or **participants**, while objects are referred to as **experimental units**.*
- What (and in what units) – *the variables recorded about each individual.*
- When – *when the data was collected.*
- Where – *where the data was collected.*
- Why – *why the data was collected. This can determine whether a variable is treated as **categorical** or **quantitative**.*
- How – *how the data was collected.*

There are two major ways to treat data: **categorical** and **quantitative**.

- A **categorical variable** names categories and is used to answer questions about how cases fall into those categories. A categorical variable may be comprised of word labels, or it may use numbers as labels.
- A **quantitative variable** is used to answer questions about the quantity of what is being measured. A quantitative variable is comprised of numeric values.

What is a statistic? *A statistic is a numerical summary of data.*

**17, 21, 44, 76**

Are the numbers listed above data? Data must have **context** to be meaningful. The numbers listed above could be test scores, ages of a group of golfers, or the uniform numbers of the starting backfield on the football team. Without **context**, data cannot be interpreted.

Suppose a Consumer Reports article (published in June 2005) on energy bars gave the brand name, flavor, price, number of calories, and grams of protein and fat. Identify the following:

- Who: *energy bars*
- What: *brand, flavor, price, calories, protein, fat*
- When: *not specified*
- Where: *not specified*
- How: *not specified (nutrition label? laboratory testing?)*
- Why: *to inform potential consumers*
- Categorical variables: *brand, flavor*

- Quantitative variables (with units): *price (US\$), number of calories (calories), protein (grams), fat (grams)*

A report on the Boston Marathon listed each runner's gender, county, age, and time. Identify the following:

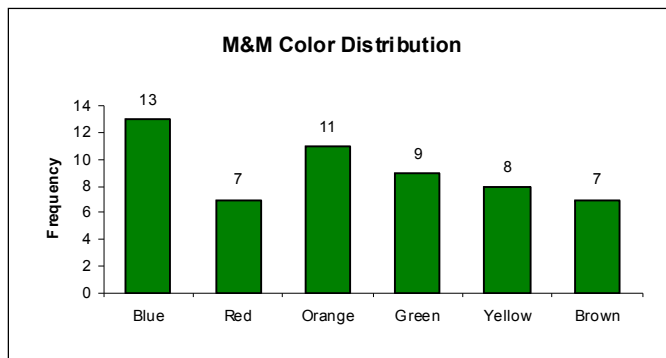
- Who: *Boston Marathon runners*
- What: *gender, county, age, time*
- When: *not specified*
- Where: *Boston*
- How: *not specified (registration information?)*
- Why: *race result reporting*
- Categorical variables: *gender, county*
- Quantitative variables (with units): *age (years), time (hours, minutes, seconds)*

### Chapter 3 – Displaying and Describing Categorical Data

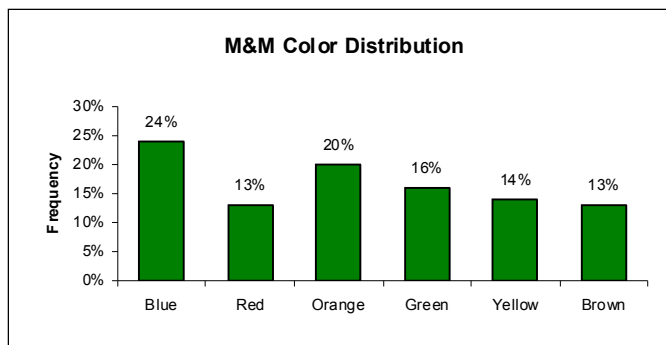
\_\_\_\_\_ are often used to organize categorical data. Frequency tables display the category names and the \_\_\_\_\_ of the number of data values in each category. \_\_\_\_\_ also display the category names, but they give the \_\_\_\_\_ rather than the counts for each category.

Color	Freq.	Rel. Freq.	Percent
Blue	13		
Red	7		
Orange	11		
Green	9		
Yellow	8		
Brown	7		
<b>TOTAL</b>	<b>55</b>	<b>1.000</b>	<b>100%</b>

A \_\_\_\_\_ is often used to display categorical data. The height of each bar represents the \_\_\_\_\_ for each category. Bars are displayed next to each other for easy comparison. When constructing a bar chart, note that the bars do not \_\_\_\_\_ one another. Categorical variables usually cannot be ordered in a meaningful way; therefore the order in which the bars are displayed is often meaningless.



A \_\_\_\_\_ bar chart displays the proportion of counts for each category.



The sum of the relative frequencies is \_\_\_\_\_.

A \_\_\_\_\_ is another type of display used to show categorical data. Pie charts show parts of a whole. Pie charts are often difficult to construct by hand.

A \_\_\_\_\_ shows two categorical variables together. The margins give the frequency distributions for each of the variables, also called the \_\_\_\_\_.

Examine the class data about gender and political view – liberal, moderate, conservative.

	Liberal	Moderate	Conservative	TOTAL
Male				
Female				
TOTAL				

- What percent of the class are girls with liberal political views?
- What percent of the liberals are girls?
- What percent of the girls are liberals?
- What is the marginal distribution of gender?
- What is the marginal distribution of political views?

A conditional distribution shows the distribution of one variable for only the individuals who satisfy some condition on another variable.

The conditional distribution of political preference, conditional on being male:

	Liberal	Moderate	Conservative	TOTAL
Male				

The conditional distribution of political preference, conditional on being female:

	Liberal	Moderate	Conservative	TOTAL
Female				

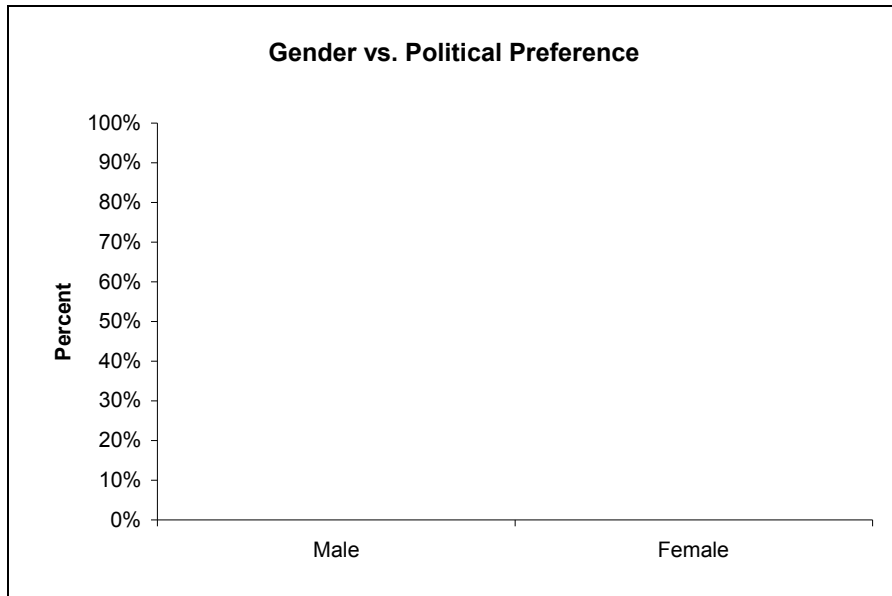
- What is the conditional relative frequency distribution of gender among conservatives?

If the conditional distributions are the same, we can conclude that the variables are not associated. Therefore, they are \_\_\_\_\_ of one another.

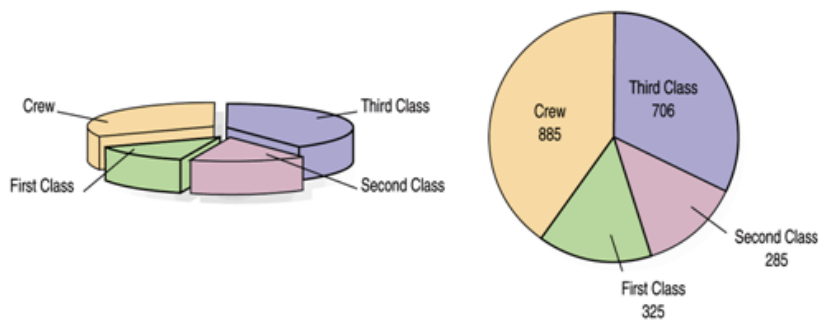
If the conditional distributions differ, we can conclude that the variables are somehow associated. Therefore, they are \_\_\_\_\_ of one another.

- Are gender and political view independent?

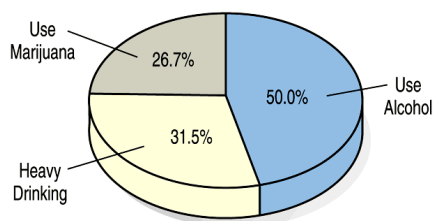
A segmented bar chart displays the same information as a pie chart, but in the form of bars instead of circles. Comparing segmented bar charts is a good way to tell if two variables are independent of one another or not.



- Explain how the graph on the left violates the “area principle.”



- Explain what is wrong with the graph below.



Averaging one variable across different levels of a second variable can lead to  
 \_\_\_\_\_. Consider the following example:

It's the last inning of an important game. Your team is a run down with the bases loaded and two outs. The pitcher is due up, so you'll be sending in a pinch-hitter. There are 2 batters available on the bench. Whom should you send in to bat?

Player	Overall
A	33 for 103
B	45 for 151

- Compare A's batting average to B's batting average. Which player appears to be the better choice?

Does it matter whether the pitcher throws right- or left-handed?

Player	Overall	vs LHP	vs RHP
A	33 for 103	28 for 81	5 for 22
B	45 for 151	12 for 32	33 for 119

- Compare A's batting average vs. a left-handed pitcher to B's. Compare A's batting average against a right-handed pitcher. Which player appears to be the better choice?

Pooling the data together loses important information and sometimes leads to the wrong conclusion. We always should take into account any factor that might matter.

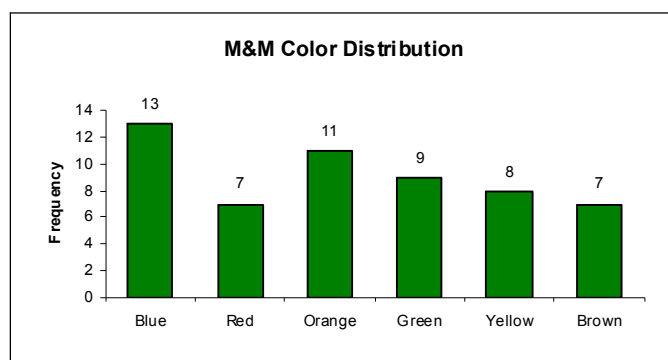


## Notes: Displaying and Describing Categorical Data

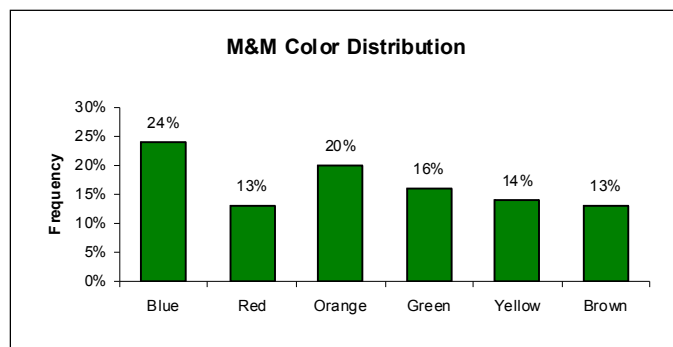
Frequency tables are often used to organize categorical data. Frequency tables display the category names and the counts of the number of data values in each category. Relative frequency tables also display the category names, but they give the percentages rather than the counts for each category.

Color	Freq.	Rel. Freq.	Percent
Blue	13	0.236	24%
Red	7	0.127	13%
Orange	11	0.200	20%
Green	9	0.164	16%
Yellow	8	0.145	14%
Brown	7	0.127	13%
<b>TOTAL</b>	<b>55</b>	<b>1.000</b>	<b>100%</b>

A bar chart is often used to display categorical data. The height of each bar represents the count for each category. Bars are displayed next to each other for easy comparison. When constructing a bar chart, note that the bars do not touch one another. Categorical variables usually cannot be ordered in a meaningful way; therefore the order in which the bars are displayed is often meaningless.



A relative frequency bar chart displays the proportion of counts for each category.



The sum of the relative frequencies is 100%.

A pie chart is another type of display used to show categorical data. Pie charts show parts of a whole. Pie charts are often difficult to construct by hand.

A contingency table shows two categorical variables together. The margins give the frequency distributions for each of the variables, also called the marginal distribution.

Examine the class data about gender and political view – liberal, moderate, conservative.

	Liberal	Moderate	Conservative	TOTAL
Male				
Female				
TOTAL				

- What percent of the class are girls with liberal political views?
- What percent of the liberals are girls?
- What percent of the girls are liberals?
- What is the marginal distribution of gender?
- What is the marginal distribution of political views?

A conditional distribution shows the distribution of one variable for only the individuals who satisfy some condition on another variable.

The conditional distribution of political preference, conditional on being male:

	Liberal	Moderate	Conservative	TOTAL
Male				

The conditional distribution of political preference, conditional on being female:

	Liberal	Moderate	Conservative	TOTAL
Female				

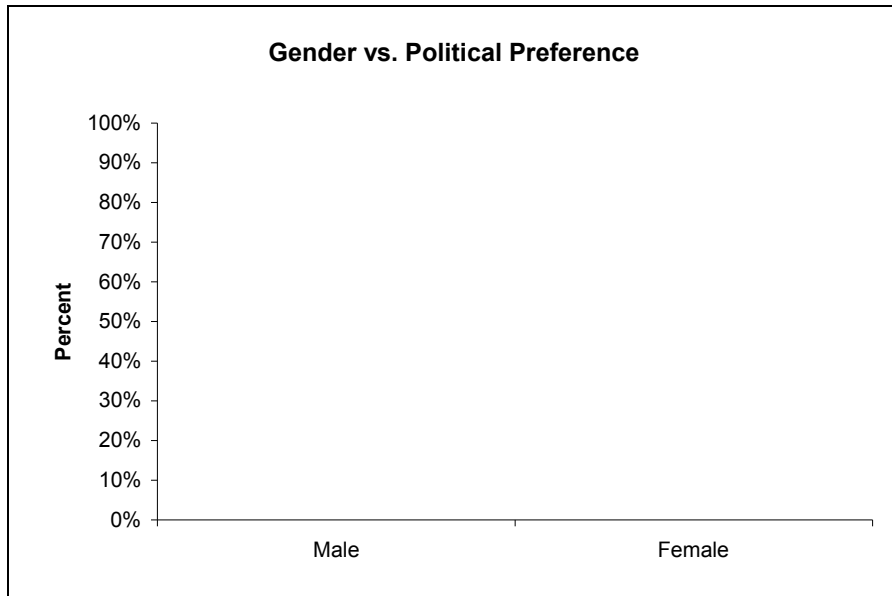
- What is the conditional relative frequency distribution of gender among conservatives?

If the conditional distributions are the same, we can conclude that the variables are not associated. Therefore, they are independent of one another.

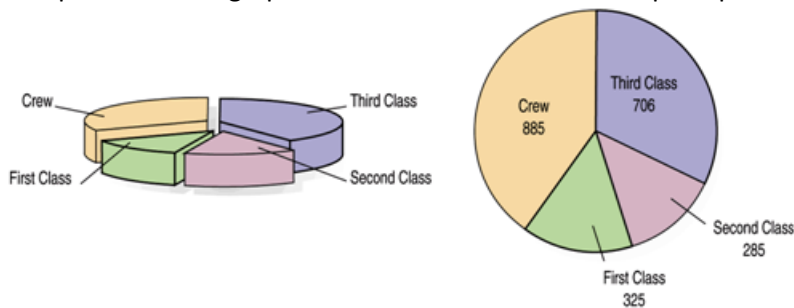
If the conditional distributions differ, we can conclude that the variables are somehow associated. Therefore, they are not independent of one another.

- Are gender and political view independent?

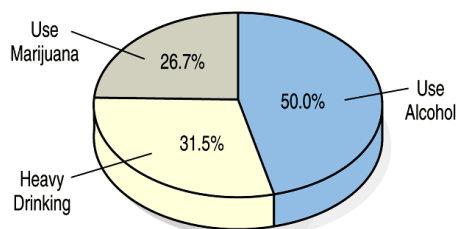
A segmented bar chart displays the same information as a pie chart, but in the form of bars instead of circles. Comparing segmented bar charts is a good way to tell if two variables are independent of one another or not.



- Explain how the graph on the left violates the “area principle.”



- Explain what is wrong with the graph below.



Averaging one variable across different levels of a second variable can lead to **Simpson's Paradox**.

Consider the following example:

It's the last inning of an important game. Your team is a run down with the bases loaded and two outs. The pitcher is due up, so you'll be sending in a pinch-hitter. There are 2 batters available on the bench. Whom should you send in to bat?

Player	Overall
A	33 for 103
B	45 for 151

- Compare A's batting average to B's batting average. Which player appears to be the better choice?

Player A has a higher batting average (0.320 vs. 0.298), so he looks like the better choice.

Does it matter whether the pitcher throws right- or left-handed?

Player	Overall	vs LHP	vs RHP
A	33 for 103	28 for 81	5 for 22
B	45 for 151	12 for 32	33 for 119

- Compare A's batting average vs. a left-handed pitcher to B's. Compare A's batting average against a right-handed pitcher. Which player appears to be the better choice?

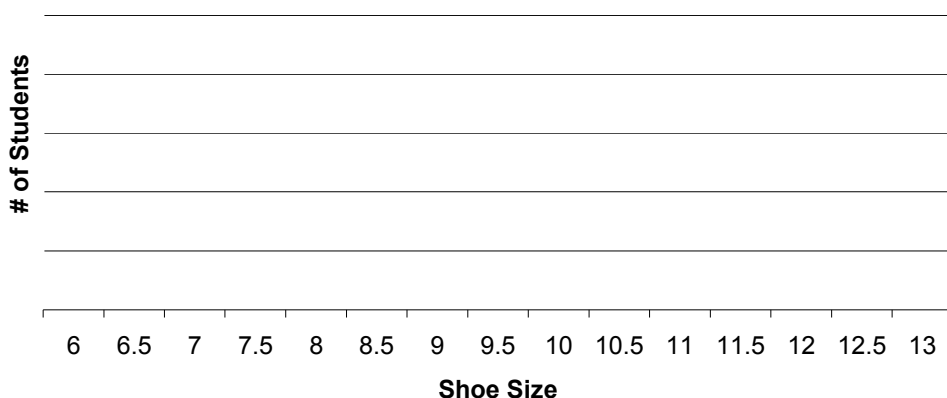
Player B has a higher batting average against both right- and left-handed pitching, even though his overall average is lower. Player B hits better against both right- and left-handed pitchers. So no matter the pitcher, B is a better choice. So why is his batting "average" lower? Because B sees a lot more right-handed pitchers than A, and (at least for these guys) right-handed pitchers are harder to hit. For some reason, A is used mostly against left-handed pitchers, so A has a higher average.

Pooling the data together loses important information and leads to the wrong conclusion. We always should take into account any factor that might matter.

**Notes: Displaying Quantitative Data**

A \_\_\_\_\_ or \_\_\_\_\_ is often used to display categorical data. These types of displays, however, are not appropriate for quantitative data. Quantitative data is often displayed using either a \_\_\_\_\_ or a \_\_\_\_\_.

In a histogram, the interval corresponding to the width of each bar is called a \_\_\_\_\_. A histogram displays the bin counts as the height of the bars (like a bar chart). Unlike a bar chart, however, the bars in a histogram \_\_\_\_\_ one another. An empty space between bars represents a \_\_\_\_\_ in data values. If a value falls on the border between two consecutive bars, it is placed in the bin on the \_\_\_\_\_.

**Shoe Sizes of AP Stat Students**

A \_\_\_\_\_ histogram displays the proportion of cases in each bin instead of the count.

Histograms are useful when \_\_\_\_\_, and they can easily be constructed using a graphing calculator. A disadvantage of histograms is that they \_\_\_\_\_.

Be sure to choose an appropriate bin width when constructing a histogram. As a general rule of thumb, your histogram should contain about \_\_\_\_\_ bars.

A \_\_\_\_\_ is similar to a histogram, but it shows \_\_\_\_\_ rather than bars. It may be necessary to \_\_\_\_\_ stems if the range of data values is small.

**Number of Pairs of Shoes Owned**

0  
0  
1  
1  
2  
2  
3  
3

KEY:

A \_\_\_\_\_ stem-and-leaf plot can be useful when \_\_\_\_\_ two distributions.

**Number of Pairs of Shoes Owned**

Male	Female
	0
	0
	1
	1
	2
	2
	3
	3

KEY:

The stems of the stem-and-leaf plot correspond to the \_\_\_\_\_ of a histogram. You may only use \_\_\_\_\_ digit for the leaves. Round or truncate your values if necessary.

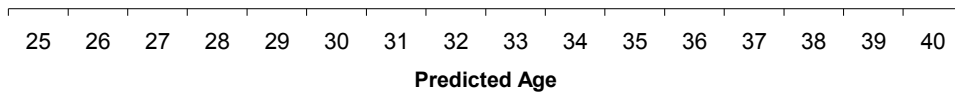
Stem-and-leaf plots are useful when working with sets of data that are \_\_\_\_\_ in size, and when you want to display \_\_\_\_\_.

How would you setup the following stem-and-leaf plots?

- ❖ quiz scores (out of 100)
- ❖ student GPA's
- ❖ student weights
- ❖ SAT scores
- ❖ weights of cattle (1000-2000 pounds)

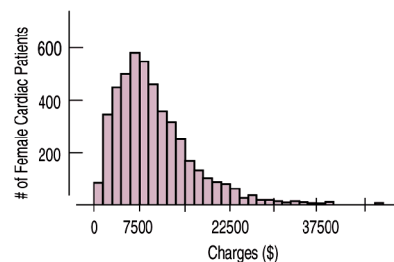
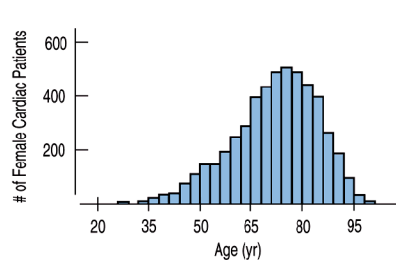
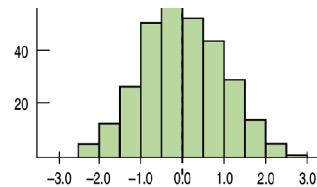
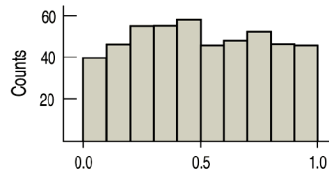
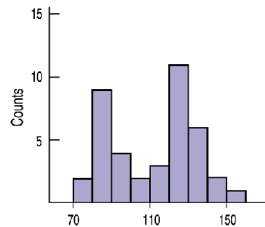
\_\_\_\_\_ may also be used to display quantitative variables. Dot plots are useful when working with \_\_\_\_\_ sets of data.

### Guess Your Teacher's Age



When describing a distribution, you should tell about three things: \_\_\_\_\_, \_\_\_\_\_, and \_\_\_\_\_. You should also mention any unusual features, like \_\_\_\_\_ or \_\_\_\_\_.

Identify the shapes of the following distributions:



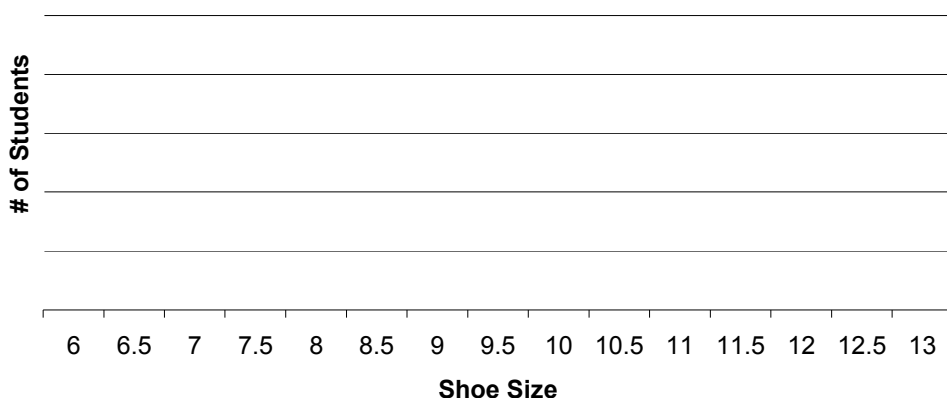
When comparing two or more distributions, compare the \_\_\_\_\_, \_\_\_\_\_, and \_\_\_\_\_, and compare any \_\_\_\_\_ features. It is important, when comparing distributions, that their graphs be constructed using the same \_\_\_\_\_.

You can sometimes make a skewed distribution appear more symmetric by \_\_\_\_\_ (or transforming) your data.

**Notes: Displaying Quantitative Data**

A [bar chart](#) or [pie chart](#) is often used to display categorical data. These types of displays, however, are not appropriate for quantitative data. Quantitative data is often displayed using either a [histogram](#), [dot plot](#), or a [stem-and-leaf plot](#).

In a histogram, the interval corresponding to the width of each bar is called a [bin](#). A histogram displays the bin counts as the height of the bars (like a bar chart). Unlike a bar chart, however, the bars in a histogram [touch](#) one another. An empty space between bars represents a [gap](#) in data values. If a value falls on the border between two consecutive bars, it is placed in the bin on the [right](#).

**Shoe Sizes of AP Stat Students**

A [relative frequency](#) histogram displays the proportion of cases in each bin instead of the count.

Histograms are useful when [working with large sets of data](#), and they can easily be constructed using a graphing calculator. A disadvantage of histograms is that they [do not show individual values](#).

Be sure to choose an appropriate bin width when constructing a histogram. As a general rule of thumb, your histogram should contain about [10](#) bars.

A [stem-and-leaf plot](#) is similar to a histogram, but it shows [individual values](#) rather than bars. It may be necessary to [split](#) stems if the range of data values is small.

**Number of Pairs of Shoes Owned**

0  
0  
1  
1  
2  
2  
3  
3

KEY:



A back-to-back stem-and-leaf plot can be useful when comparing two distributions.

Number of Pairs of Shoes Owned		
Male	Female	.
	0	
	0	
	1	
	1	
	2	
	2	
	3	
	3	

KEY:

The stems of the stem-and-leaf plot correspond to the bins of a histogram. You may only use one digit for the leaves. Round or truncate your values if necessary.

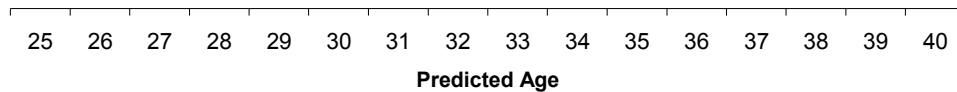
Stem-and-leaf plots are useful when working with sets of data that are small to moderate in size, and when you want to display individual values.

How would you setup the following stem-and-leaf plots?

- ❖ quiz scores (out of 100)
- ❖ student GPA's
- ❖ student weights
- ❖ SAT scores
- ❖ weights of cattle (1000-2000 pounds)

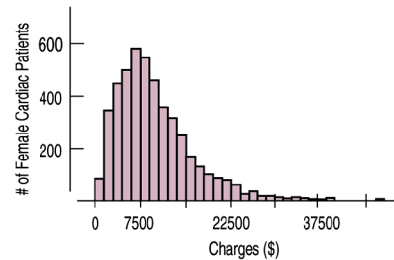
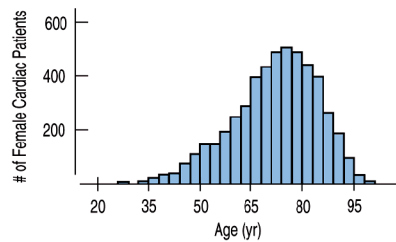
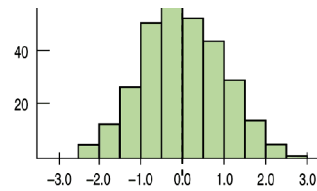
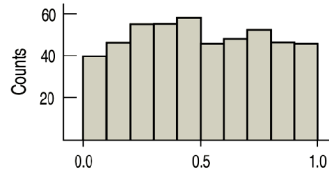
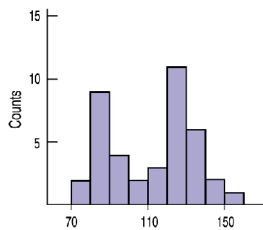
Dot plots may also be used to display quantitative variables. Dot plots are useful when working with small sets of data.

### Guess Your Teacher's Age



When describing a distribution, you should tell about three things: shape, center, and spread. You should also mention any unusual features, like outliers or gaps.

Identify the shapes of the following distributions:



When comparing two or more distributions, compare the shapes, centers, and spreads, and compare any unusual features. It is important, when comparing distributions, that their graphs be constructed using the same scale.

You can sometimes make a skewed distribution appear more symmetric by re-expressing (or transforming) your data.

**Notes: Describing Distributions Numerically**

When describing distributions, we need to discuss \_\_\_\_\_, \_\_\_\_\_, and \_\_\_\_\_. How we measure the center and spread of a distribution depends on its \_\_\_\_\_. The center of a distribution is a “typical” value. If the shape is unimodal and symmetric, a “typical” value is in the \_\_\_\_\_. If the shape is skewed, however, a “typical” value is not necessarily in the middle.

For \_\_\_\_\_ distributions, use the \_\_\_\_\_ to determine the \_\_\_\_\_ of the distribution and the \_\_\_\_\_ to describe the \_\_\_\_\_ of the distribution.

The median:

- is the \_\_\_\_\_ data value (when the data have been \_\_\_\_\_) that divides the histogram into two equal \_\_\_\_\_
- has the same \_\_\_\_\_ as the data
- is \_\_\_\_\_ to outliers (extreme data values)

The range:

- is the difference between the \_\_\_\_\_ value and the \_\_\_\_\_ value
- is a \_\_\_\_\_, NOT an \_\_\_\_\_
- is \_\_\_\_\_ to outliers

The interquartile range (IQR):

- contains the \_\_\_\_\_ of the data
- is the difference between the \_\_\_\_\_ and \_\_\_\_\_ quartiles
- is a \_\_\_\_\_, NOT an \_\_\_\_\_
- is \_\_\_\_\_ to outliers

The \_\_\_\_\_ gives: \_\_\_\_\_

A graphical display of the five-number summary is called a \_\_\_\_\_.

How many hours, on average, do you spend watching TV per week? \_\_\_\_\_ Collect data from the entire class and record the values in order from smallest to largest. Calculate the five-number summary:

Construct both a histogram and a boxplot (using the same scale). Compare the displays.

## Average Number of Hours per Week Spent Watching TV

For \_\_\_\_\_ distributions, use the \_\_\_\_\_ to determine the \_\_\_\_\_ of the distribution and the \_\_\_\_\_ to describe the \_\_\_\_\_ of the distribution.

The mean:

- is the arithmetic \_\_\_\_\_ of the data values
- is the \_\_\_\_\_ of a histogram
- has the same \_\_\_\_\_ as the data
- is \_\_\_\_\_ to outliers
- is given by the formula

The standard deviation:

- measures the “typical” distance each data value is from the \_\_\_\_\_
- Because some values are above the mean and some are below the mean, finding the sum is not useful (positives cancel out negatives); therefore we first \_\_\_\_\_ the deviations, then calculate an \_\_\_\_\_. This is called the \_\_\_\_\_. This statistics does not have the same units as the data, since we squared the deviations. Therefore, the final step is to take the \_\_\_\_\_ of the variance, which gives us the \_\_\_\_\_.
- is given by the formula
- is \_\_\_\_\_ to outliers, since its calculation involves the \_\_\_\_\_

Find the mean and standard deviation of the average number of hours spent watching TV per week for this class.

**Notes: Describing Distributions Numerically**

When describing distributions, we need to discuss shape, center, and spread. How we measure the center and spread of a distribution depends on its shape. The center of a distribution is a “typical” value. If the shape is unimodal and symmetric, a “typical” value is in the middle. If the shape is skewed, however, a “typical” value is not necessarily in the middle.

For skewed distributions, use the median to determine the center of the distribution and the interquartile range to describe the spread of the distribution.

The median:

- is the middle data value (when the data have been ordered) that divides the histogram into two equal areas
- has the same units as the data
- is resistant to outliers (extreme data values)

The range:

- is the difference between the maximum value and the minimum value
- is a number, NOT an interval
- is sensitive to outliers

The interquartile range (IQR):

- contains the middle 50% of the data
- is the difference between the lower (Q1) and upper (Q3) quartiles
- is a number, NOT an interval
- is resistant to outliers

The Five-Number Summary gives: minimum, lower quartile, median, upper quartile, maximum.

A graphical display of the five-number summary is called a boxplot.

How many hours, on average, do you spend watching TV per week? Collect data from the entire class and record the values in order from smallest to largest. Calculate the five-number summary:

Construct both a histogram and a boxplot (using the same scale). Compare the displays.

## Average Number of Hours per Week Spent Watching TV

For symmetric distributions, use the mean to determine the center of the distribution and the standard deviation to describe the spread of the distribution.

The mean:

- is the arithmetic average of the data values
- is the balancing point of a histogram
- has the same units as the data
- is sensitive to outliers
- is given by the formula  $\bar{x} = \frac{\sum x}{n}$

The standard deviation:

- measures the “typical” distance each data value is from the mean
- Because some values are above the mean and some are below the mean, finding the sum is not useful (positives cancel out negatives); therefore we first square the deviations, then calculate an adjusted average. This is called the variance. This statistics does not have the same units as the data, since we squared the deviations. Therefore, the final step is to take the square root of the variance, which gives us the standard deviation.
- is given by the formula  $s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$
- is sensitive to outliers, since its calculation involves the mean

Find the mean and standard deviation of the average number of hours spent watching TV per week for this class.

**Notes: Standard Deviation and the Normal Model**

Standard deviation is a measure of spread, or \_\_\_\_\_. The smaller the standard deviation, the \_\_\_\_\_ variability is present in the data. The larger the standard deviation, the \_\_\_\_\_ variability is present in the data.

Standard deviation can be used as a ruler for measuring how an individual compares to a \_\_\_\_\_.

To measure how far above or below the mean any given data value is, we find its \_\_\_\_\_, or \_\_\_\_\_.

$$z =$$

To standardize a value, subtract the \_\_\_\_\_ and divide by the \_\_\_\_\_.

Measure your height in inches. Calculate the standardized value for your height given that the average height for women is 64.5 inches with a standard deviation of 2.5 inches and for men is 69 inches with a standard deviation of 2.5 inches. *Are you tall?*

$$z_{\text{height}} =$$

Suppose the average woman's shoe size is 8.25 with a standard deviation 1.15 and the average male shoe size is 10 with a standard deviation of 1.5. *Do you have big feet?*

$$z_{\text{shoe}} =$$

Suppose Sharon wears a size 9 shoe and Andrew wears a size 9. *Does Sharon have big feet? Does Andrew?*

$$z_{\text{Sharon}} =$$

$$z_{\text{Andrew}} =$$

In order to compare values that are measured using different scales, you must first \_\_\_\_\_ the values. The standardized values have no \_\_\_\_\_ and are called \_\_\_\_\_. Z-scores represent how far the value is above the \_\_\_\_\_ (if \_\_\_\_\_) or below the \_\_\_\_\_ (if \_\_\_\_\_).

Example:  $z = 1$  means the value is \_\_\_\_\_ standard deviation \_\_\_\_\_ the mean  
 $z = -0.5$  means the value is \_\_\_\_\_ of a standard deviation \_\_\_\_\_ the mean

The \_\_\_\_\_ the z-score, the more unusual it is.

Standardized values, because they have no units, are therefore useful when comparing values that are measured on different \_\_\_\_\_, with different \_\_\_\_\_, or from different \_\_\_\_\_.

Adding a constant to all of the values in a set of data adds the same constant to the measures of \_\_\_\_\_. It does not, however, affect the \_\_\_\_\_.

*Example:* Add 5 to each value in the given set of data (on the left) to form a new set of data (on the right). Then find the indicated measures of center and spread.

**{5, 5, 10, 35, 45}.**

**{\_\_\_\_, \_\_\_\_ , \_\_\_\_ , \_\_\_\_ , \_\_\_\_}.**

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

IQR =

SD =

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

IQR =

SD =

Multiplying a constant to all of the values in a set of data multiplies the same constant to the measures of \_\_\_\_\_ and \_\_\_\_\_.

*Example:* Multiply each value in the given set of data (on the left) by 2 to form a new set of data (on the right). Then find the indicated measures of center and spread.

**{5, 5, 10, 35, 45}.**

**{\_\_\_\_, \_\_\_\_ , \_\_\_\_ , \_\_\_\_ , \_\_\_\_}.**

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

IQR =

SD =

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

IQR =

SD =

By standardizing values, we shift the distribution so that the mean is \_\_\_\_\_, and rescale it so that the standard deviation is \_\_\_\_\_. Standardizing does not change the \_\_\_\_\_ of the distribution.

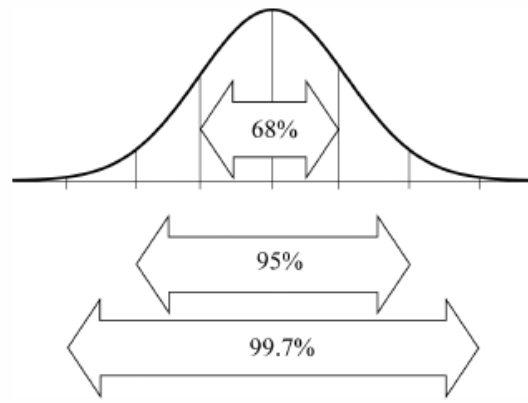
The Normal model:

❖ is \_\_\_\_\_ and \_\_\_\_\_.

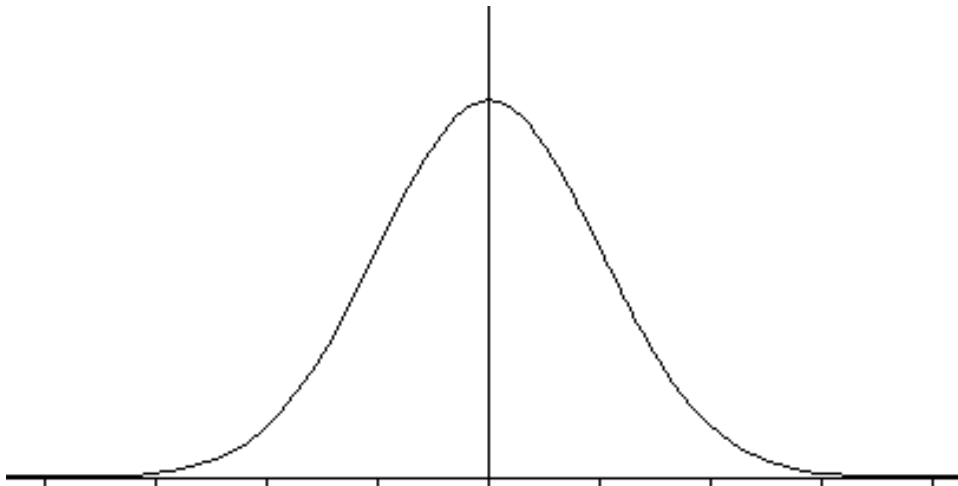
❖ follows the \_\_\_\_\_

- About \_\_\_\_\_ of the values fall within \_\_\_\_\_ standard deviation of the mean.
- About \_\_\_\_\_ of the values fall within \_\_\_\_\_ standard deviations of the mean.
- About \_\_\_\_\_ (almost all) of the values fall within \_\_\_\_\_ standard deviations of the mean.

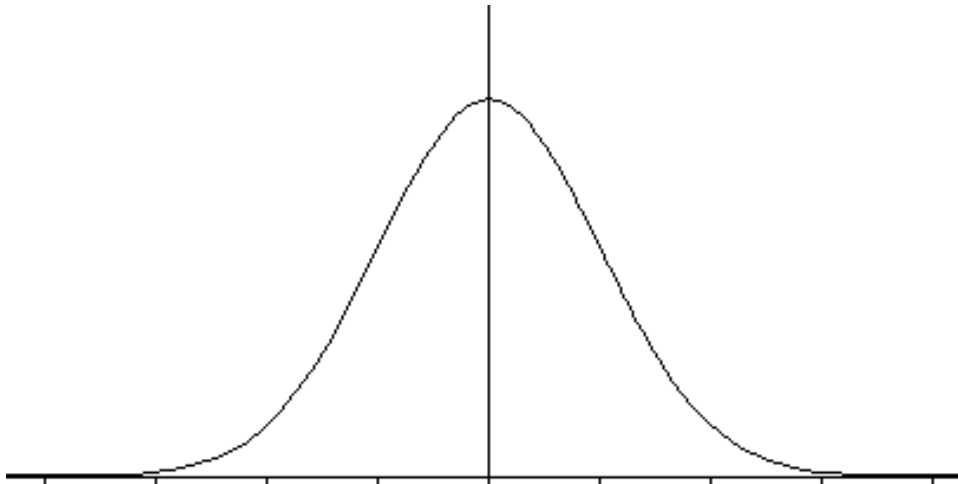




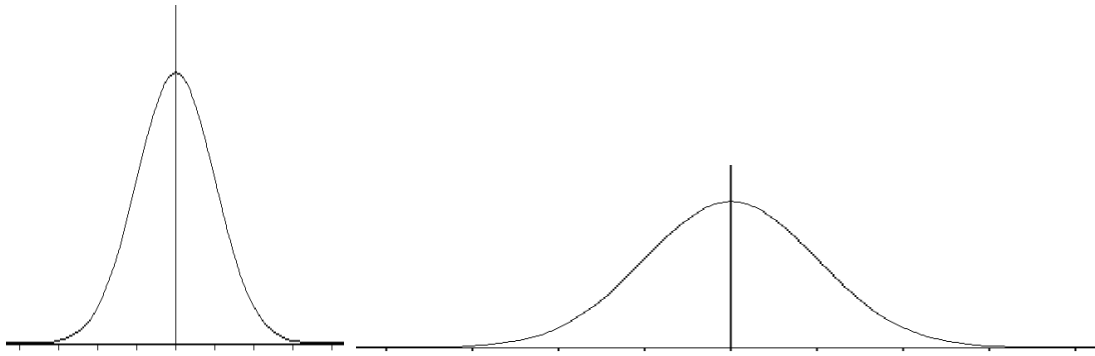
The standard Normal model has mean \_\_\_\_\_ and standard deviation \_\_\_\_\_.



The Normal model is determined by \_\_\_\_\_ and \_\_\_\_\_. We use the Greek letters sigma and mu because this is a \_\_\_\_\_; it does not come from actual \_\_\_\_\_. Sigma and mu are the \_\_\_\_\_ that specify the model.



The larger sigma, the \_\_\_\_\_ spread out the normal model appears. The inflection points occur a distance of \_\_\_\_\_ on either side of \_\_\_\_\_.



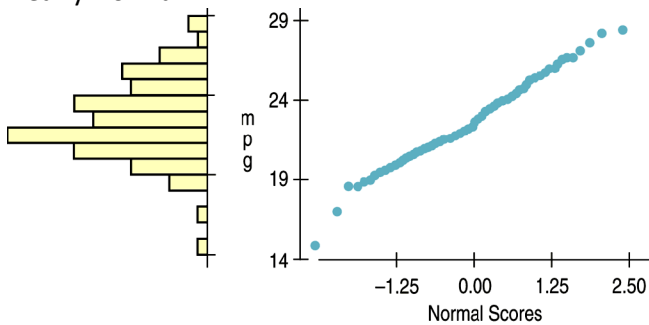
To standardize Normal data, subtract the \_\_\_\_\_ (\_\_\_\_\_) and divide by the \_\_\_\_\_ (\_\_\_\_\_).

$$z =$$

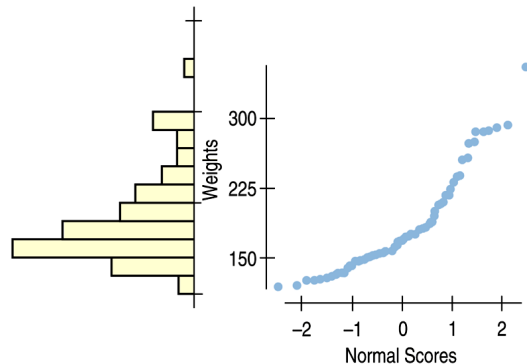
To assess normality:

- ❖ Examine the \_\_\_\_\_ of the histogram or stem-and-leaf plot. A normal model is \_\_\_\_\_ about the mean and \_\_\_\_\_.
- ❖ Compare the mean and median. In a Normal model, the mean and median are \_\_\_\_\_.
- ❖ Verify that the \_\_\_\_\_ holds.
- ❖ Construct a \_\_\_\_\_. If the graph is linear, the model is Normal.

Nearly Normal:



Skewed distribution:



**Notes: Standard Deviation and the Normal Model**

Standard deviation is a measure of spread, or **variability**. The smaller the standard deviation, the **less** variability is present in the data. The larger the standard deviation, the **more** variability is present in the data.

Standard deviation can be used as a ruler for measuring how an individual compares to a **group**.

To measure how far above or below the mean any given data value is, we find its **standardized value**, or **z-score**.

$$z = \frac{y - \bar{y}}{s}$$

To standardize a value, subtract the **mean** and divide by the **standard deviation**.

Measure your height in inches. Calculate the standardized value for your height given that the average height for women is 64.5 inches with a standard deviation of 2.5 inches and for men is 69 inches with a standard deviation of 2.5 inches. *Are you tall?*

$$z_{\text{height}} =$$

Suppose the average woman's shoe size is 8.25 with a standard deviation 1.15 and the average male shoe size is 10 with a standard deviation of 1.5. *Do you have big feet?*

$$z_{\text{shoe}} =$$

Suppose Sharon wears a size 9 shoe and Andrew wears a size 9. *Does Sharon have big feet? Does Andrew?*

$$z_{\text{Sharon}} =$$

$$z_{\text{Andrew}} =$$

In order to compare values that are measured using different scales, you must first **standardize** the values. The standardized values have no **units** and are called **z-scores**. Z-scores represent how far the value is above the **mean** (if **positive**) or below the mean (if **negative**).

Ex:      $z = 1$  means the value is **one** standard deviation **above** the mean  
            $z = -0.5$  means the value is **one-half** of a standard deviation **below** the mean

The **larger** the z-score, the more unusual it is.

Standardized values, because they have no units, are therefore useful when comparing values that are measured on different **scales**, with different **units**, or from different **populations**.

Adding a constant to all of the values in a set of data adds the same constant to the measures of center and percentiles. It does not, however, affect the spread.

*Example:* Add 5 to each value in the given set of data (on the left) to form a new set of data (on the right). Then find the indicated measures of center and spread.

**{5, 5, 10, 35, 45}.**

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

IQR =

SD =

**{10, 10, 15, 40, 50}.**

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

IQR =

SD =

Multiplying a constant to all of the values in a set of data multiplies the same constant to the measures of center and spread.

*Example:* Multiply each value in the given set of data (on the left) by 2 to form a new set of data (on the right). Then find the indicated measures of center and spread.

**{5, 5, 10, 35, 45}.**

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

IQR =

SD =

**{10, 10, 20, 70, 90}.**

Center:

$\bar{x}$  =

M =

Mode =

Spread:

Range =

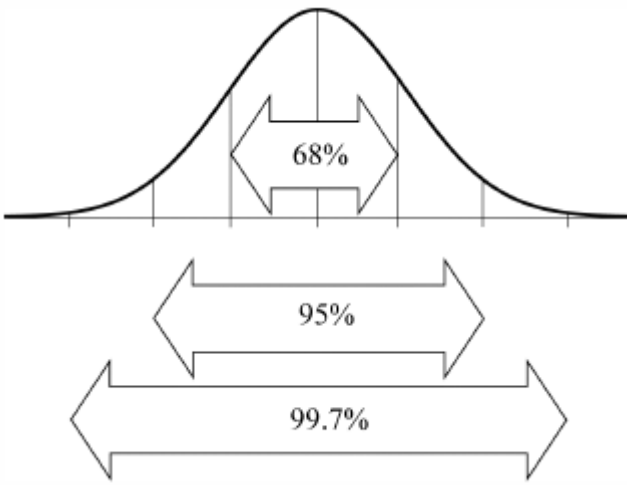
IQR =

SD =

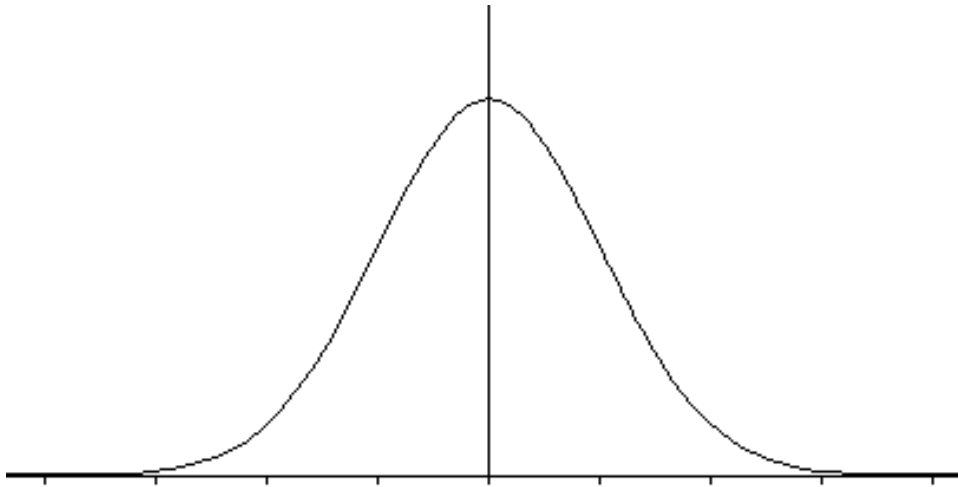
By standardizing values, we shift the distribution so that the mean is 0, and rescale it so that the standard deviation is 1. Standardizing does not change the shape of the distribution.

The Normal model:

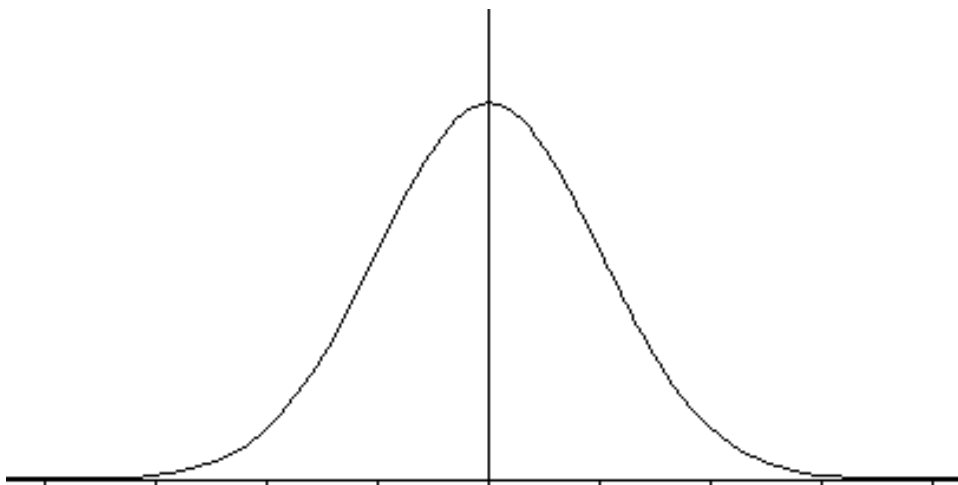
- ❖ is symmetric and bell-shaped.
- ❖ follows the 68-95-99.7 Rule
  - About 68% of the values fall within one standard deviation of the mean.
  - About 95% of the values fall within two standard deviations of the mean.
  - About 99.7% (almost all) of the values fall within three standard deviations of the mean.



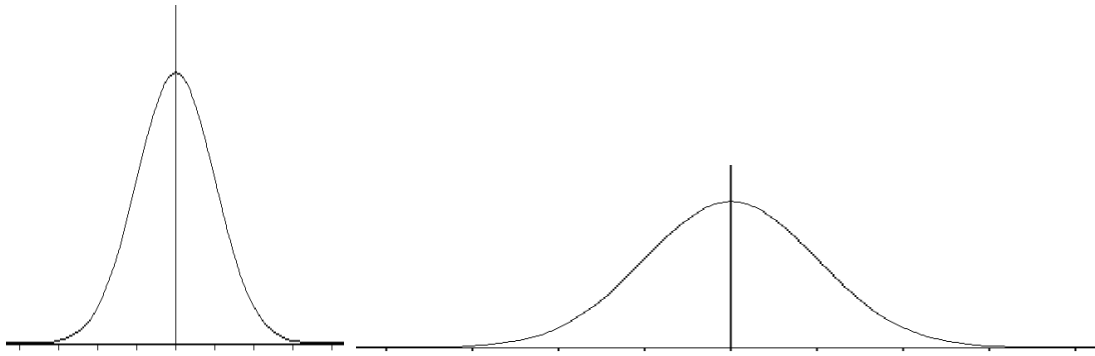
The standard Normal model has mean 0 and standard deviation 1.



The Normal model is determined by sigma and mu. We use the Greek letters sigma and mu because this is a model; it does not come from actual data. Sigma and mu are the parameters that specify the model.



The larger sigma, the **more** spread out the normal model appears. The inflection points occur a distance of **sigma** on either side of **mu**.



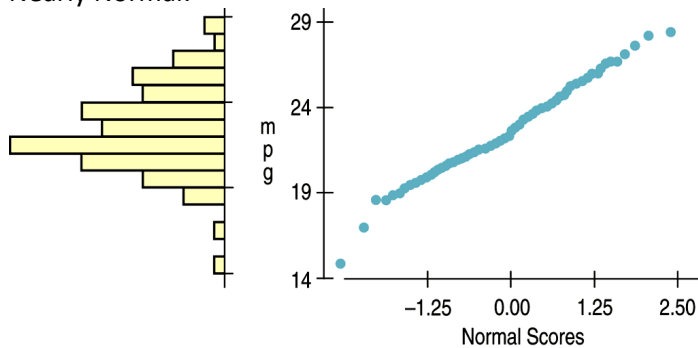
To standardize Normal data, subtract the **mean (mu)** and divide by the **standard deviation (sigma)**.

$$z = \frac{y - \mu}{\sigma}$$

To assess normality:

- ❖ Examine the **shape** of the histogram or stem-and-leaf plot. A normal model is **symmetric** about the mean and **bell-shaped**.
- ❖ Compare the mean and median. In a Normal model, the mean and median are **equal**.
- ❖ Verify that the **68-95-99.7 Rule** holds.
- ❖ Construct a **normal probability plot**. If the graph is linear, the model is Normal.

Nearly Normal:



Skewed distribution:

