

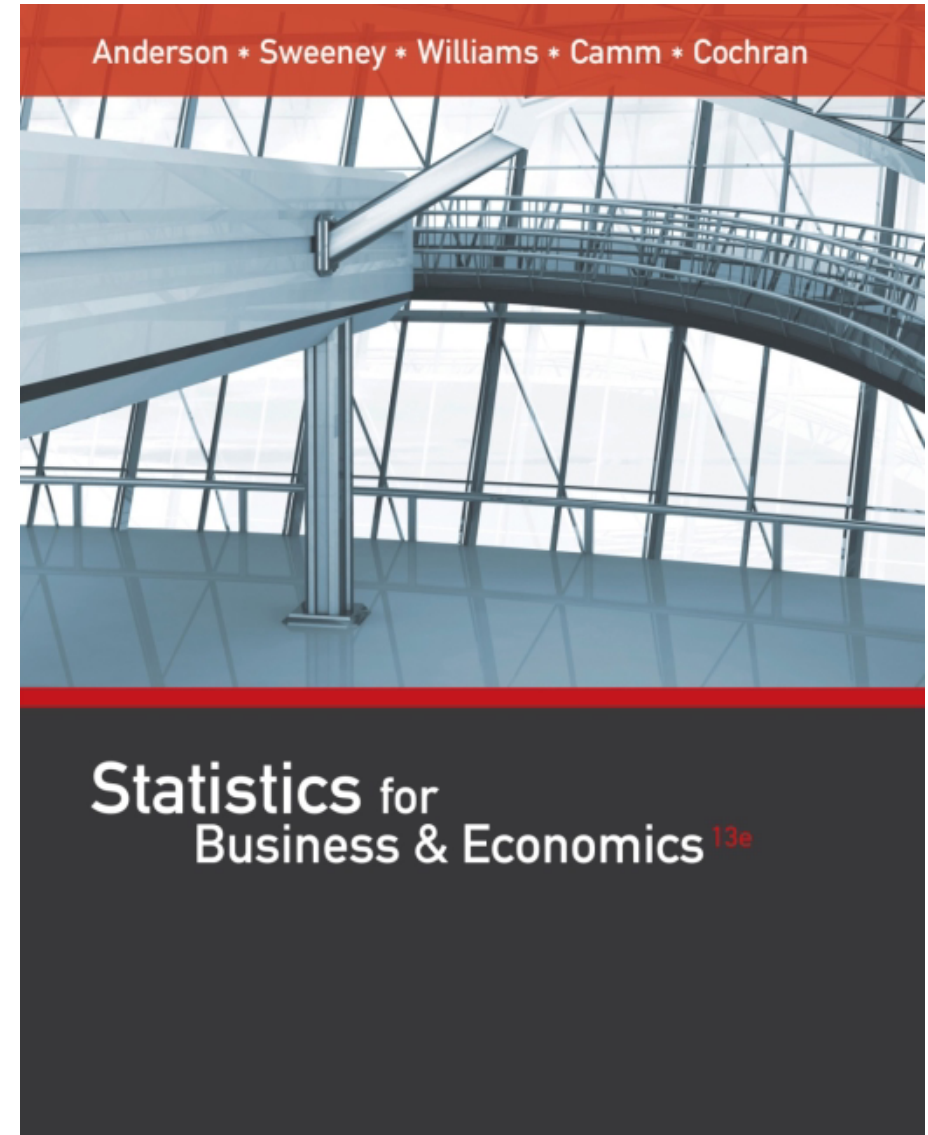
Statistics for Business and Economics (13e)

Anderson, Sweeney, Williams, Camm, Cochran

© 2017 Cengage Learning

Slides by John Loucks

St. Edwards University



Chapter 3, Part B

Descriptive Statistics: Numerical Measures

- Measures of Distribution Shape, Relative Location, and Detecting Outliers
- Five-Number Summaries and Box Plots
- Measures of Association Between Two Variables
- Data Dashboards: Adding Numerical Measures to Improve Effectiveness

Measures of Distribution Shape, Relative Location, and Detecting Outliers

- Distribution Shape
- z-Scores
- Chebyshev's Theorem
- Empirical Rule
- Detecting Outliers

Distribution Shape: Skewness

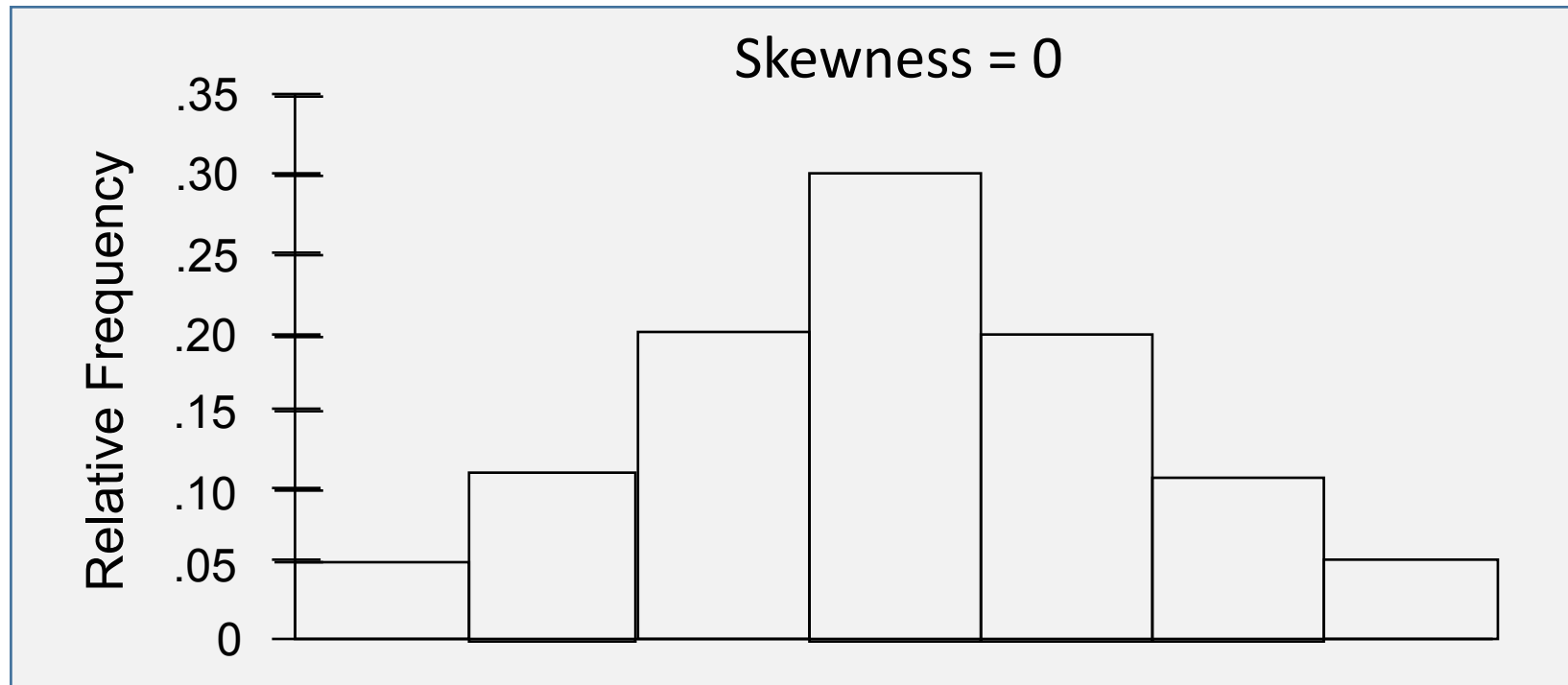
- An important numerical measure of the shape of a distribution is called skewness.
- The formula for the skewness of sample data is

$$\text{Skewness} = \frac{n}{(n-1)(n-2)} \sum \left[\frac{x_i - \bar{x}}{s} \right]^3$$

- Skewness can be easily computed using statistical software.

Distribution Shape: Skewness

- Symmetric (not skewed)
 - Skewness is zero.
 - Mean and median are equal.



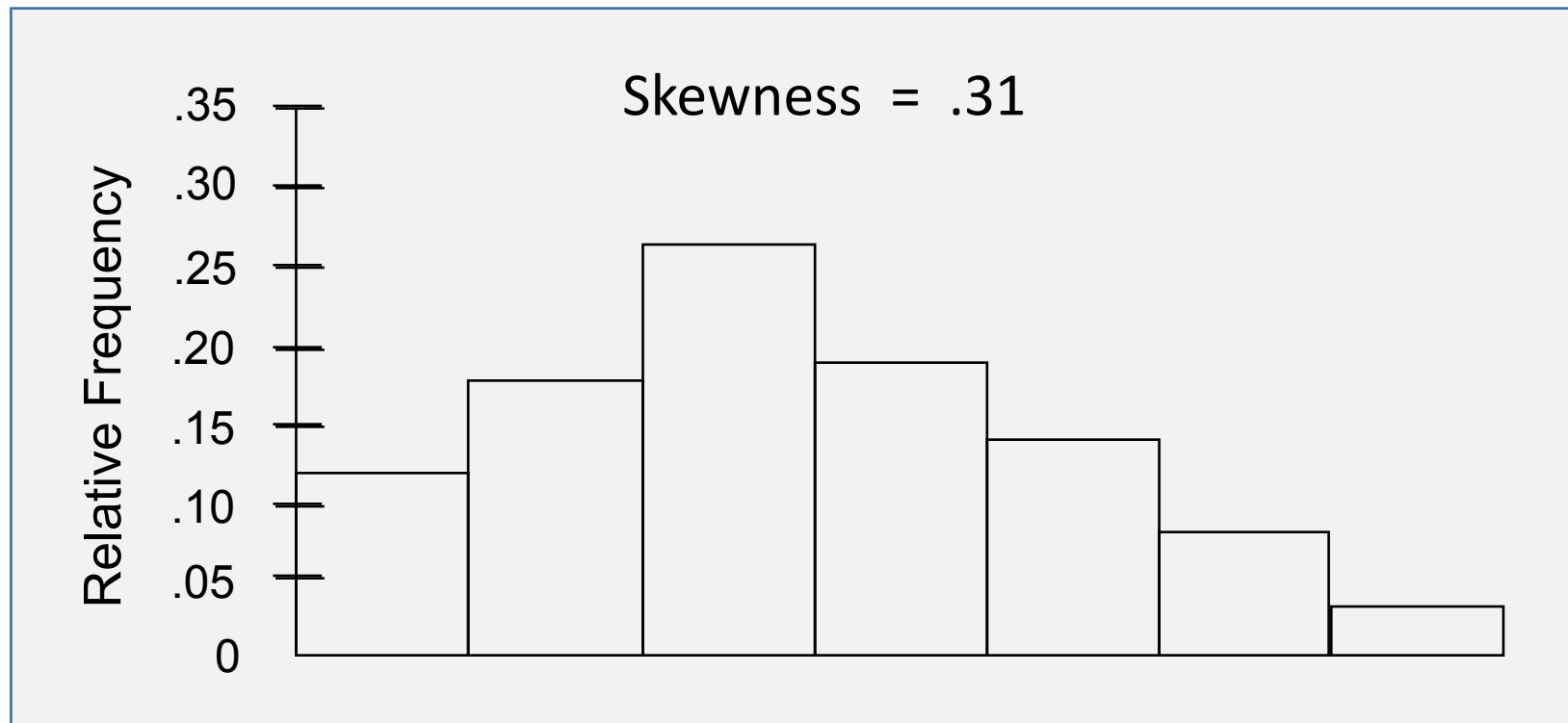
Distribution Shape: Skewness

- Moderately Skewed Left
 - Skewness is negative.
 - Mean will usually be less than the median.



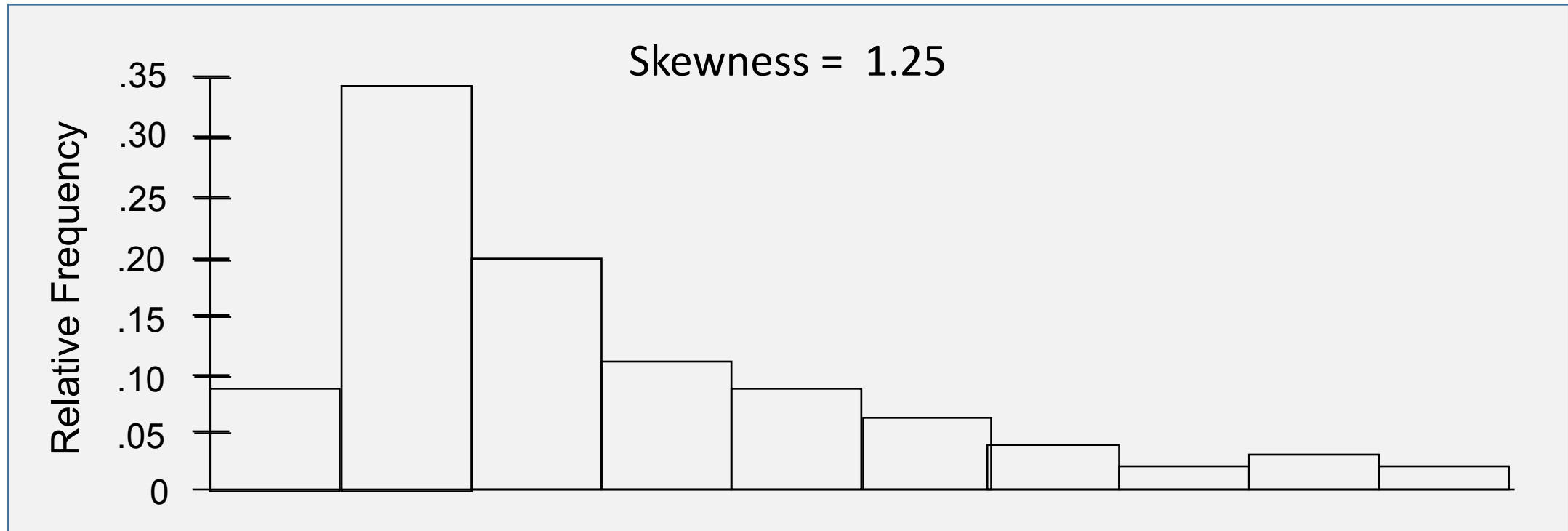
Distribution Shape: Skewness

- Moderately Skewed Right
 - Skewness is positive.
 - Mean will usually be more than the median.



Distribution Shape: Skewness

- Highly Skewed Right
 - Skewness is positive (often above 1.0).
 - Mean will usually be more than the median.



Distribution Shape: Skewness

- Example: Apartment Rents

Seventy efficiency apartments were randomly sampled in a college town. The monthly rent prices for the apartments are listed below in ascending order.

525	530	530	535	535	535	535	535	540	540
540	540	540	545	545	545	545	545	550	550
550	550	550	550	550	560	560	560	565	565
565	570	570	572	575	575	575	580	580	580
580	585	590	590	590	600	600	600	600	610
610	615	625	625	625	635	649	650	670	670
675	675	680	690	700	700	700	700	715	715

Distribution Shape: Skewness

- Example: Apartment Rents



z-Scores

- The z-score is often called the standardized value.
- It denotes the number of standard deviations a data value x_i is from the mean.

$$Z_i = \frac{x_i - \bar{x}}{s}$$

z-Scores

- An observation's z-score is a measure of the relative location of the observation in a data set.
- A data value less than the sample mean will have a z-score less than zero.
- A data value greater than the sample mean will have a z-score greater than zero.
- A data value equal to the sample mean will have a z-score of zero.

z-Scores

- Example: Apartment Rents
 - z-Score of Smallest Value (525)

$$Z_i = \frac{x_i - \bar{x}}{s} = \frac{525 - 590.80}{54.74} = -1.20$$

Standardized Values for Apartment Rents

-1.20	-1.11	-1.11	-1.02	-1.02	-1.02	-1.02	-1.02	-0.93	-0.93
-0.93	-0.93	-0.93	-0.84	-0.84	-0.84	-0.84	-0.84	-0.75	-0.75
-0.75	-0.75	-0.75	-0.75	-0.75	-0.56	-0.56	-0.56	-0.47	-0.47
-0.47	-0.38	-0.38	-0.34	-0.29	-0.29	-0.29	-0.20	-0.20	-0.20
-0.20	-0.11	-0.01	-0.01	-0.01	0.17	0.17	0.17	0.17	0.35
0.35	0.44	0.62	0.62	0.62	0.81	1.06	1.08	1.45	1.45
1.54	1.54	1.63	1.81	1.99	1.99	1.99	1.99	2.27	2.27

Chebyshev's Theorem

- At least $(1 - 1/z^2)$ of the data values must be within z standard deviations of the mean, where z is any value greater than 1.
- Chebyshev's theorem requires $z > 1$; but z need not be an integer.

Chebyshev's Theorem

- At least 75% of the data values must be within $z = 2$ standard deviations of the mean.
- At least 89% of the data values must be within $z = 3$ standard deviations of the mean.
- At least 94% of the data values must be within $z = 4$ standard deviations of the mean.

Chebyshev's Theorem

- Example: Apartment Rents

Let $z = 1.5$ with $\bar{x} = 590.80$ and $s = 54.74$

At least $(1 - 1/(1.5)^2) = 1 - 0.44 = 0.56$ or **56%**

of the rent values must be between

$$\bar{x} - z(s) = 590.80 - 1.5(54.74) = \mathbf{509}$$

and

$$\bar{x} + z(s) = 590.80 + 1.5(54.74) = \mathbf{673}$$

(Actually, 86% of the rent values
are between 509 and 673.)

Empirical Rule

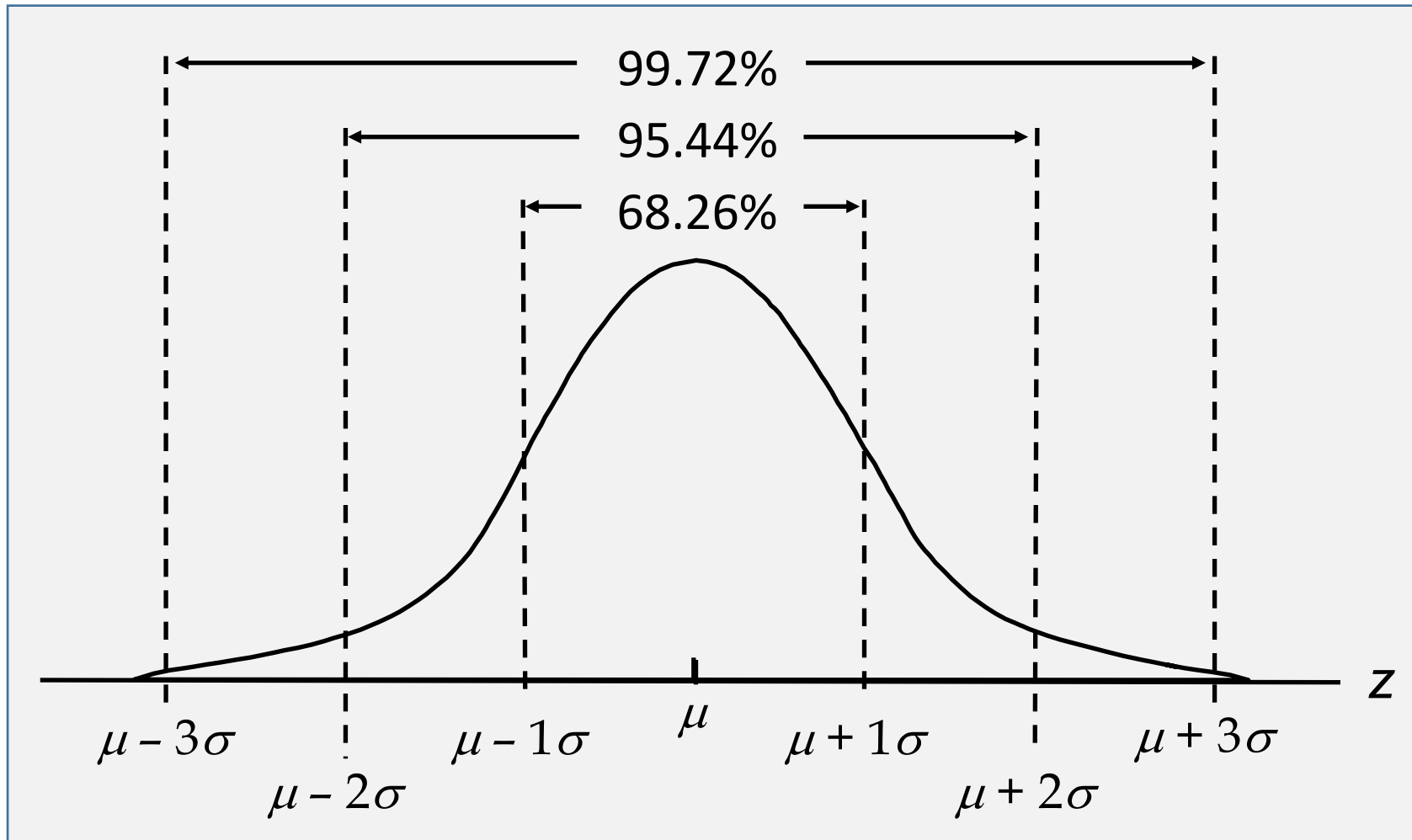
- When the data are believed to approximate a bell-shaped distribution:
 - The empirical rule can be used to determine the percentage of data values that must be within a specified number of standard deviations of the mean.
 - The empirical rule is based on the normal distribution, which is covered in Chapter 6.

Empirical Rule

For data having a bell-shaped distribution:

- Approximately 68% of the data values will be within one standard deviation of the mean.
- Approximately 95% of the data values will be within two standard deviations of the mean.
- Almost all of the data values will be within three standard deviations of the mean.

Empirical Rule



Detecting Outliers

- An outlier is an unusually small or unusually large value in a data set.
- A data value with a z-score less than -3 or greater than +3 might be considered an outlier.
- It might be:
 - an incorrectly recorded data value
 - a data value that was incorrectly included in the data set
 - a correctly recorded data value that belongs in the data set

Empirical Rule

- Example: Apartment Rents
 - The most extreme z-scores are -1.20 and 2.27.
 - Using $|z| \geq 3$ as the criterion for an outlier, there are no outliers in this data set.

Standardized Values for Apartment Rents

-1.20	-1.11	-1.11	-1.02	-1.02	-1.02	-1.02	-1.02	-0.93	-0.93
-0.93	-0.93	-0.93	-0.84	-0.84	-0.84	-0.84	-0.84	-0.75	-0.75
-0.75	-0.75	-0.75	-0.75	-0.75	-0.56	-0.56	-0.56	-0.47	-0.47
-0.47	-0.38	-0.38	-0.34	-0.29	-0.29	-0.29	-0.20	-0.20	-0.20
-0.20	-0.11	-0.01	-0.01	-0.01	0.17	0.17	0.17	0.17	0.35
0.35	0.44	0.62	0.62	0.62	0.81	1.06	1.08	1.45	1.45
1.54	1.54	1.63	1.81	1.99	1.99	1.99	1.99	2.27	2.27

Five-Number Summaries and Box Plots

- Summary statistics and easy-to-draw graphs can be used to quickly summarize large quantities of data.
- Two tools that accomplish this are five-number summaries and box plots.

Five-Number Summary

1. Smallest Value
2. First Quartile
3. Median
4. Third Quartile
5. Largest Value

Five-Number Summary

- Example: Apartment Rents

Lowest Value = 525

First Quartile = 545

Median = 575

Third Quartile = 625

Largest Value = 715

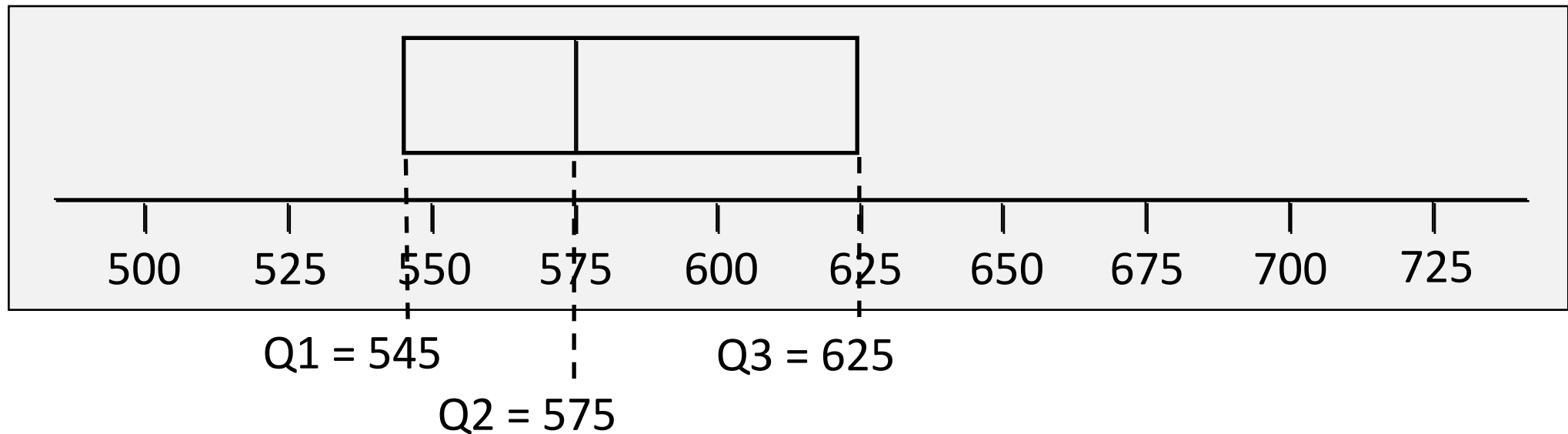
525	530	530	535	535	535	535	535	540	540
540	540	540	545	545	545	545	545	550	550
550	550	550	550	550	560	560	560	565	565
565	570	570	572	575	575	575	580	580	580
580	585	590	590	590	600	600	600	600	610
610	615	625	625	625	635	649	650	670	670
675	675	680	690	700	700	700	700	715	715

Box Plot

- A box plot is a graphical display of data that is based on a five-number summary.
- A key to the development of a box plot is the computation of the median and the quartiles Q_1 and Q_3 .
- Box plots provide another way to identify outliers.

Box Plot

- Example: Apartment Rents
 - A box is drawn with its ends located at the first and third quartiles.
 - A vertical line is drawn in the box at the location of the median (second quartile).



Box Plot

- Limits are located (not drawn) using the interquartile range (IQR).
- Data outside these limits are considered outliers.
- The location of each outlier is shown with the symbol * .

Box Plot

- Example: Apartment Rents

- The lower limit is located $1.5(\text{IQR})$ below $Q1$.

$$\text{Lower Limit: } Q1 - 1.5(\text{IQR}) = 545 - 1.5(80) = 425$$

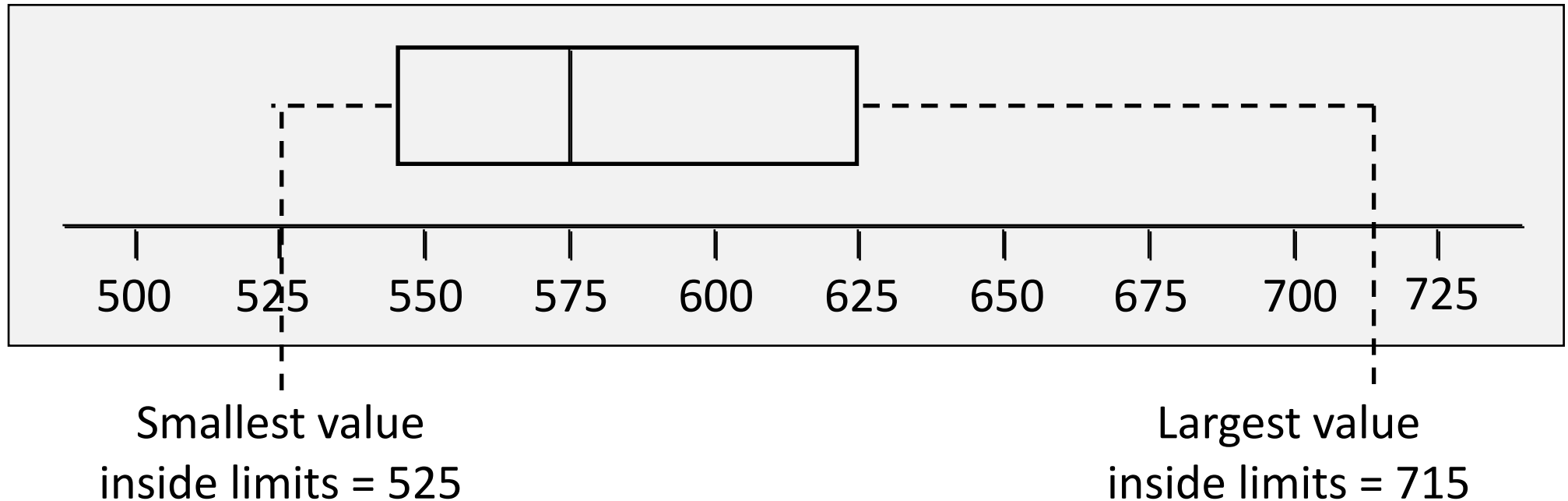
- The upper limit is located $1.5(\text{IQR})$ above $Q3$.

$$\text{Upper Limit: } Q3 + 1.5(\text{IQR}) = 625 + 1.5(80) = 745$$

- There are no outliers (values less than 425 or greater than 745) in the apartment rent data.

Box Plot

- Example: Apartment Rents
 - Whiskers (dashed lines) are drawn from the ends of the box to the smallest and largest data values inside the limits.



Measures of Association Between Two Variables

- Thus far we have examined numerical methods used to summarize the data for one variable at a time.
- Often a manager or decision maker is interested in the relationship between two variables.
- Two descriptive measures of the relationship between two variables are covariance and correlation coefficient.

Covariance

- The covariance is a measure of the linear association between two variables.
- Positive values indicate a positive relationship.
- Negative values indicate a negative relationship.

Covariance

- The covariance is computed as follows:

For samples:
$$S_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

For populations:
$$\sigma_{xy} = \frac{\sum(x_i - \mu_x)(y_i - \mu_y)}{N}$$

Correlation Coefficient

- Correlation is a measure of linear association and not necessarily causation.
- Just because two variables are highly correlated, it does not mean that one variable is the cause of the other.

Correlation Coefficient

- The correlation coefficient is computed as follows:

For samples:
$$r_{xy} = \frac{s_{xy}}{s_x s_y}$$

For populations:
$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

Correlation Coefficient

- The coefficient can take on values between -1 and +1.
- Values near -1 indicate a strong negative linear relationship.
- Values near +1 indicate a strong positive linear relationship.
- The closer the correlation is to zero, the weaker the relationship.

Covariance and Correlation Coefficient

- Example: Golfing Study

A golfer is interested in investigating the relationship, if any, between driving distance and 18-hole score.

<u>Average Driving Distance (yds.)</u>	<u>Average 18-Hole Score</u>
277.6	69
259.5	71
269.1	70
267.0	70
255.6	71
272.9	69

Covariance and Correlation Coefficient

- Example: Golfing Study

	x	y	$(x_i - \bar{x})$	$(y_i - \bar{y})$	$(x_i - \bar{x})(y_i - \bar{y})$
	277.6	69	10.65	-1.0	-10.65
	259.5	71	-7.45	1.0	-7.45
	269.1	70	2.15	0	0
	267.0	70	0.05	0	0
	255.6	71	-11.35	1.0	-11.35
	272.9	69	5.95	-1.0	-5.95
Average	267.0	70.0			Total -35.40
Std. Dev.	8.2192	.8944			

Covariance and Correlation Coefficient

- Example: Golfing Study
 - Sample Covariance

$$s_{xy} = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{n-1} = \frac{-35.40}{6-1} = -7.08$$

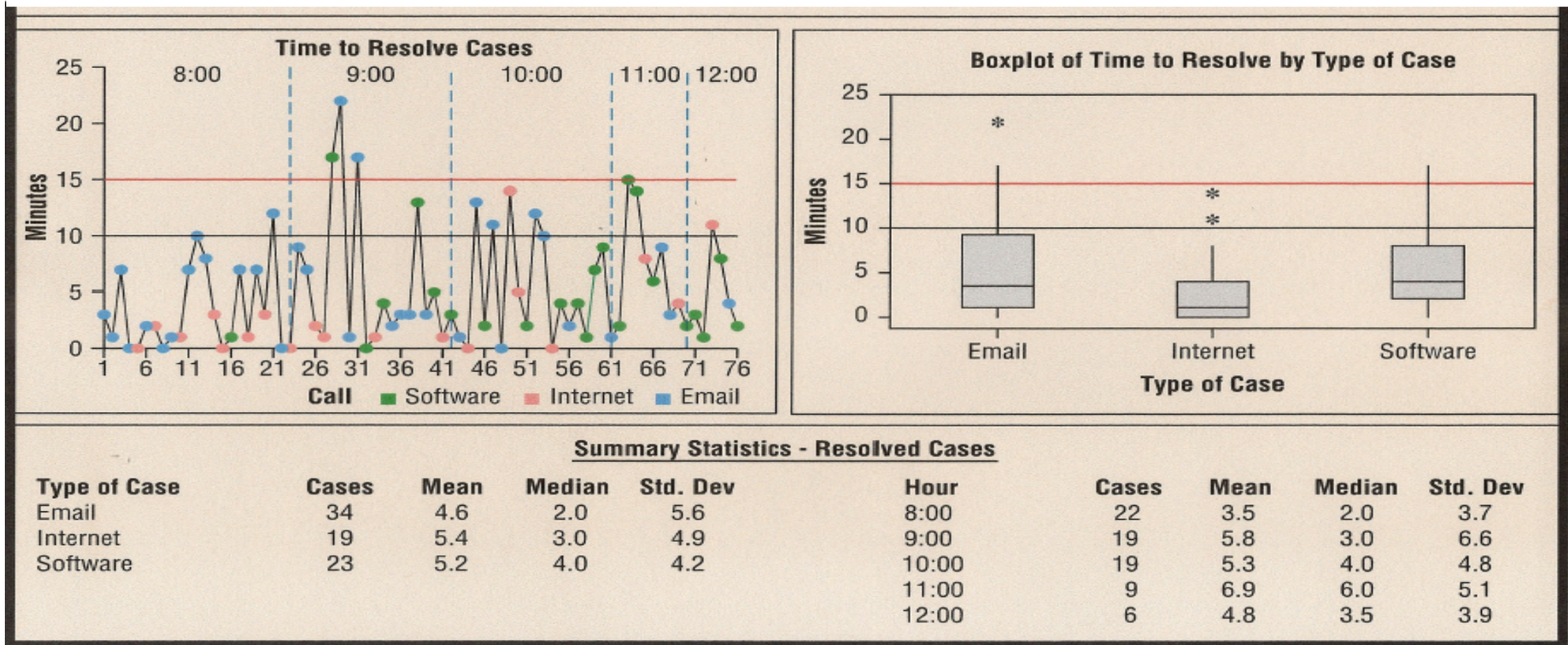
- Sample Correlation Coefficient

$$r_{xy} = \frac{s_{xy}}{s_x s_y} = \frac{-7.08}{(8.2192)(.8944)} = -.9631$$

Data Dashboards: Adding Numerical Measures to Improve Effectiveness

- Data dashboards are not limited to graphical displays.
- The addition of numerical measures, such as the mean and standard deviation of KPIs, to a data dashboard is often critical.
- Dashboards are often interactive.
- Drilling down refers to functionality in interactive dashboards that allows the user to access information and analyses at an increasingly detailed level.

Data Dashboards: Adding Numerical Measures to Improve Effectiveness



End of Chapter 3, Part B

