# Chapter 9 Slide 1

9-1 Overview

9-2 Correlation

9-3 Regression

9-4 Variation and Prediction Intervals

**9-5 Multiple Regression** 

9-6 Modeling

## Section 9-1 & 9-2 Overview and Correlation and Regression

Created by Erin Hodgess, Houston, Texas







## Paired Data (p.506)

Is there a relationship?

## If so, what is the equation?

Use that equation for prediction.

## Definition



# A correlation exists between two variables when one of them is related to the other in some way.

## Definition



A Scatterplot (or scatter) diagram) is a graph in which the paired (x, y) sample data are plotted with a horizontal x-axis and a vertical y-axis. Each individual (x, y) pair is plotted as a single point.



## Scatter Diagram of Paired Data (p.507)



Chapter 9, Triola, Elementary Statistics, MATH 1342

## Positive Linear Correlation (p.498)







(a) Positive correlation between *x* and *y* 

(**b**) Strong positive correlation between *x* and *y* 

(c) Perfect positive correlation between *x* and *y* 

#### Figure 9-2 Scatter Plots

## Negative Linear Correlation







(d) Negative correlation between x and y

(e) Strong negative correlation between *x* and *y* 

(f) Perfect negative correlation between *x* and *y* 

#### Figure 9-2 Scatter Plots

## **No Linear Correlation**



(g) No correlation between x and y

(h) Nonlinear relationship between *x* and *y* 

Slide 9

#### Figure 9-2 Scatter Plots

## **Definition (p.509)**



# The linear correlation coefficient *r* measures strength of the linear relationship between paired *x* and *y* values in a sample.



- 1. The sample of paired data (*x*, *y*) is a random sample.
- 2. The pairs of (*x*, *y*) data have a bivariate normal distribution.

#### Notation for the Linear Correlation Coefficient



- n = number of pairs of data presented
- $\Sigma$  denotes the addition of the items indicated.
- $\Sigma x$  denotes the sum of all *x*-values.
- $\Sigma x^2$  indicates that each *x*-value should be squared and then those squares added.
- $(\Sigma x)^2$  indicates that the *x*-values should be added and the total then squared.
- $\Sigma xy$  indicates that each x-value should be first multiplied by its corresponding y-value. After obtaining all such products, find their sum.
- *r* represents linear correlation coefficient for a <u>sample</u>
- $\rho$  represents linear correlation coefficient for a <u>population</u>

## Definition



The linear correlation coefficient *r* measures the strength of a linear relationship between the paired values in a sample.

$$r = \frac{n\Sigma xy - (\Sigma x)(\Sigma y)}{\sqrt{n(\Sigma x^2) - (\Sigma x)^2}\sqrt{n(\Sigma y^2) - (\Sigma y)^2}}$$

Formula 9-1

#### Calculators can compute r

## $\rho$ (rho) is the linear correlation coefficient for <u>all</u> paired data in the population.

Rounding the Linear Correlation Coefficient r



Round to three decimal places so that it can be compared to critical values in Table A-6. (see p.510)

Use calculator or computer if possible.

## Calculating r



#### Data

<i>x</i>	1	1	3	5	
у	2	8	6	4	

#### This data is from exercise #7 on p.521.

## Calculating r



Table 9-2	-2 Finding Statistics Used to Calculate <i>r</i>								
	x	У	x·y	<i>x</i> <sup>2</sup>	y <sup>2</sup>				
	1	2	2	1	4				
	1	8	8	1	64				
	3	6	18	9	36				
	5	4	20	25	16				
Total	<b>10</b> ↑ Σ <i>x</i>	<b>20</b> ↑ Σ <i>y</i>	<b>48</b> ↑ Σ <i>xy</i>	<b>36</b> ↑ Σ <i>x</i> <sup>2</sup>	<b>120</b> ↑ Σy²				

#### Calculating r



#### Data 1 3 5 x 2 8 6 4 y $n\Sigma xy - (\Sigma x)(\Sigma y)$ r = $\sqrt{n(\Sigma x^2) - (\Sigma x)^2} \sqrt{n(\Sigma y^2) - (\Sigma y)^2}$ 4(48) - (10)(20)r = $\sqrt{4(36) - (10)^2} \sqrt{4(120) - (20)^2}$ = -0.13559.329

## Interpreting the Linear Slide 18 Correlation Coefficient (p.511)

If the absolute value of r exceeds the value in Table A - 6, conclude that there is a significant linear correlation.

Otherwise, there is not sufficient evidence to support the conclusion of significant linear correlation.

#### Example: Boats and Manatees



Given the sample data in Table 9-1, find the value of the linear correlation coefficient r, then refer to Table A-6 to determine whether there is a significant linear correlation between the number of registered boats and the number of manatees killed by boats.

Using the same procedure previously illustrated, we find that r = 0.922.

Referring to Table A-6, we locate the row for which n=10. Using the critical value for  $\alpha=5$ , we have 0.632. Because r = 0.922, its absolute value exceeds 0.632, so we conclude that there is a significant linear correlation between number of registered boats and number of manatee deaths from boats.



- 1.  $-1 \le r \le 1$  (see also p.512)
- 2. Value of *r* does not change if all values of either variable are converted to a different scale.
- 3. The *r* is not affected by the choice of *x* and *y*. interchange *x* and *y* and the value of *r* will not change.
- 4. *r* measures strength of a linear relationship.

## Interpreting *r*: Explained Variation



The value of  $r^2$  is the proportion of the variation in y that is explained by the linear relationship between x and y. (p.503 and p.533)

#### Example: Boats and Manatees



Using the boat/manatee data in Table 9-1, we have found that the value of the linear correlation coefficient r = 0.922. What proportion of the variation of the manatee deaths can be explained by the variation in the number of boat registrations?

With r = 0.922, we get  $r^2 = 0.850$ .

We conclude that 0.850 (or about 85%) of the variation in manatee deaths can be explained by the linear relationship between the number of boat registrations and the number of manatee deaths from boats. This implies that 15% of the variation of manatee deaths cannot be explained by the number of boat registrations.



#### Common Errors Involving Correlation (pp.503-504)

- 1. Causation: It is wrong to conclude that correlation implies causality.
- 2. Averages: Averages suppress individual variation and may inflate the correlation coefficient.
- 3. Linearity: There may be <u>some relationship</u> between x and y even when there is no significant linear correlation.

## Common Errors Involving Correlation





#### **FIGURE 9-3**

Scatterplot of Distance above Ground and Time for Object Thrown Upward

## Formal Hypothesis Test (p.504)

## We wish to determine whether there is a significant linear correlation between two variables.

We present two methods.

#### **\*Both methods let** $H_0$ : $\rho = 0$ (no significant linear correlation) $H_1$ : $\rho \neq 0$ (significant linear correlation)





Chapter 9, Triola, Elementary Statistics, MATH 1342

## Method 1: Test Statistic is *t* (follows format of earlier chapters)

**Test statistic:** 

$$t = \frac{r}{\sqrt{\frac{1-r^2}{n-2}}}$$

**Critical values:** 

# Use Table A-3 with degrees of freedom = n - 2

## Method 2: Test Statistic is r (uses fewer calculations)



#### Test statistic: r

#### Critical values: Refer to Table A-6 (no degrees of freedom)

#### Example: Boats and Manatees



Using the boat/manatee data in Table 9-1, test the claim that there is a linear correlation between the number of registered boats and the number of manatee deaths from boats. Use Method 1.

TA

$$t = \frac{1}{\sqrt{\frac{1-r^2}{n-2}}}$$
$$t = \frac{0.922}{\sqrt{\frac{1-0.922}{10-2}^2}} = 6.735$$

#### Method 1: <u>Slide 30</u> Test Statistic is *t* (follows format of earlier chapters)



Chapter 9, Triola, Elementary Statistics, MATH 1342

#### Example: Boats and Manatees



Using the boat/manatee data in Table 9-1, test the claim that there is a linear correlation between the number of registered boats and the number of manatee deaths from boats. Use Method 2.

The test statistic is r = 0.922. The critical values of  $r = \pm 0.632$  are found in Table A-6 with n = 10 and  $\alpha = 0.05$ .

## Method 2: Test Statistic is r (uses fewer calculations)



Test statistic: r

## Critical values: Refer to Table A-6 (10 degrees of freedom)



#### Example: Boats and Manatees



Using the boat/manatee data in Table 9-1, test the claim that there is a linear correlation between the number of registered boats and the number of manatee deaths from boats. Use both (a) Method 1 and (b) Method 2.

Using either of the two methods, we find that the absolute value of the test statistic does exceed the critical value (Method 1: 6.735 > 2.306. Method 2: 0.922 > 0.632); that is, the test statistic falls in the critical region.

We therefore reject the null hypothesis. There is sufficient evidence to support the claim of a linear correlation between the number of registered boats and the number of manatee deaths from boats.

#### Justification for r Formula



#### Formula 9-1 is developed from



Chapter 9, Triola, *Elementary Statistics*, MATH 1342



Created by Erin Hodgess, Houston, Texas



## Regression



## Definition

#### Regression Equation

The regression equation expresses a relationship between x (called the independent variable, predictor variable or explanatory variable, and y (called the dependent variable or response variable.

The typical equation of a straight line y = mx + b is expressed in the form  $y = b_0 + b_1 x$ , where  $b_0$  is the yintercept and  $b_1$  is the slope.
## Assumptions



- 1. We are investigating only linear relationships.
- 2. For each x-value, y is a random variable having a normal (bell-shaped) distribution. All of these y distributions have the same variance. Also, for a given value of x, the distribution of y-values has a mean that lies on the regression line. (Results are not seriously affected if departures from normal distributions and equal variances are not too extreme.)

## Regression



## Definition

## Regression Equation

Given a collection of paired data, the regression equation

$$\hat{y} = b_0 + b_1 x$$

algebraically describes the relationship between the two variables

Regression Line The graph of the regression equation is called the regression line (or line of best fit, or least squares line).



## Notation for Regression Equation

	Population Parameter	<u>Sample</u> <u>Statistic</u>
y-intercept of regression equation	$eta_{0}$	b <sub>0</sub>
Slope of regression equation	$\beta_1$	b <sub>1</sub>
Equation of the regression line	$y = \beta_0 + \beta_1 x$	$\hat{y} = b_0 + b_1 \mathbf{X}$



Formula 9-2 
$$b_1 = \frac{n(\Sigma xy) - (\Sigma x)(\Sigma y)}{n(\Sigma x^2) - (\Sigma x)^2}$$
 (slope)

Formula 9-3 
$$b_0 = \overline{y} - b_1 \overline{x}$$
 (y-intercept)

#### calculators or computers can compute these values



$$b_0 = \bar{y} - b_1 \bar{x}$$

Formula 9-4

## Can be used for Formula 9-2, where $\overline{y}$ is the mean of the y-values and $\overline{x}$ is the mean of the x values



## The regression line fits the sample points best.



## Rounding the y-intercept $b_0$ and the slope $b_1$



If you use the formulas 9-2 and 9-3, try not to round intermediate values. (see p.527)

	Slide 44				
Data				-	
x	1	1	3	5	
у	2	8	6	4	

In Section 9-2, we used these values to find that the linear correlation coefficient of r = -0.135. Use this sample to find the regression equation.



Chapter 9, Triola, Elementary Statistics, MATH 1342

## **Calculating the Regression Equation**



Data

x	1	1	3	5	
у	2	8	6	4	

n=4 $\Sigma x = 10$  $\Sigma y = 20$  $\Sigma x^2 = 36$ 

 $b_0 = \overline{y} - b_1 \overline{x}$ 5 - (-0.181818)(2.5) = 5.45 $\Sigma y^2 = 120$ 

 $\Sigma xy = 48$ 

# Calculating the<br/>Regression EquationSlide 47Datax1135y2864

n = 4The estimated equation of the regression line is: $\Sigma x = 10$  $\widehat{\Sigma} y = 20$  $\Sigma y = 20$  $\widehat{Y} = 5.45 - 0.182x$  $\Sigma y^2 = 120$  $\Sigma xy = 48$ 



Given the sample data in Table 9-1, find the regression equation. (from pp.507-508)

Using the same procedure as in the previous example, we find that  $b_1 = 2.27$  and  $b_0 = -113$ . Hence, the estimated regression equation is:

$$\hat{y} = -113 + 2.27x$$



## Given the sample data in Table 9-1, find the regression equation.



## Predictions



In predicting a value of y based on some given value of x ...

- 1. If there is not a significant linear correlation, the best predicted y-value is  $\overline{y}$ .
- 2. If there is a significant linear correlation, the best predicted *y*-value is found by substituting the *x*-value into the regression equation.



(p.522)







We must consider whether there is a linear correlation that justifies the use of that equation. We do have a significant linear correlation (with r = 0.922).



Given the sample data in Table 9-1, we found that the regression equation is \$ = -113 + 2.27x. Given that x = 85, find the best predicted value of y, the number of manatee deaths from boats.

$$\overset{\wedge}{y} = -113 + 2.27x \\ -113 + 2.27(85) = 80.0$$

The predicted number of manatee deaths is 80.0. The actual number of manatee deaths in 2001 was 82, so the predicted value of 80.0 is quite close.

## **Guidelines for Using The Regression Equation (p.523)**



- 1. If there is no significant linear correlation, don't use the regression equation to make predictions.
- 2. When using the regression equation for predictions, stay within the scope of the available sample data.
- 3. A regression equation based on old data is not necessarily valid now.
- 4. Don't make predictions about a population that is different from the population from which the sample data was drawn.



- Marginal Change: The marginal change is the amount that a variable changes when the other variable changes by exactly one unit.
- Outlier: An outlier is a point lying far away from the other data points.
- Influential Points: An influential point strongly affects the graph of the regression line.

## Residuals and the Slide 57 Least-Squares Property Definitions (p.525)

#### Residual

for a sample of paired (x, y) data, the difference (y - y)between an observed sample y-value and the value of y, which is the value of y that is predicted by using the regression equation.

### Least-Squares Property

A straight line satisfies this property if the sum of the squares of the residuals is the smallest sum possible.

## Residuals and the Least-Squares Property



Chapter 9, Triola, Elementary Statistics, MATH 1342

## Section 9-4 Variation and Prediction Intervals

Created by Erin Hodgess, Houston, Texas





We consider different types of variation that can be used for two major applications:

1. To determine the proportion of the variation in y that can be explained by the linear relationship between x and y.

2. To construct interval estimates of predicted *y*-values. Such intervals are called prediction intervals.



**Total Deviation** The total deviation from the mean of the particular point (x, y) is the vertical distance  $y - \overline{y}$ , which is the distance between the point (x, y) and the horizontal line passing through the sample mean  $\overline{y}$ .

#### **Explained Deviation is**

the vertical distance  $\hat{y} - \overline{y}$ , which is the distance between the predicted *y*-value and the horizontal line passing through the sample mean  $\overline{y}$ .



#### **Unexplained Deviation is**

the vertical distance  $y - \hat{y}$ , which is the vertical distance between the point (x, y) and the regression line. (The distance  $y - \hat{y}$  is also called a *residual*, as defined in Section 9-3.).





#### Figure 9-10 Unexplained, Explained, and Total Deviation



#### (total deviation) = (explained deviation) + (unexplained deviation)

$$(y - \overline{y}) = (\hat{y} - \overline{y}) + (y - \hat{y})$$

#### (total variation) = (explained variation) + (unexplained variation)

$$\Sigma (y - \overline{y})^2 = \Sigma (\hat{y} - \overline{y})^2 + \Sigma (y - \hat{y})^2$$
  
Formula 9-4



## **Coefficient of determination** the amount of the variation in y that is explained by the regression line

= explained variation.

total variation

#### or

#### simply square *r* (determined by Formula 9-1, section 9-2)

## **Prediction Intervals**



## Definition

The standard error of estimate is a measure of the differences (or distances) between the observed sample y values and the predicted values  $\hat{y}$  that are obtained using the regression equation.





5



$$\frac{\sum y^2 - b_0 \sum y - b_1 \sum xy}{n - 2}$$
 Formula 9-



Given the sample data in Table 9-1, we found that the regression equation is  $\hat{y} = -113 + 2.27x$ . Find the standard error of estimate  $s_e$  for the boat/manatee data.





Given the sample data in Table 9-1, we found that the regression equation is  $\hat{y} = -113 + 2.27x$ . Find the standard error of estimate  $s_e$  for the boat/manatee data.

n = 10  $\Sigma y^2 = 33456$   $\Sigma y = 558$   $\Sigma xy = 42214$   $b_0 = -112.70989$  $b_1 = 2.27408$ 

$$s_{e} = 6.61234 = 6.61$$

Prediction Interval for an Individual y



$$\hat{y} - E < y < \hat{y} + E$$

#### where

$$E = t_{\alpha/2} s_{e} / \sqrt{1 + \frac{1}{n} + \frac{n(x_{0} - \bar{x})^{2}}{n(\Sigma x^{2}) - (\Sigma x)^{2}}}$$

 $x_0$  represents the given value of x $t_{\alpha/2}$  has n - 2 degrees of freedom



Given the sample data in Table 9-1, we found that the regression equation is  $\hat{y} = -113 + 2.27x$ . We have also found that when x = 85, the predicted number of manatee deaths is 80.0. Construct a 95% prediction interval given that x = 85.

$$E = t_{\alpha/2} s_{e} \sqrt{1 + \frac{1}{n} + \frac{n(x_{0} - \overline{x})^{2}}{n(\Sigma x^{2}) - (\Sigma x)^{2}}}$$
$$E = (2.306)(6.6123) \sqrt{1 + \frac{1}{10} + \frac{10(85 - 74)^{2}}{10(55289) - (741)^{2}}}$$

#### *E* = 18.1



Given the sample data in Table 9-1, we found that the regression equation is  $\hat{y} = -113 + 2.27x$ . We have also found that when x = 85, the predicted number of manatee deaths is 80.0. Construct a 95% prediction interval given that x = 85.

## $\hat{y} - E < y < \hat{y} + E$ 80.6 - 18.1 < y < 80.6 + 18.1 62.5 < y < 98.7
# Section 9-5 Multiple Regression

Created by Erin Hodgess, Houston, Texas



#### **Multiple Regression**



#### Definition Multiple Regression Equation

A linear relationship between a dependent variable y and two or more independent variables  $(x_1, x_2, x_3, \dots, x_k)$ 

#### $\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \ldots + b_k x_k$

# Notation



 $\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_k x_k$ (General form of the estimated multiple regression equation)

- n = sample size
- k = number of independent variables
- $\hat{y}$  = predicted value of the dependent variable y
- $x_1, x_2, x_3 \dots, x_k$  are the independent variables

# Notation



 $\beta_0$  = the y-intercept, or the value of y when all of the predictor variables are 0  $b_0$  = estimate of  $\beta_0$  based on the sample data  $\beta_1, \beta_2, \beta_3, \ldots, \beta_k$  are the coefficients of the independent variables  $x_1, x_2, x_3, \ldots, x_k$  $b_1, b_2, b_3, \ldots, b_k$  are the sample estimates the coefficients  $\beta_1, \beta_2, \beta_3, \ldots, \beta_k$ of

#### Assumption



#### Use a statistical software package such as









For reasons of safety, a study of bears involved the collection of various measurements that were taken after the bears were anesthetized. Using the data in Table 9-3, find the multiple regression equation in which the dependent variable is weight and the independent variables are head length and total overall length.



Table 9-3	Data fro	Data from Anesthetized Male Bears									
Variable	Minitab Column	Name	Sample Data								
У	C1	WEIGHT	80	344	416	348	262	360	332	34	
<i>x</i> <sub>2</sub>	C2	AGE	19	55	81	115	56	51	68	8	
<i>x</i> <sub>3</sub>	C3	HEADLEN	11.0	16.5	15.5	17.0	15.0	13.5	16.0	9.0	
<i>x</i> <sub>4</sub>	C4	HEADWDTH	5.5	9.0	8.0	10.0	7.5	8.0	9.0	4.5	
<i>x</i> <sub>5</sub>	C5	NECK	16.0	28.0	31.0	31.5	26.5	27.0	29.0	13.0	
<i>x</i> <sub>6</sub>	C6	LENGTH	53.0	67.5	72.0	72.0	73.5	68.5	73.0	37.0	
<i>x</i> <sub>7</sub>	C7	CHEST	26	45	54	49	41	49	44	19	



Minitab											
The regression equation is Multiple											
WEIGHT = $-374 + 18.8$ HEADLEN + 5.87 LENGTH $\leftarrow 1$ regression											
D	0	C	<b>C</b> + 1		equation						
Predictor	Coe	Ξ	Stdev	t-ratio	Р						
Constant	-374.	3	134.1	-2.79	0.038						
HEADLEN	18.8	32	23.15	0.81	0.453						
LENGTH	5.87	′ 5	5.065	1.16	0.299						
s = 68.56	R-sq = 82	.8% R-	sq(adj) =	75.9%							
Analysis of Variance $R^2 = 0.828$ ② Adjusted $R^2 = 0.759$											
SOURCE	DF	SS	MS	F	р						
Regression	2	113142	56571	12.03	0.012						
Error	5	23506	4701		1						
Total	7	③ Overall significance of multiple regression equation									



The regression equation is:

WEIGHT = -374 + 18.8 HEADLEN + 5.87 LENGTH  $y = -374 + 18.8x_3 + 5.87x_6$ 

# Adjusted R<sup>2</sup>



# Definitions

The multiple coefficient of determination is a measure of how well the multiple regression equation fits the sample data.

The Adjusted coefficient of determination R<sup>2</sup> is modified to account for the number of variables and the sample size.

# Adjusted R<sup>2</sup>



# Adjusted $R^2 = 1 - \frac{(n-1)}{[n-(k+1)]}(1-R^2)$

Formula 9-6

#### where n =sample size k =number of independent (x) variables

#### Finding the Best Multiple Regression Equation



- 1. Use common sense and practical considerations to include or exclude variables.
- 2. Instead of including almost every available variable, include relatively few independent (*x*) variables, weeding out independent variables that don't have an effect on the dependent variable.
- 3. Select an equation having a value of adjusted  $R^2$  with this property: If an additional independent variable is included, the value of adjusted  $R^2$  does not increase by a substantial amount.
- 4. For a given number of independent (x) variables, select the equation with the largest value of adjusted  $R^2$ .
- 5. Select an equation having overall significance, as determined by the *P*-value in the computer display.



Created by Erin Hodgess, Houston, Texas



# Definition



## **Mathematical Model**

A mathematical model is a mathematical function that 'fits' or describes real-world data.

## **TI-83 Generic Models**



- Linear: Quadratic: Logarithmic: Exponential: Power: Logistic:
  - y = a + bx $y = ax^2 + bx + c$  $y = a + b \ln x$  $y = ab^x$  $y = ax^b$  $y = \frac{c}{1 + ae^{-bx}}$











# Logarithmic: $y = 1 + 2\ln x$













#### Development of a Good Mathematics Model



- Look for a Pattern in the Graph: Examine the graph of the plotted points and compare the basic pattern to the known generic graphs.
- Find and Compare Values of R<sup>2</sup>: Select functions that result in larger values of R<sup>2</sup>, because such larger values correspond to functions that better fit the observed points.
- Think: Use common sense. Don't use a model that lead to predicted values known to be totally unrealistic.