

## CHAPTER PROBLEM

## Do women really talk more than men?

A common belief is that women talk more than men. Is that belief founded in fact, or is it a myth? Do men actually talk more than women? Or do men and women talk about the same amount?

In the book *The Female Brain*, neuropsychiatrist Louann Brizendine stated that women speak 20,000 words per day, compared to only 7,000 for men. She deleted that statement after complaints from linguistics experts who said that those word counts were not substantiated.

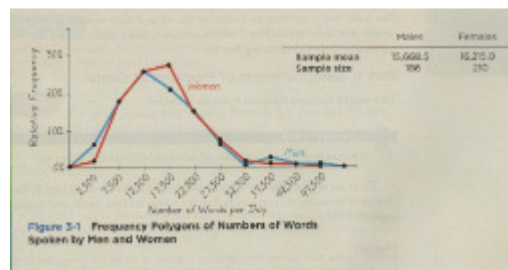
Researchers conducted a study in an attempt to address the issue of words spoken by men and women. Their findings were published in the article "Are Women Really More Talkative Than Men?" (by Mehl, Vazire, Ramirez-Esparza, Slatcher, and Pennebaker, *Science*, Vol. 317, No. 5834). The study involved 396 subjects who each wore a voice recorder that collected samples of conversations over several days.

Researchers then analyzed those conversations and counted the number of spoken words for each of the subjects.

Data Set 8 in Appendix B includes

male/female word counts from each of the six different sample groups, but if we combine all of the male word counts and all of the female word counts in Data Set 8, we get two sets of sample data that can be compared. A good way to begin to explore the data is to construct a graph that allows us to visualize the samples. See the relative frequency polygon shown below. Based on that figure, the samples of word counts from men and women appear to be very close, with no substantial differences.

When comparing the word counts of the samples of women, one step is to compare the **means** from the two samples. Shown below are the values of the means and the sample sizes. The graph and the sample means give us considerable insight into a comparison of the numbers of words spoken by men and women. In this section, we introduce other common statistical methods that are helpful in making comparisons. Using the methods of this chapter and of other chapters, we will determine whether women actually do talk more than men, or whether that is just a myth.



## MATH 103 CHAPTER 3 HOMEWORK

3.1	NA	3.3	1-5, 7, 9, 11, 13, 17, 20, 21, 25, 29, 31, 33, 35
3.2	1-5, 7, 9, 11, 17, 20, 21, 25, 29, 31, 33	3.4	1, 3, 4, 6, 7, 8, 10

## 3.1 REVIEW AND PREVIEW

Chapter 1 discussed methods of collecting \_\_\_\_\_ data, and Chapter 2 presented the \_\_\_\_\_ distribution as a tool for \_\_\_\_\_ data. Chapter 2 also presented graphs designed to help us understand some \_\_\_\_\_ of the data, including the \_\_\_\_\_.

We noted in Chapter 2 that when \_\_\_\_\_, \_\_\_\_\_, and \_\_\_\_\_ data sets, these characteristics are usually extremely important: (1) \_\_\_\_\_, (2) \_\_\_\_\_, (3) \_\_\_\_\_, (4) \_\_\_\_\_, and (5) \_\_\_\_\_ characteristics of data over time. Upon completing this chapter, you should be able to find the \_\_\_\_\_, \_\_\_\_\_, standard \_\_\_\_\_, and \_\_\_\_\_ from a data set, and you should be able to clearly understand and \_\_\_\_\_ such values.

## 3.2 MEASURES OF CENTER

Key Concept...

In this section, we discuss the characteristic of \_\_\_\_\_.

In particular, we present measures of center, including \_\_\_\_\_

and \_\_\_\_\_, as tools for \_\_\_\_\_ data.

### DEFINITION

A measure of center is a value at the \_\_\_\_\_ or \_\_\_\_\_ of a data set.

### DEFINITION

The arithmetic mean (aka mean) of a set of data is the \_\_\_\_\_ of \_\_\_\_\_ found by \_\_\_\_\_ the \_\_\_\_\_ values and \_\_\_\_\_ the total by the \_\_\_\_\_ of data values.

$$\text{mean} = \frac{\sum x}{n} = \underline{\hspace{10em}}$$

\*\*One advantage of the mean is that it is relatively \_\_\_\_\_, so that when samples are selected from the same population, sample means tend to be more consistent than other measures of center. Another advantage of the mean is that it takes every data value into account. However, because the mean is \_\_\_\_\_ to every value, just one \_\_\_\_\_ value can affect it dramatically. Because of this fact, we say the mean is not a \_\_\_\_\_ measure of center.

**NOTATION** $\Sigma$  $x$  $n$  $N$ 

$$\bar{x} = \frac{\Sigma x}{n}$$

 $\tilde{x}$ 

$$\mu = \frac{\Sigma x}{N}$$

Example 1: Find the mean of the following numbers:

17 23 17 22 21 34 27

**DEFINITION**

The median of a data set is the measure of center that is the

\_\_\_\_\_ value when the original data values are arranged

in \_\_\_\_\_ of increasing (or decreasing) magnitude. The

median is often denoted \_\_\_\_\_ (pronounced "x-tilde"). To find the

median, first \_\_\_\_\_ the values, then follow one of these two procedures:

1. If the number of data values is \_\_\_\_\_, the median is the number located in the exact \_\_\_\_\_ of the list.

2. If the number of data values is \_\_\_\_\_, the median is the \_\_\_\_\_ of the \_\_\_\_\_ two numbers.

\*\*The median is a \_\_\_\_\_ measure of center,

because it does not change by \_\_\_\_\_ amounts due to the

presence of just a few \_\_\_\_\_ values.

Example 2:

a. Find the median of the following numbers:

17 23 17 22 21 34 27

b. Find the median of the following numbers

17 23 17 22 34 27

## DEFINITION

The mode of a data set is the value that occurs with the greatest

\_\_\_\_\_ . A data set can have more than one mode, or no mode.

$\pi$  When two data values occur with the same greatest frequency, each one is

a \_\_\_\_\_ and the data set is \_\_\_\_\_ .

$\pi$  When more than two data values occur with the same greatest frequency,

each is a \_\_\_\_\_ and the data set is said to be

\_\_\_\_\_ .

$\pi$  When no data value is repeated, we say there is no \_\_\_\_\_ .

\*\*The mode is the only measure of center that can be used with data at the

\_\_\_\_\_ level of measurement.

Example 3:

a. Find the mode of the following numbers:

17 23 17 22 21 34 27

b. Find the mode of the following numbers

17 23 17 22 21 34 27 22

### DEFINITION

The midrange of a data set is the measure of center that is the value

\_\_\_\_\_ between the \_\_\_\_\_

and \_\_\_\_\_ values in the original data set. It is found by adding the maximum data value to the minimum data value and then dividing the sum by two.

midrange = \_\_\_\_\_

\*\*The midrange is rarely used because it is too sensitive to extremes since it uses only the minimum and maximum data values.

Example 4: Find the midrange of the following numbers:

17 23 17 22 21 34 27

### ROUND-OFF RULE FOR THE MEAN, MEDIAN, AND MIDRANGE

Carry \_\_\_\_\_ more decimal place than is present in the original data set. Because values of the mode are the same as some of the original data values, they can be left without any rounding.

## MEAN FROM A FREQUENCY DISTRIBUTION

When working with data summarized in a frequency distribution, we don't know the \_\_\_\_\_ values falling in a particular \_\_\_\_\_. To make calculations possible, we assume that all sample values in each class are equal to the class \_\_\_\_\_. We can then add the \_\_\_\_\_ from each \_\_\_\_\_ to find the total of all sample values, which we can then \_\_\_\_\_ by the sum of the frequencies,  $\sum f$ .

$$\text{mean from frequency distribution: } \bar{x} = \frac{\sum (f \cdot x)}{\sum f}$$

Example 5: Find the mean of the data summarized in the given frequency distribution.

Tar (mg) in nonfiltered cigarettes	Frequency
10-13	1
14-17	0
18-21	15
22-25	7
26-29	2



**WEIGHTED MEAN**

When data values are assigned different weights, we can compute a weighted mean.

$$\text{weighted mean: } \bar{x} = \frac{\sum (w \cdot x)}{\sum w}$$

Example 6: A student earned grades of 92, 83, 77, 84, and 82 on her regular tests. She earned grades of 88 on the final and 95 on her class project. Her combined homework grade was 77. The five regular tests count for 60% of the final grade, the final exam counts for 10%, the project counts for 15%, and homework counts for 15%. What is her weighted mean grade? What letter grade did she earn?

## SKEWNESS

A comparison of the \_\_\_\_\_, \_\_\_\_\_, and \_\_\_\_\_ can reveal information about the characteristic of skewness. A distribution of data is said to be \_\_\_\_\_ if it is not \_\_\_\_\_ and extends more to one side than the other.

### 3.3 MEASURES OF VARIATION

Key Concept...

In this section, we discuss the characteristic of \_\_\_\_\_.

In particular, we present measures of variation, such as

\_\_\_\_\_, \_\_\_\_\_, as tools for  
\_\_\_\_\_ data.

**DEFINITION**

The range of a set of data values is the \_\_\_\_\_  
 between the \_\_\_\_\_ and the  
 \_\_\_\_\_ data value.

**DEFINITION**

The standard deviation of a set of sample values, denoted by  $s$ , is a measure of  
 \_\_\_\_\_ of values about the \_\_\_\_\_.

It is a type of \_\_\_\_\_ deviation of values from the mean  
 that is calculated by using either of the following formulas:

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

or

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

$\pi$  The standard deviation is a measure of \_\_\_\_\_ of all

values from the \_\_\_\_\_.

$\pi$  The value of the standard deviation is usually \_\_\_\_\_.

- It is zero only when all of the data values are the same

- \_\_\_\_\_.

- It is never \_\_\_\_\_.
- π Larger values of the standard deviation indicate \_\_\_\_\_ amounts of \_\_\_\_\_.
- π The value of the standard deviation can increase dramatically with the inclusion of one or more \_\_\_\_\_.
- π The units of the standard deviation are the same units as the original \_\_\_\_\_ values.

### General Procedure for Finding Standard Deviation (1<sup>st</sup> formula)

### Specific Example Using the Following Numbers: 2, 4, 5, 16

**Step 1:** Compute the mean  $\bar{x}$

**Step 2:** Subtract the mean from each individual sample value

**Step 3:** Square each of the deviations obtained from Step 2.

**Step 4:** Add all of the squares obtained from Step 3.

**Step 5:** Divide the total from Step 4 by the number  $n - 1$ , which is one less than the total number of sample values present.

**Step 6:** Find the square root of the result from Step 5. The result is the standard deviation.

## STANDARD DEVIATION OF A POPULATION

The definition of standard deviation and the previous formulas apply to the standard deviation of \_\_\_\_\_ data. A slightly different formula is used to calculate the standard deviation  $\sigma$  of a \_\_\_\_\_: instead of dividing by  $n - 1$ , we divide by the population size  $N$ .

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}}$$

**DEFINITION**

The variance (aka dispersion aka spread) of a set of values is a measure of \_\_\_\_\_ equal to the \_\_\_\_\_ of the \_\_\_\_\_.

Sample variance:  $s^2$

Population variance:  $\sigma^2$

\*\*The sample variance is an unbiased estimator of the \_\_\_\_\_ variance, which means that values of  $s^2$  tend to target the value  $\sigma^2$  of instead of systematically tending to \_\_\_\_\_ or underestimate  $\sigma^2$ .

**USING AND UNDERSTANDING STANDARD DEVIATION**

One simple tool for understanding standard deviation is the \_\_\_\_\_ of \_\_\_\_\_, which is based on the principle that for many data sets, the vast majority (such as 95%) lie within \_\_\_\_\_ standard deviations of the \_\_\_\_\_.

**RANGE RULE OF THUMB**

Interpreting a known value of the standard deviation: We informally defined \_\_\_\_\_ values in a data set to be those that are typical and not too \_\_\_\_\_. If the standard deviation of a collection of data is \_\_\_\_\_, use it to find rough estimates of the

\_\_\_\_\_ and \_\_\_\_\_ values as follows:

minimum "usual " value = (mean) - 2 x (standard deviation)

maximum "usual " value = (mean) + 2 x (standard deviation)

**Estimating a value of the standard deviation  $s$ :** To roughly estimate the standard deviation from a collection of \_\_\_\_\_ sample data, use

$$s \approx \frac{\text{range}}{4}$$

Example 1: Use the range rule of thumb to estimate the ages of all instructors at MiraCosta if the ages of instructors are between 24 and 60.

### **EMPIRICAL (OR 68-95-99.7) RULE FOR DATA WITH A BELL-SHAPED DISTRIBUTION**

Another concept that is helpful in interpreting the value of a standard deviation is

the \_\_\_\_\_ rule. This rule states that for data sets

having a \_\_\_\_\_ that is approximately

\_\_\_\_\_, the following properties apply:

$\pi$  About 68% of all values fall within 1 standard deviation of the mean

$\pi$  About 95% of all values fall within 2 standard deviations of the mean

$\pi$  About 99.7% of all values fall within 3 standard deviations of the mean

Example 2: The author's *Generac* generator produces voltage amounts with a mean of 125.0 volts and a standard deviation of 0.3 volt, and the voltages have a bell-shaped distribution. Use the empirical to find the approximate percentage of voltage amounts between

a. 124.4 volts and 125.6 volts

b. 124.1 volts and 125.9 volts

### CHEBYSHEV'S THEOREM

The \_\_\_\_\_ (or fraction) of any data set lying within  $K$  standard deviations of the mean is always \_\_\_\_\_

$1 - \frac{1}{K^2}$ ,  $K \geq 1$ . For  $K = 2$  or  $K = 3$ , we get the following statements:

$\pi$  At least  $\frac{3}{4}$  or 75% of all values lie within 2 standard deviations of the mean.

$\pi$  At least  $\frac{8}{9}$  or 89% of all values lie within 3 standard deviations of the mean.



## COMPARING VARIATION IN DIFFERENT POPULATIONS

When comparing \_\_\_\_\_ in \_\_\_\_\_ different sets of \_\_\_\_\_, the \_\_\_\_\_ deviations should be compared only if the two sets of data use the same \_\_\_\_\_ and \_\_\_\_\_ and they have approximately the same \_\_\_\_\_.

### DEFINITION

The **coefficient of variation (aka CV)** for a set of nonnegative sample or population data, expressed as a percent, describes the standard deviation

\_\_\_\_\_ to the \_\_\_\_\_, and is given by the following:

$$\text{Sample : } CV = \frac{s}{\bar{x}} \cdot 100\%$$

$$\text{Population: } CV = \frac{\sigma}{\mu} \cdot 100\%$$

Example 3: Find the coefficient of variation for each of the two sets of data, then compare the variation.

The trend of thinner Miss America winners has generated charges that the contest encourages unhealthy diet habits among young women. Listed below are body mass indexes (BMI) for Miss America winners from two different time periods.

BMI (from the 1920s and 1930s): 20.4 21.9 22.1 22.3 20.3 18.8 18.9 19.4 18.4 19.1

BMI (from recent winners): 19.5 20.3 19.6 20.2 17.8 17.9 19.1 18.8 17.6 16.8

### 3.4 MEASURES OF RELATIVE STANDING AND BOXPLOTS

Key Concept...

In this section, we introduce measures of \_\_\_\_\_

\_\_\_\_\_, which are numbers showing the

\_\_\_\_\_ of data values \_\_\_\_\_

to the other values within a data set. The most important concept is the

\_\_\_\_\_, which will be used often in following

chapters. We will also discuss \_\_\_\_\_ and

\_\_\_\_\_, which are common statistics, as well as

a statistical graph called a \_\_\_\_\_.

### BASICS OF Z-SCORES, PERCENTILES, QUARTILES, AND BOXPLOTS

A \_\_\_\_\_ (aka standard value) is found by converting a

value to a \_\_\_\_\_ scale.

**DEFINITION**

The **z score (aka standard value)** is the number of \_\_\_\_\_ deviations a given value  $x$  is above or below the \_\_\_\_\_. The z score is calculated by using one of the following:

$$\text{Sample: } z = \frac{x - \bar{x}}{s} \quad \text{Population: } z = \frac{x - \mu}{\sigma}$$

**ROUND-OFF RULE FOR Z SCORES**

Round z scores to \_\_\_\_\_ decimal places. This rule is due to the fact that the standard table of z scores (Table A-2 in Appendix A) has z scores with two decimal places.

**Z SCORES, UNUSUAL VALUES, AND OUTLIERS**

In Section 3.3 we used the \_\_\_\_\_ of \_\_\_\_\_ to conclude that a value is \_\_\_\_\_ if it is more than 2 standard deviations away from the \_\_\_\_\_. It follows that unusual values have z scores less than \_\_\_\_\_ or greater than \_\_\_\_\_.

Example 1: The U.S. Army requires women's heights to be between 58 inches and 80 inches. Women have heights with a mean of 63.6 inches and a standard deviation of 2.5 inches. Find the z score corresponding to the minimum height requirement and find the z score corresponding to the maximum height requirement. Determine whether the minimum and maximum heights are unusual.

**PERCENTILES**

Percentiles are one type of \_\_\_\_\_ or \_\_\_\_\_

which \_\_\_\_\_ data into groups with roughly the

\_\_\_\_\_ number of values in each group.

**DEFINITION**

**Percentiles** are measures of \_\_\_\_\_, denoted \_\_\_\_\_, which divide a set of data into \_\_\_\_\_ groups with about \_\_\_\_\_ of the values in each group.

The process of finding the percentile that corresponds to a particular data value  $x$  is given by the following:

Percentile of  $x =$  -----

Example 2: Use the given sorted values, which are the number of points scored in the Super Bowl for a recent period of 24 years.

36 37 37 39 39 41 43 44 44 47 50 53 54 55 56 56 57 59 61 61 65 69 69 75

- a. Find the percentile corresponding to the given number of points.
  - i. 65

ii. 41

b. Find the indicated percentile or quartile.

i.  $Q_1$

ii.  $P_{80}$

iii.  $P_{95}$

## NOTATION

$n$

$k$

$L$

$P_k$

**DEFINITION**

Quartiles are measures of \_\_\_\_\_, denoted \_\_\_\_\_, which divide a set of data into \_\_\_\_\_ groups with about \_\_\_\_\_ of the values in each group.

**FIRST QUARTILE:**

**SECOND QUARTILE:**

**THIRD QUARTILE:**

**5 NUMBER SUMMARY AND BOXPLOT**

The values of the three \_\_\_\_\_ are used for the \_\_\_\_\_ and \_\_\_\_\_ the construction of \_\_\_\_\_ graphs.

**DEFINITION**

For a set of data, the **5-number summary** consists of the \_\_\_\_\_ value, the \_\_\_\_\_, the \_\_\_\_\_ (aka \_\_\_\_\_), the \_\_\_\_\_, and the \_\_\_\_\_ value.

A **boxplot (aka box-and-whisker diagram)** is a graph of a data set that consists of a \_\_\_\_\_ extending from the \_\_\_\_\_ value to the \_\_\_\_\_ value, and a \_\_\_\_\_ with lines drawn at the \_\_\_\_\_, the \_\_\_\_\_, and the \_\_\_\_\_.

**OUTLIERS**

When \_\_\_\_\_ data, it is important to \_\_\_\_\_ and \_\_\_\_\_ outliers because they can strongly affect values of some important statistics, such as the \_\_\_\_\_ and \_\_\_\_\_. In \_\_\_\_\_, a data value is an \_\_\_\_\_ if it is...

above quartile 3 by an amount greater than  $1.5 \times$  inner quartile range or below quartile 1 by an amount greater than  $1.5 \times$  inner quartile range

\_\_\_\_\_ are called  
\_\_\_\_\_ or \_\_\_\_\_ boxplots,  
which represent \_\_\_\_\_ as special points. A  
**modified boxplot** is a boxplot constructed with these modifications: (1) A special  
symbol, such as an \_\_\_\_\_ or point is used to identify  
\_\_\_\_\_ and (2) the solid horizontal line extends only as  
far as the minimum and maximum values which are not outliers.

Example 3: Use the given sorted values, which are the number of points scored in the Super Bowl for a recent period of 24 years to construct a boxplot. Are there any outliers?

36 37 37 39 39 41 43 44 44 47 50 53 54 55 56 56 57 59 61 61 65 69 69 75



**PUTTING IT ALL TOGETHER**

We have discussed several basic tools commonly used in statistics. When designing

an \_\_\_\_\_, \_\_\_\_\_ data, reading an article in a professional journal, or doing anything else with data, it is important to consider certain key factors, such as:

$\pi$  \_\_\_\_\_ of the data

$\pi$  \_\_\_\_\_ of the data

$\pi$  \_\_\_\_\_ method

$\pi$  Measures of \_\_\_\_\_

$\pi$  Measures of \_\_\_\_\_

$\pi$  \_\_\_\_\_

$\pi$  \_\_\_\_\_

$\pi$  Changing \_\_\_\_\_ over \_\_\_\_\_

$\pi$  \_\_\_\_\_ implications