Chapter Three

Data Description

3-1 Introduction

- Measures of average are also called <u>measures of central</u> <u>tendency</u>, which is used to summarize data, and include <u>mean</u>, <u>median</u>, <u>mode</u>, <u>midrange</u>, and <u>weighted mean</u>.
- Measures that determine the spread of data values are called <u>measures of variation</u> and include <u>range</u>, <u>variance</u>, and <u>standard deviation</u>.
- Measures of position tell where a specific data value falls within a data set or its relative position in comparison with other data values and include standard scores and <u>quartiles</u>.

3-1 Introduction

- The measures of central tendency, variation, and position are part of what is called <u>traditional statistics</u>.
- Another type of statistics is called <u>exploratory data analysis</u>, which includes the <u>box plot</u> and <u>five-number summary</u>.
- A <u>statistic</u> is a measure calculated using the data values of a sample.
- A <u>parameter</u> is a measure calculated using all the data values of a specific population.

- The <u>mean</u> is the sum of the values divided by the total number of values.
 - The Greek letter μ (mu) is used to represent the <u>population</u> <u>mean</u>.

$$\mu = \frac{\sum X}{N} = \frac{X_1 + X_2 + \dots + X_N}{N}$$

• The symbol **x** (x-bar) represents the **sample mean**.

$$\overline{x} = \frac{\sum x}{n} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

• Example:

 Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the mean.

 $\overline{x} = \frac{\sum x}{n} = \frac{61 + 11 + 1 + 3 + 2 + 30 + 18 + 3 + 7}{9} = \frac{136}{9} = 15.1$

- The median (MD) is the halfway point in a data set.
- The median is found by arranging the data in order and selecting the middle point.
- Example if *n* is odd:
 - Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the median.

$$1, 2, 3, 3, 7, 11, 18, 30, 61$$

 $MD = 7$

- Example *n* is even:
 - Suppose that the number of burglaries reported in a specific year for eight cities are 11, 1, 3, 2, 30, 18, 3, 7 Find the median.

1, 2, 3,
$$3, 7, 11, 18, 30$$

 $MD = \frac{3+7}{2} = 5$

- The **mode** is the value that occurs most often in a data set.
 - A data set with one value that occurs with greatest frequency is said to be <u>unimodal</u>.
 e.g., 3, 4, 2, 6, 4, 1, 5 → mode=4
 - A data set with two values that occur with greatest frequency is said to be <u>bimodal</u>.
 e.g., 3, 4, 2, 6, 4, 1, 2 → mode=2, 4
 - A data set with more than two values that occur with greatest frequency is said to be <u>multimodal</u>.

e.g., 6, 3, 4, 2, 6, 4, 1, 2, 5, 6, 4, 2 → *mode=2, 4,6*

 When all the values in a data set occur with the same frequency is said to have <u>no mode</u>.

> e.g., 3, 4, 2, 3, 6, 4, 1, 2, 1, 6 → *no mode* 3, 4, 2, 6, 5, 1, 7, 8 → *no mode*

• Example:

• Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the mode.

61, 11, 1, 3, 2, 30, 18, 3, 7

mode = 3

 The <u>midrange</u> (*MR*) is a rough estimate of the middle and defined as the sum of the lowest and highest values in a data set divided by 2.

$$MR = \frac{Max - Min}{2}$$

• Example:

• Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the midrange.

$$\underbrace{61}, 11, \underbrace{1}, 3, 2, 30, 18, 3, 7$$
$$MR = \frac{Min + Max}{2} = \frac{1 + 61}{2} = \frac{62}{2} = 31$$

- The <u>weighted mean</u> is used when the values in a data set are not all equally represented.
- The weighted mean of a variable X is found by multiplying each value by its corresponding weight and dividing the sum of the products by the sum of the weights.

$$\overline{x}_{w} = \frac{x_{1}w_{1} + x_{2}w_{2} + \dots + x_{n}w_{n}}{w_{1} + w_{2} + \dots + w_{n}} = \frac{\sum xw}{\sum w}$$

Where w_1, w_2, \dots, w_n are the weights for x_1, x_2, \dots, x_n

• Example:

A student received 90 in English (3 credits), 70 in Statistics (3 credits), 80 in Biology (4 credits) and 60 in physical education (2 credits), find the student's average grade.

$$\overline{x}_{w} = \frac{(90)(3) + (70)(3) + (80)(4) + (60)(2)}{3 + 3 + 4 + 2}$$
$$= \frac{270 + 210 + 320 + 120}{12} = \frac{920}{12} = 76.67$$

3-3 Properties of Central Tendency Measures

- The <u>mean</u> is affected by extremely high or low values and may not be the appropriate average.
- The <u>median</u> is <u>affected less</u> than the <u>mean</u> by extremely high or extremely low values.
- The <u>mode</u> can be used for categorical data, such as religious preference or gender.
- The <u>midrange</u> is affected by extremely high or low values in a data set.

3-4 Distribution Shapes

 In a <u>positively skewed</u> or <u>right skewed distribution</u>, the majority of the data values falls to the left of the mean and cluster at the lower end of the distribution.

mode < median < mean



3-4 Distribution Shapes

• In a **symmetrical distribution**, the data values are evenly distributed on both sides of the mean.

mean = median = mode



3-4 Distribution Shapes

 In a <u>negatively skewed</u> or <u>left skewed distribution</u>, the majority of the data values falls to the right of the mean and cluster at the upper end of the distribution.

mean < median < mode



 The <u>range</u> is the highest value minus the lowest value in a data set.

$$R = Max - Min$$

- <u>Example</u>:
 - Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the range.

$$61, 11, 11, 3, 2, 30, 18, 3, 7$$
$$R = Max - Min = 61 - 1 = 60$$

- The <u>variance</u> is the average of the squares of the distance each value is from the mean.
 - The symbol for the **population variance** is σ^2

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

• The symbol for the **sample variance** is **S**²

$$S^{2} = \frac{\sum (x - \overline{x})^{2}}{n - 1} = \frac{\sum x^{2} - \left[\frac{(\sum x)^{2}}{n}\right]}{n - 1}$$

- The **standard deviation** is the square root of the variance.
 - The symbol for the population standard deviation is σ

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

• The symbol for the **sample standard deviation** is **S**

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

• <u>Example</u>:

 Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the variance and the standard deviation.



- Variance and Standard Deviation
 - Variances and standard deviations can be used to determine the spread of the data. If the variance or standard deviation is large, the data are more dispersed.
 - The measures of variance and standard deviation are used to determine the consistency of a variable, to determine the number of data values that fall within a specified interval in a distribution and in comparing two or more data sets to determine which is more variable.

- The <u>coefficient of variation</u> is the standard deviation divided by the mean expressed as a percentage.
 - For populations

$$CVar = \frac{\sigma}{\mu} * 100 \%$$

• For Samples

$$CVar = \frac{S}{\overline{x}} * 100 \%$$

- The <u>coefficient of variation</u> is used to compare standard deviations of two variables or more when the units or the values of the means are different.
- Large coefficient of variation means large variability.

• <u>Example</u>:

• The average age of the employees at certain company is 30 years with a standard deviation of 5 years; the average salary of the employees is \$40,000 with a standard deviation of \$5000. Which one has more variation age or income?

$$CVar(age) = \frac{S}{\overline{x}} * 100 = \frac{5}{30} * 100 = 16.67\%$$

$$CVar(income) = \frac{S}{\overline{x}} * 100 = \frac{5000}{40000} * 100 = 12.5\%$$

Age is more variable than income.

- A <u>standard score</u> or <u>z score</u> is used when direct comparison of raw scores is impossible.
- The <u>z score</u> represents the number of standard deviations a data value falls above or below the mean.

$$Z=\frac{x-\overline{x}}{S}$$

• Example:

• A student scored 65 on a statistics exam that had a mean of 50 and a standard deviation of 10. Compute the *z*-score.

$$Z = \frac{x - \overline{x}}{S} = \frac{65 - 50}{10} = 1.5$$

That is, the score of 65 is 1.5 standard deviations above the mean. Above - since the *z*-score is positive.

• Example:

• Which of the following exam scores has a better relative position?

a. A score of 42 on an exam with $\overline{x} = 39$ and S = 4

$$Z = \frac{x - \overline{x}}{S} = \frac{42 - 39}{4} = \frac{3}{4} = 0.75$$

b. A score of 76 on an exam with $\overline{x} = 71$ and S = 3

$$Z = \frac{x - \overline{x}}{S} = \frac{76 - 71}{3} = \frac{5}{3} = 1.67$$

So, a score of 76 has a better relative position.

• <u>Quartiles</u> divide the distribution into four groups, denoted by Q_1, Q_2, Q_3 . Note that Q_1 is the same as the 25th percentile, Q_2 is the same as the 50th percentile or the median and Q_3 corresponds to the 75th percentile.

- <u>Quartiles</u> can be found as follow
 - 1. Arrange the data in order from lowest to highest.
 - 2. Find the median of the data values (Q_2) .
 - 3. Find the median of the data values that fall bellow $Q_2(Q_1)$.
 - 4. Find the median of the data values that fall above $Q_2(Q_3)$.

• Example:

 Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the first, second and third quartile.

$$1, [2, 3], 3, [7], 11, [18, 30], 62$$

$$Q_2 = 7$$

$$Q_1 = 2.5$$

$$Q_3 = 24$$

Outliers

- An <u>outlier</u> is an extremely high or an extremely low data value when compared with the rest of the data values.
- **<u>Outliers</u>** can be identified as follows:
 - 1. Arrange the data in order and find Q_1 and Q_3 .
 - 2. Find the **interquartile range**: *IQR=Q₃-Q₁*.
 - 3. The values that are <u>smaller</u> than Q₁-(1.5)(IQR) or <u>larger</u> than Q₃+(1.5)(IQR) are called outliers.
- Outliers can be the result of measurement or observational error.

• <u>Example</u>:

 Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the outlier values if any.

1, 2, 3, 3, 7, 11, 18, 30, 61

$$Q_2 = 7$$

 $Q_1 = 2.5$ $Q_3 = 24$

 $IQR = Q_3 - Q_1 = 24 - 2.5 = 21.5$ $Q_3 + (1.5)(IQR) = 24 + (1.5)(21.5) = 24 + 32.25 = 56.25$ $Q_1 - (1.5)(IQR) = 2.5 - (1.5)(21.5) = 2.5 - 32.25 = -29.75$ 1, 2, 3, 3, 7, 11, 18, 30, 61

61 is an outlier value in this data.

3-7 Exploratory Data Analysis

- The purpose of <u>exploratory data analysis</u> is to examine data in order to find out what information can be discovered such as the center and the spread.
- <u>Boxplots</u> are graphical representations of a <u>five-number</u> <u>summary</u> of a data set. The five specific values that make up a five-number summary are

 $Min, Q_{1'}, Q_{2'}Q_{3'}, Max$

3-7 Exploratory Data Analysis

- Information obtained from a boxplot
 - Using the box:
 - If the median is near the center of the box, the distribution is approximately symmetric.
 - If the median is to the left of the box, the distribution is positively skewed.
 - If the median is to the right of the box, the distribution is negatively skewed.
 - Using the lines:
 - If the lines are about the same length, the distribution is approximately symmetric.
 - If the right line is taller, the distribution is positively skewed.
 - If the left line is taller, the distribution is negatively skewed.

3-7 Exploratory Data Analysis

• <u>Example</u>:

 Suppose that the number of burglaries reported in a specific year for nine cities are 61, 11, 1, 3, 2, 30, 18, 3, 7 Find the fivenumber summary and comment on the skewness of the data.

$$\begin{array}{c} 1, 2, 3, 3, 7, 11, 18, 30, 61 \\ \hline \\ Q_2 = 7 \\ \hline \\ Min = 1 \\ Q_1 = 2.5 \\ \hline \\ Q_3 = 24 \\ \hline \\ Max = 61 \\ \hline \\ \\ Max = 61 \\ \hline \\ \end{array}$$



• The distribution is positively skewed.