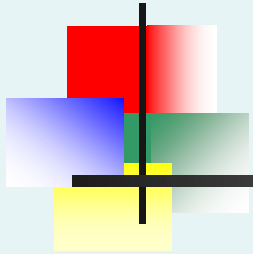


*Statistics for Managers Using
Microsoft Excel*
7th Edition



Chapter 13

Simple Linear Regression



Learning Objectives

In this chapter, you learn:

- How to use regression analysis to predict the value of a dependent variable based on an independent variable
- The meaning of the regression coefficients b_0 and b_1
- How to evaluate the assumptions of regression analysis and know what to do if the assumptions are violated
- To make inferences about the slope and correlation coefficient
- To estimate mean values and predict individual values



Correlation vs. Regression

DCOVA A

- A **scatter plot** can be used to show the relationship between two variables
- **Correlation** analysis is used to measure the strength of the association (linear relationship) between two variables
 - Correlation is only concerned with strength of the relationship
 - No causal effect is implied with correlation
 - Scatter plots were first presented in Ch. 2
 - Correlation was first presented in Ch. 3



Introduction to Regression Analysis

DCOVA A

- **Regression analysis** is used to:
 - Predict the value of a dependent variable based on the value of at least one independent variable
 - Explain the impact of changes in an independent variable on the dependent variable

Dependent variable: the variable we wish to predict or explain

Independent variable: the variable used to predict or explain the dependent variable



Simple Linear Regression Model

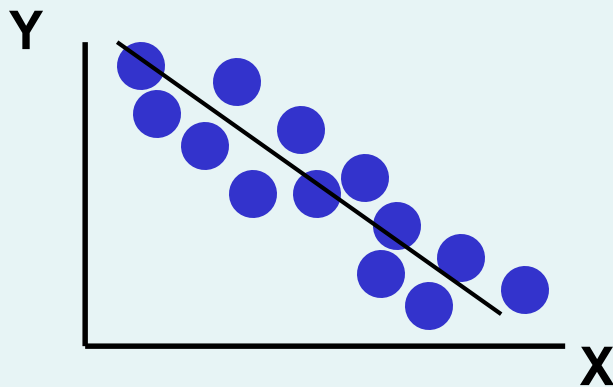
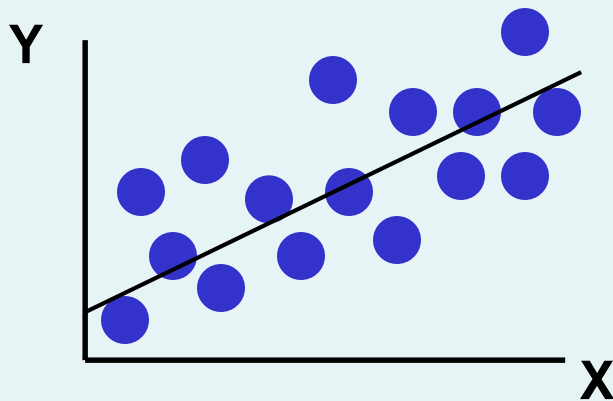
DCOVA A

- Only **one** independent variable, X
- Relationship between X and Y is described by a linear function
- Changes in Y are assumed to be related to changes in X

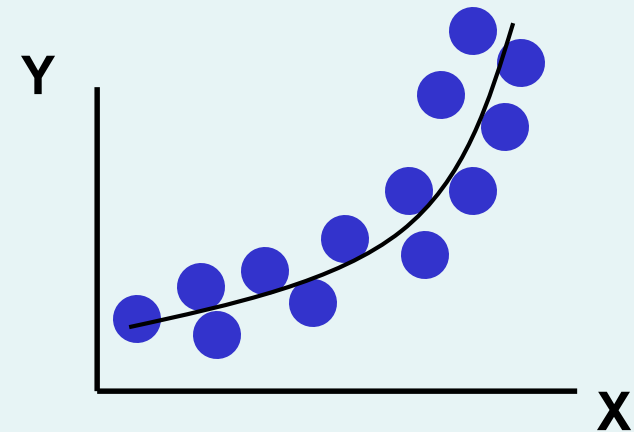
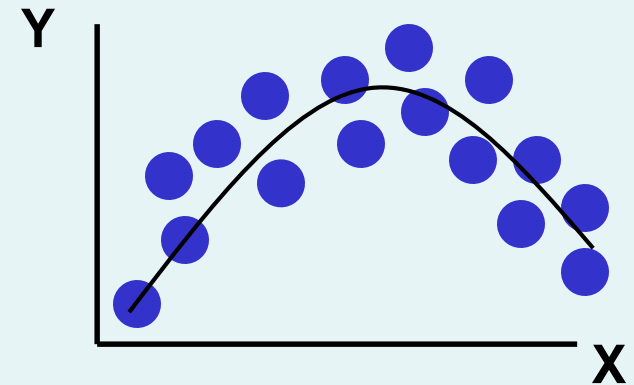
Types of Relationships

DCOVA

Linear relationships



Curvilinear relationships

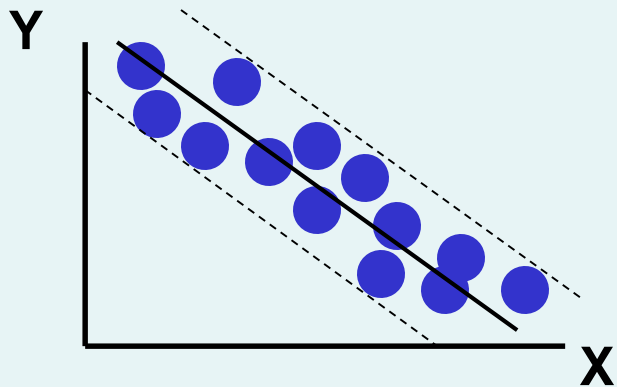
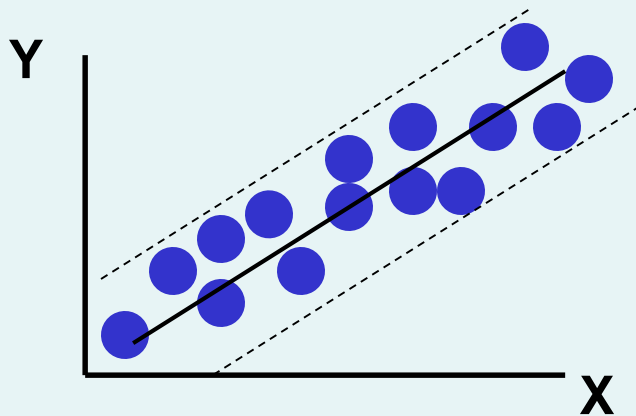


Types of Relationships

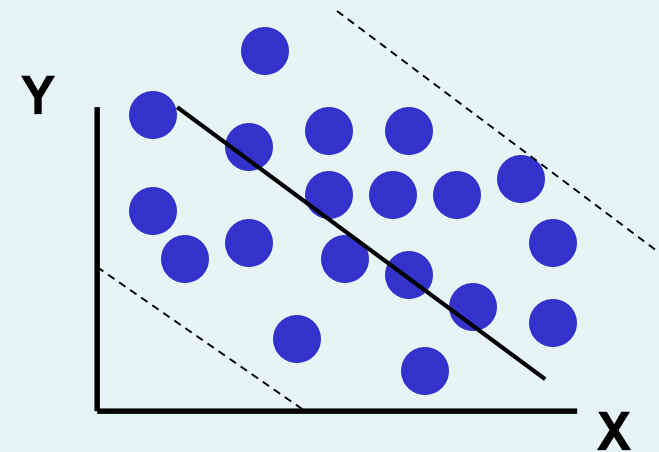
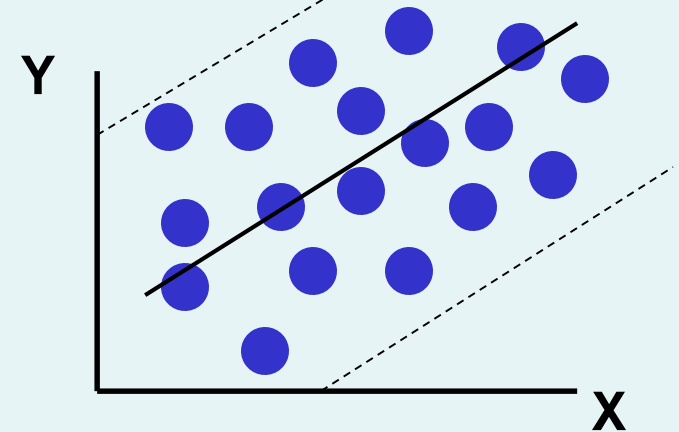
DCOVA

(continued)

Strong relationships



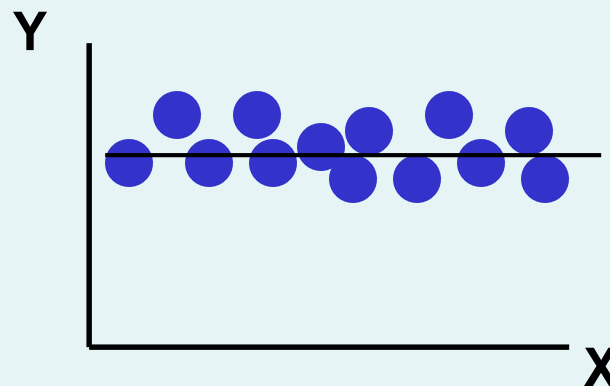
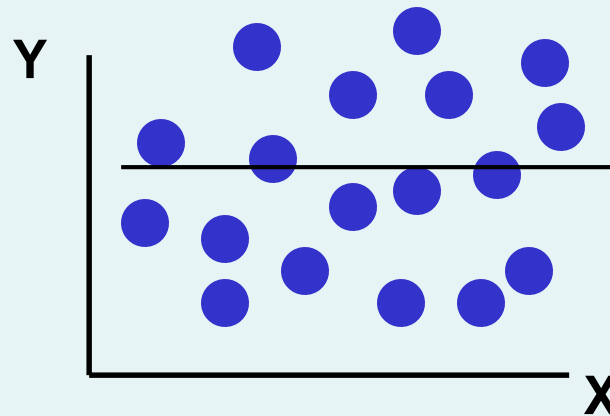
Weak relationships



Types of Relationships

DCOVA
(continued)

No relationship



Simple Linear Regression Model

DCOVA A

The diagram illustrates the Simple Linear Regression Model equation: $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$. The equation is enclosed in a light orange box. Labels with arrows point to each term: Y_i is labeled 'Dependent Variable', β_0 is 'Population Y intercept', β_1 is 'Population Slope Coefficient', X_i is 'Independent Variable', and ϵ_i is 'Random Error term'. Below the equation, two blue brackets group the terms: the first bracket under $\beta_0 + \beta_1 X_i$ is labeled 'Linear component', and the second bracket under ϵ_i is labeled 'Random Error component'.

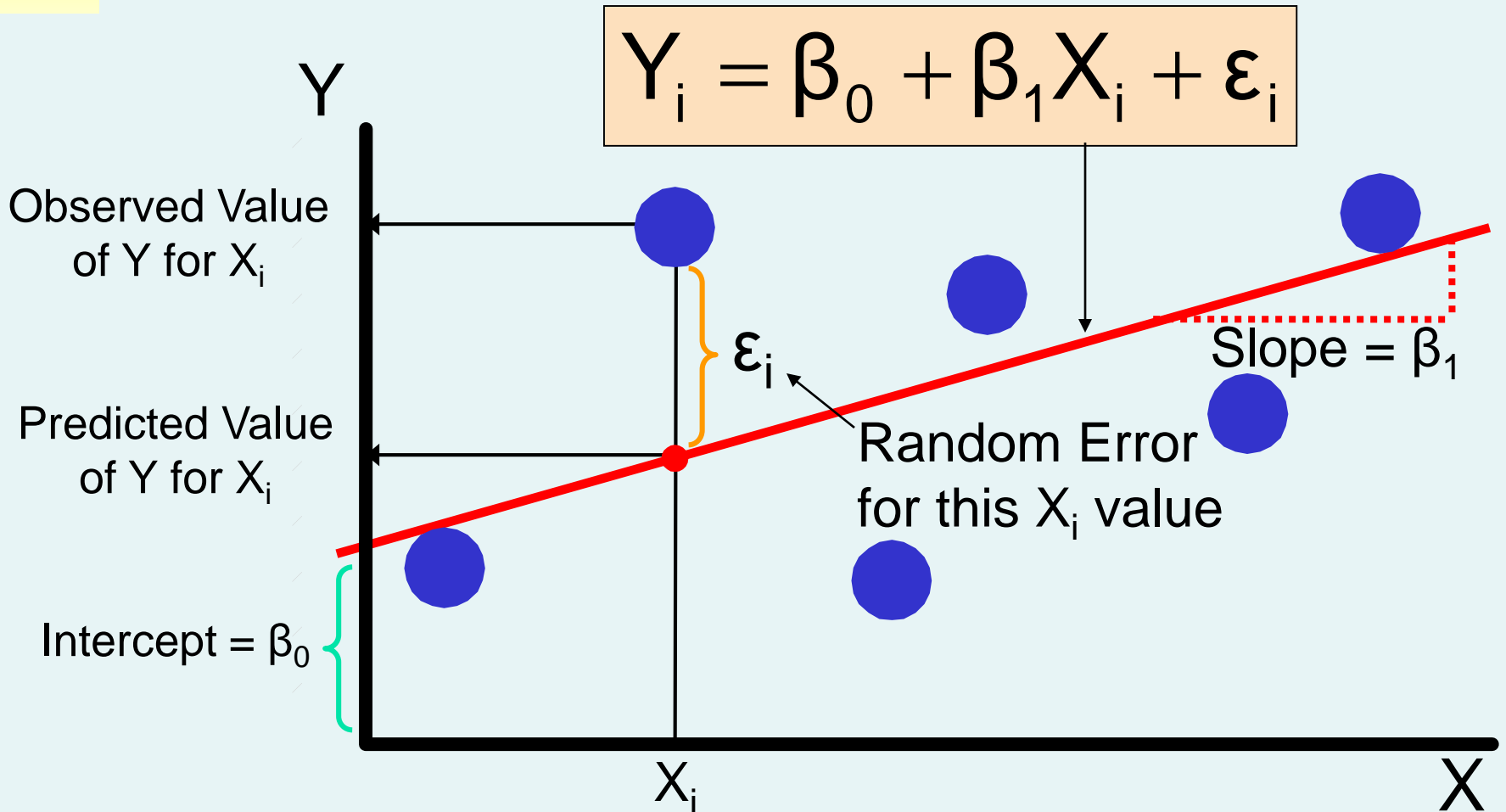
$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

Labels and components:

- Dependent Variable: Y_i
- Population Y intercept: β_0
- Population Slope Coefficient: β_1
- Independent Variable: X_i
- Random Error term: ϵ_i
- Linear component: $\beta_0 + \beta_1 X_i$
- Random Error component: ϵ_i

Simple Linear Regression Model

DCOVA
(continued)



Simple Linear Regression Equation (Prediction Line)

DCOVA

The simple linear regression equation provides an **estimate** of the population regression line

Estimated
(or predicted)
Y value for
observation i

Estimate of
the regression
intercept

Estimate of the
regression slope

Value of X for
observation i

$$\hat{Y}_i = b_0 + b_1 X_i$$



The Least Squares Method

DCOVA A

b_0 and b_1 are obtained by finding the values of that minimize the sum of the squared differences between Y and \hat{Y} :

$$\min \sum (Y_i - \hat{Y}_i)^2 = \min \sum (Y_i - (b_0 + b_1 X_i))^2$$

Finding the Least Squares Equation

- The coefficients b_0 and b_1 , and other regression results in this chapter, will be found using Excel

Formulas are shown in the text for those who are interested



Interpretation of the Slope and the Intercept

DCOVA A

- b_0 is the estimated average value of Y when the value of X is zero
- b_1 is the estimated change in the average value of Y as a result of a one-unit increase in X

Simple Linear Regression Example

- A real estate agent wishes to examine the relationship between the selling price of a home and its size (measured in square feet)
- A random sample of 10 houses is selected
 - Dependent variable (Y) = house price in \$1000s
 - Independent variable (X) = square feet



Simple Linear Regression

Example: Data

DCOVA

House Price in \$1000s (Y)	Square Feet (X)
245	1400
312	1600
279	1700
308	1875
199	1100
219	1550
405	2350
324	2450
319	1425
255	1700

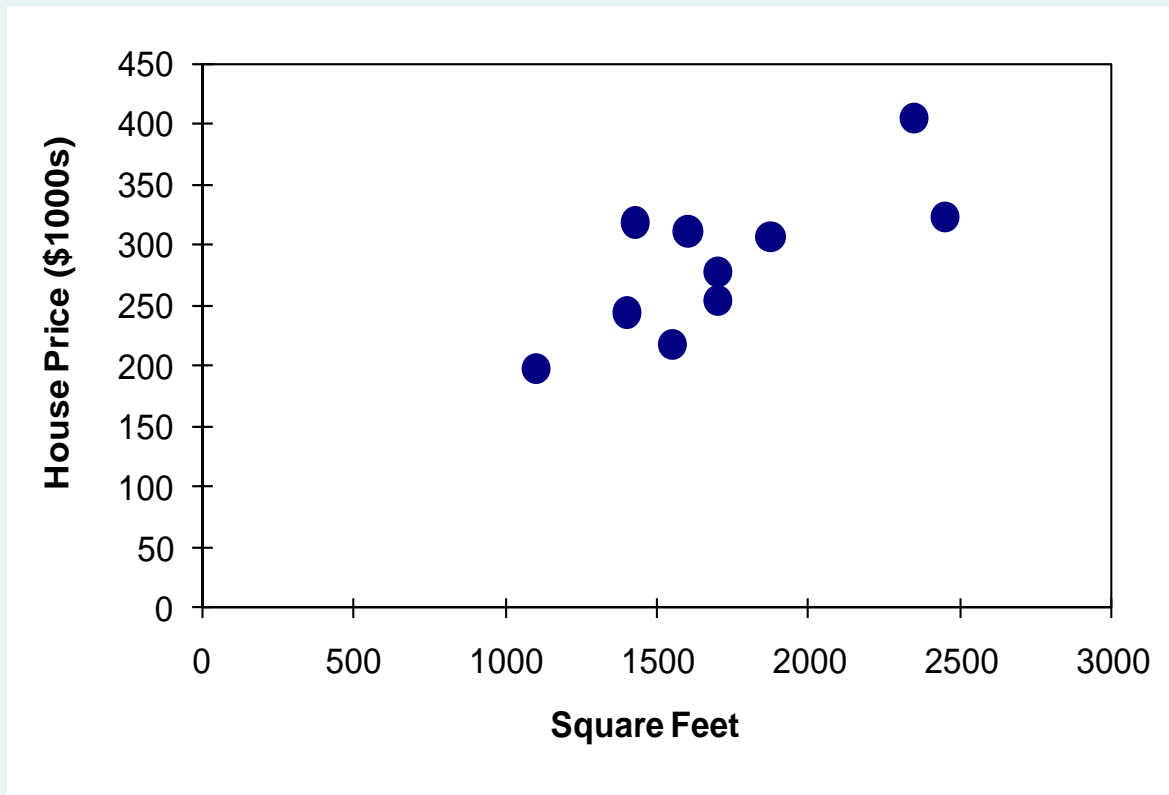


Simple Linear Regression

Example: Scatter Plot

DCOVA

House price model: Scatter Plot



Simple Linear Regression Example: Using Excel Data Analysis Function

DCOVA

1. Choose Data

2. Choose Data Analysis

3. Choose Regression

The screenshot shows the Microsoft Excel interface with the Data Analysis toolpak installed. The Data tab is active, and the Data Analysis button is visible in the ribbon. The Data Analysis dialog box is open, showing a list of analysis tools. The 'Regression' option is selected. The spreadsheet data is as follows:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R
1	House Price	Square Feet																
2		245	1400															
3		312	1600															
4		279	1700															
5		308	1875															
6		199	1100															
7		219	1550															
8		405	2350															
9		324	2450															
10		319	1425															
11		255	1700															
12																		



Simple Linear Regression Example: Using Excel Data Analysis Function

(continued)

Enter Y range and X range and desired options

DCOVA A

	A	B	C	D	E	F	G	H	I
1	House Price	Square Feet							
2	245	1400							
3	312	1600							
4	279	1700							
5	308	1875							
6	199	1100							
7	219	1550							
8	405	2350							
9	324	2450							
10	319	1425							
11	255	1700							
12									
13									
14									
15									
16									
17									
18									
19									
20									

Regression

Input

Input Y Range:

Input X Range:

Labels Constant is Zero

Confidence Level: %

Output options

Output Range:

New Worksheet Ply:

New Workbook

Residuals

Residuals Residual Plots

Standardized Residuals Line Fit Plots

Normal Probability

Normal Probability Plots

OK Cancel Help



Simple Linear Regression

Example: Using PHStat

Add-Ins: PHStat: Regression: Simple Linear Regression

The screenshot displays the PHStat add-in menu in Microsoft Excel. The 'Regression' option is selected, and the 'Simple Linear Regression...' option is highlighted. A red arrow points from this option to the 'Simple Linear Regression' dialog box. The dialog box shows the Y Variable Cell Range as 'Sheet1!\$A\$2:\$A\$11' and the X Variable Cell Range as 'Sheet1!\$B\$2:\$B\$11'. The 'Regression Statistics Table' and 'ANOVA and Coefficients Table' options are checked. The 'Confidence level for regression coefficients' is set to 95%.

	A	B	C	D	F	G	H
1	House Price	Square Feet					
2	245	1400					
3	312	1600					
4	279	1700					
5	308	1875					
6	199	1100					
7	219	1550					
8	405	2350					
9	324	2450					
10	319	1425					
11	255	1700					
12							
13							
14							
15							
16							
17							
18							
19							
20							
21							



Simple Linear Regression Example: Excel Output

Regression Statistics

Multiple R	0.76211
R Square	0.58082
Adjusted R Square	0.52842
Standard Error	41.33032
Observations	10

The regression equation is:

$$\widehat{\text{house price}} = 98.24833 + 0.10977 (\text{square feet})$$

ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	1	18934.9348	18934.9348	11.0848	0.01039
Residual	8	13665.5652	1708.1957		
Total	9	32600.5000			

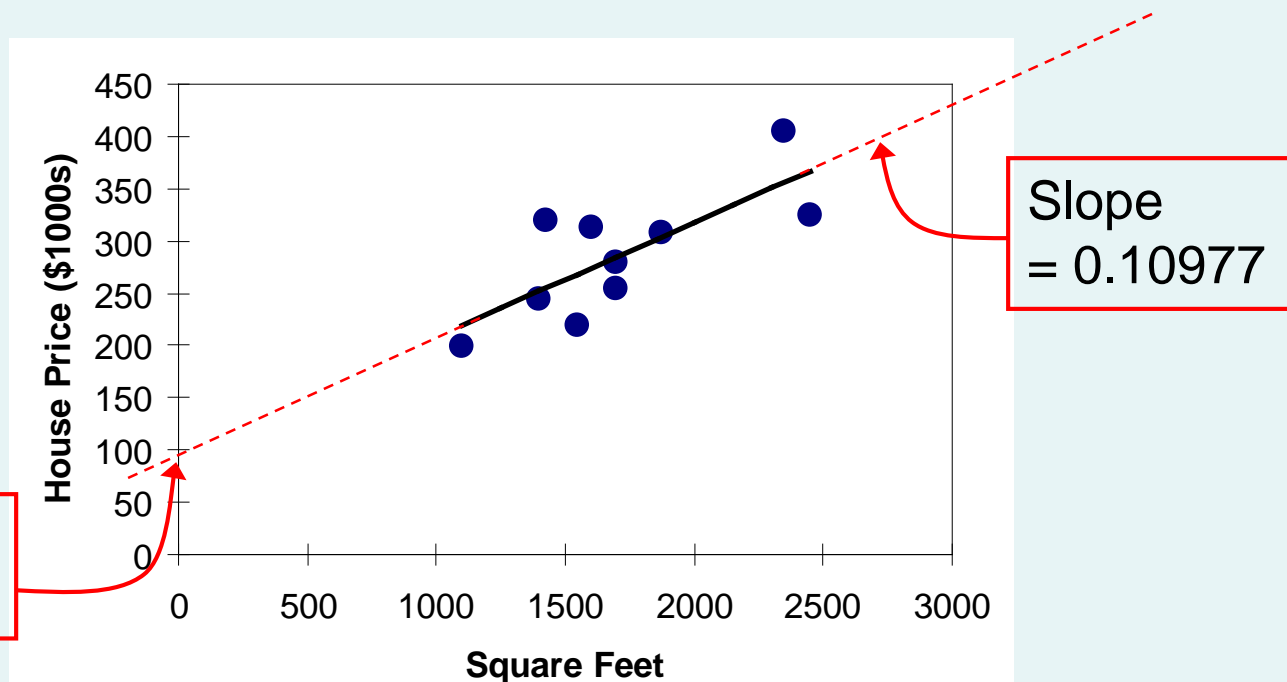
	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	98.24833	58.03348	1.69296	0.12892	-35.57720	232.07386
Square Feet	0.10977	0.03297	3.32938	0.01039	0.03374	0.18580



Simple Linear Regression Example: Graphical Representation

DCOVA

House price model: Scatter Plot and Prediction Line



$$\widehat{\text{house price}} = 98.24833 + 0.10977 (\text{square feet})$$

Simple Linear Regression

Example: Interpretation of b_0

DCOVA

$$\widehat{\text{house price}} = 98.24833 + 0.10977 (\text{square feet})$$

- b_0 is the estimated average value of Y when the value of X is zero (if $X = 0$ is in the range of observed X values)
- Because a house cannot have a square footage of 0, b_0 has no practical application



Simple Linear Regression

Example: Interpreting b_1

DCOVA

$$\widehat{\text{house price}} = 98.24833 + 0.10977(\text{square feet})$$

- b_1 estimates the change in the average value of Y as a result of a one-unit increase in X
 - Here, $b_1 = 0.10977$ tells us that the mean value of a house increases by $.10977(\$1000) = \109.77 , on average, for each additional one square foot of size



Simple Linear Regression

Example: Making Predictions

DCOVA

Predict the price for a house with 2000 square feet:

$$\begin{aligned}\widehat{\text{house price}} &= 98.25 + 0.1098 (\text{sq.ft.}) \\ &= 98.25 + 0.1098(2000) \\ &= 317.85\end{aligned}$$

The predicted price for a house with 2000 square feet is $317.85(\$1,000\text{s}) = \$317,850$



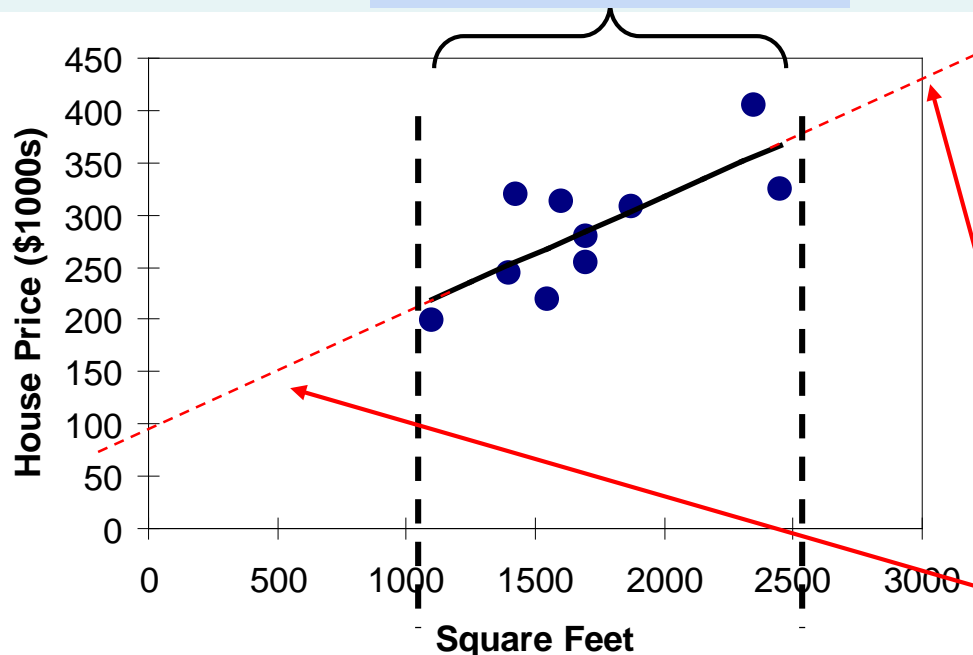
Simple Linear Regression

Example: Making Predictions

DCOVA

- When using a regression model for prediction, only predict within the relevant range of data

Relevant range for interpolation



Do not try to extrapolate beyond the range of observed X's

Measures of Variation

- Total variation is made up of two parts:

$$SST = SSR + SSE$$

Total Sum of Squares

Regression Sum of Squares

Error Sum of Squares

$$SST = \sum (Y_i - \bar{Y})^2$$

$$SSR = \sum (\hat{Y}_i - \bar{Y})^2$$

$$SSE = \sum (Y_i - \hat{Y}_i)^2$$

where:

\bar{Y} = Mean value of the dependent variable

Y_i = Observed value of the dependent variable

\hat{Y}_i = Predicted value of Y for the given X_i value

Measures of Variation

(continued)

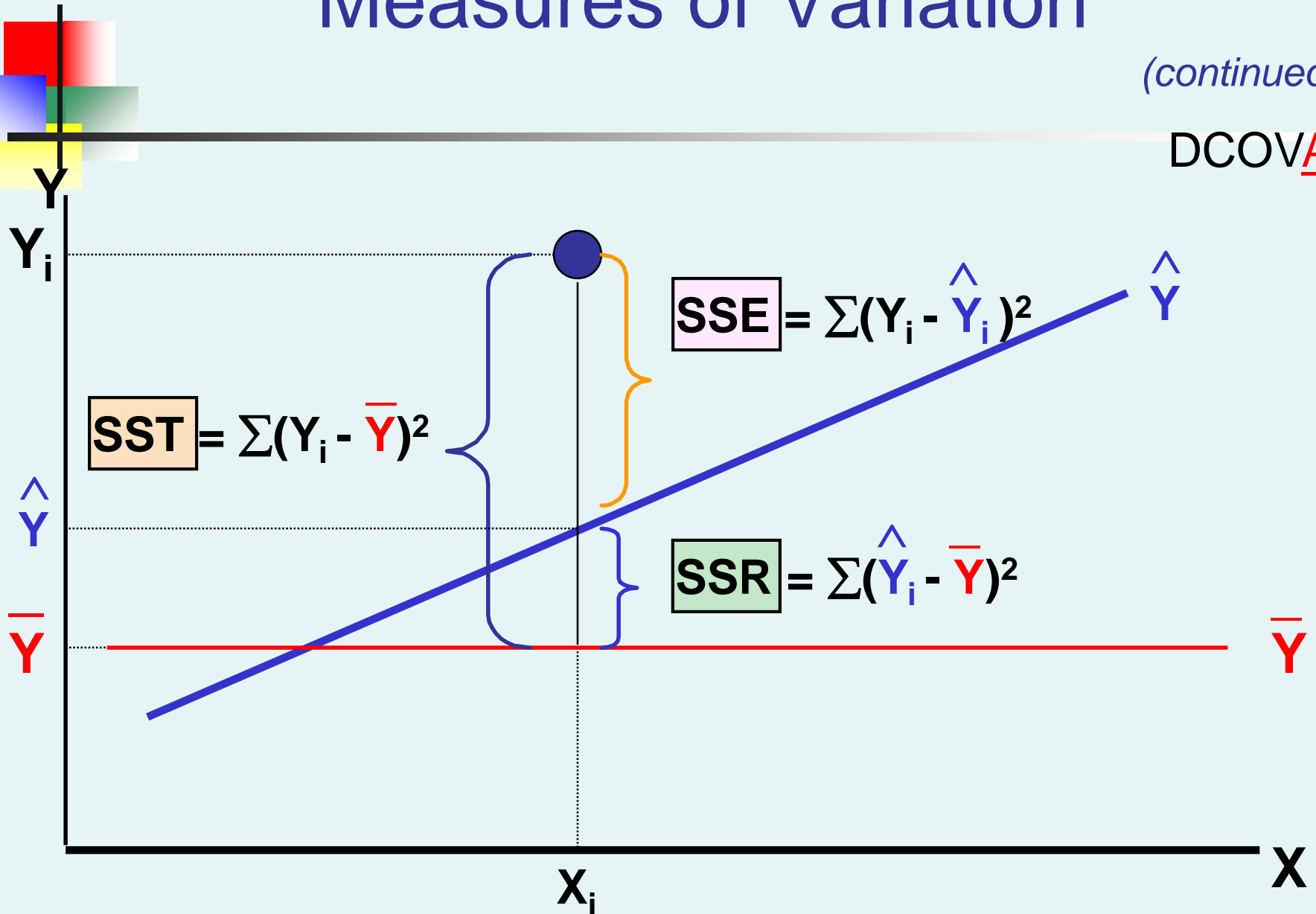
DCOVA

- SST = total sum of squares (Total Variation)
 - Measures the variation of the Y_i values around their mean \bar{Y}
- SSR = regression sum of squares (Explained Variation)
 - Variation attributable to the relationship between X and Y
- SSE = error sum of squares (Unexplained Variation)
 - Variation in Y attributable to factors other than X

Measures of Variation

(continued)

DCOVA





Coefficient of Determination, r^2

DCOVA

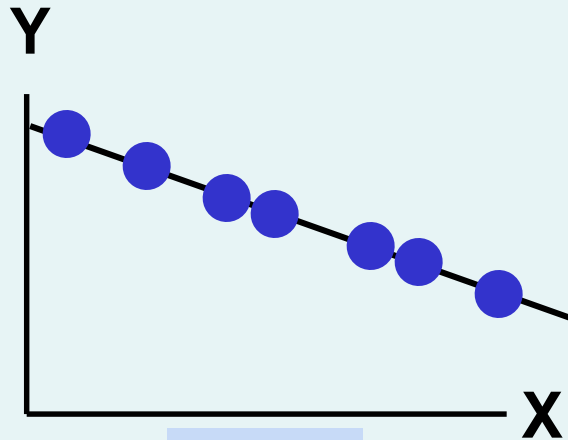
- The **coefficient of determination** is the portion of the total variation in the dependent variable that is explained by variation in the independent variable
- The coefficient of determination is also called **r-squared** and is denoted as r^2

$$r^2 = \frac{SSR}{SST} = \frac{\text{regression sum of squares}}{\text{total sum of squares}}$$

note:

$$0 \leq r^2 \leq 1$$

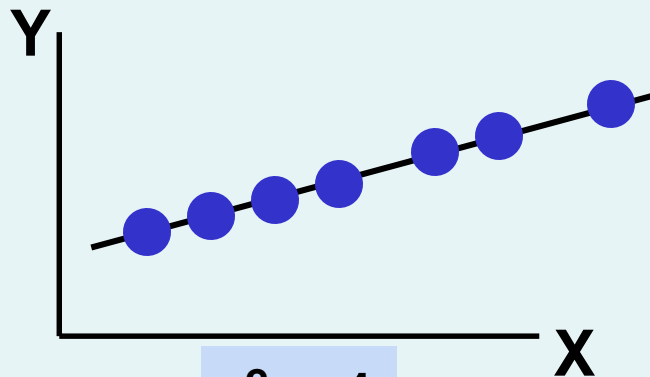
Examples of Approximate r^2 Values



$$r^2 = 1$$

$$r^2 = 1$$

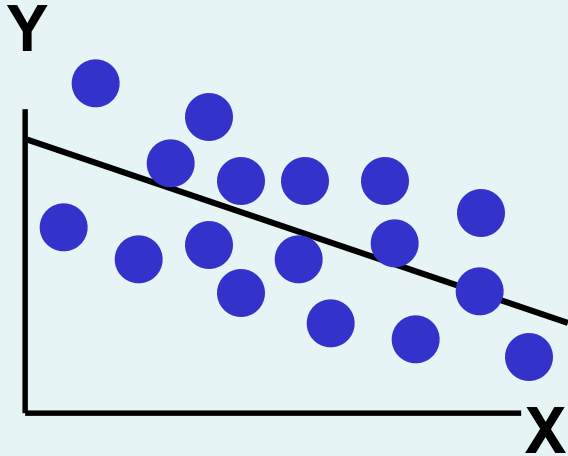
Perfect linear relationship between X and Y:



$$r^2 = 1$$

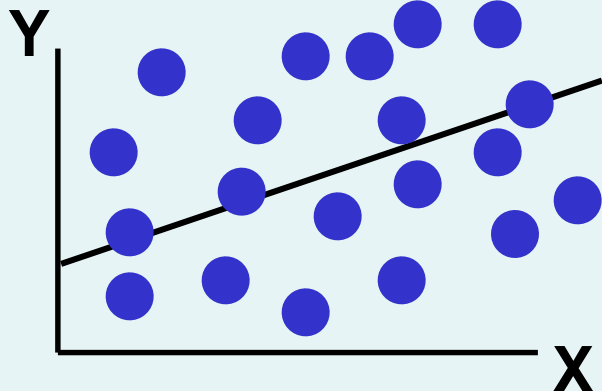
100% of the variation in Y is explained by variation in X

Examples of Approximate r^2 Values



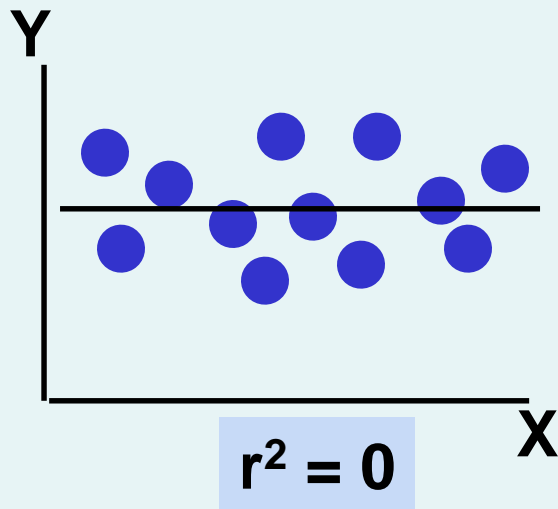
$$0 < r^2 < 1$$

Weaker linear relationships between X and Y:



Some but not all of the variation in Y is explained by variation in X

Examples of Approximate r^2 Values



$$r^2 = 0$$

No linear relationship between X and Y:

The value of Y does not depend on X. (None of the variation in Y is explained by variation in X)

Simple Linear Regression Example: Coefficient of Determination, r^2 in Excel

DCOVA

Regression Statistics

Multiple R	0.76211
R Square	0.58082
Adjusted R Square	0.52842
Standard Error	41.33032
Observations	10

$$r^2 = \frac{SSR}{SST} = \frac{18934.9348}{32600.5000} = 0.58082$$

58.08% of the variation in house prices is explained by variation in square feet

ANOVA

	<i>df</i>	SS	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	1	18934.9348	18934.9348	11.0848	0.01039
Residual	8	13665.5652	1708.1957		
Total	9	32600.5000			

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	98.24833	58.03348	1.69296	0.12892	-35.57720	232.07386
Square Feet	0.10977	0.03297	3.32938	0.01039	0.03374	0.18580



Standard Error of Estimate

DCOVA

- The standard deviation of the variation of observations around the regression line is estimated by

$$S_{YX} = \sqrt{\frac{SSE}{n-2}} = \sqrt{\frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-2}}$$

Where

SSE = error sum of squares

n = sample size

Simple Linear Regression Example: Standard Error of Estimate in Excel

DCOVA

Regression Statistics

Multiple R	0.76211
R Square	0.58082
Adjusted R Square	0.52842
Standard Error	41.33032
Observations	10

$$S_{YX} = 41.33032$$

ANOVA

	<i>df</i>	<i>SS</i>	<i>MS</i>	<i>F</i>	<i>Significance F</i>
Regression	1	18934.9348	18934.9348	11.0848	0.01039
Residual	8	13665.5652	1708.1957		
Total	9	32600.5000			

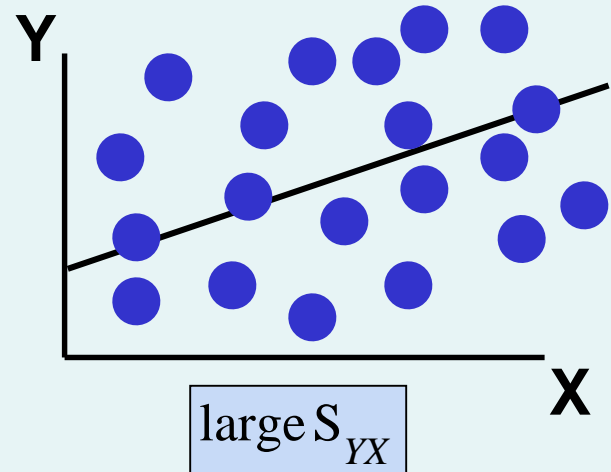
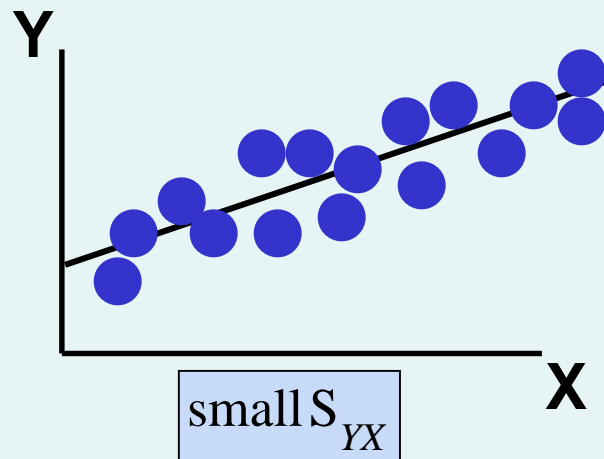
	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	98.24833	58.03348	1.69296	0.12892	-35.57720	232.07386
Square Feet	0.10977	0.03297	3.32938	0.01039	0.03374	0.18580



Comparing Standard Errors

DCOVA

S_{YX} is a measure of the variation of observed Y values from the regression line



The magnitude of S_{YX} should always be judged relative to the size of the Y values in the sample data

i.e., $S_{YX} = \$41.33K$ is moderately small relative to house prices in the \$200K - \$400K range

Assumptions of Regression

L.I.N.E

DCOVA

- Linearity
 - The relationship between X and Y is linear
- Independence of Errors
 - Error values are statistically independent
- Normality of Error
 - Error values are normally distributed for any given value of X
- Equal Variance (also called homoscedasticity)
 - The probability distribution of the errors has constant variance

Residual Analysis

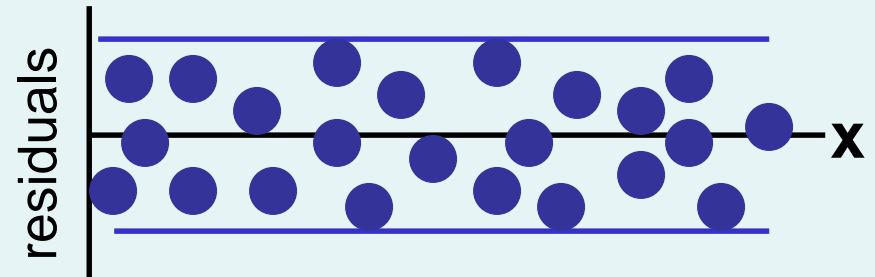
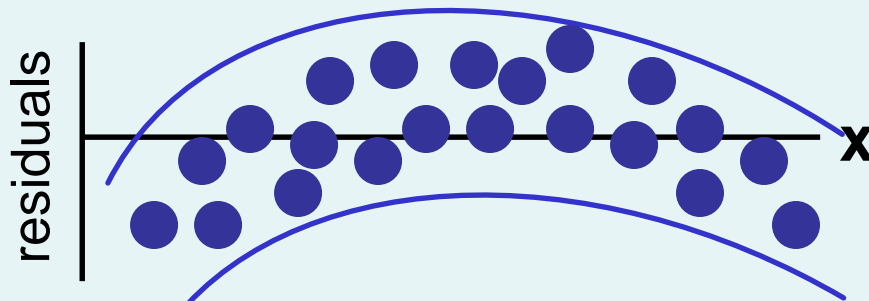
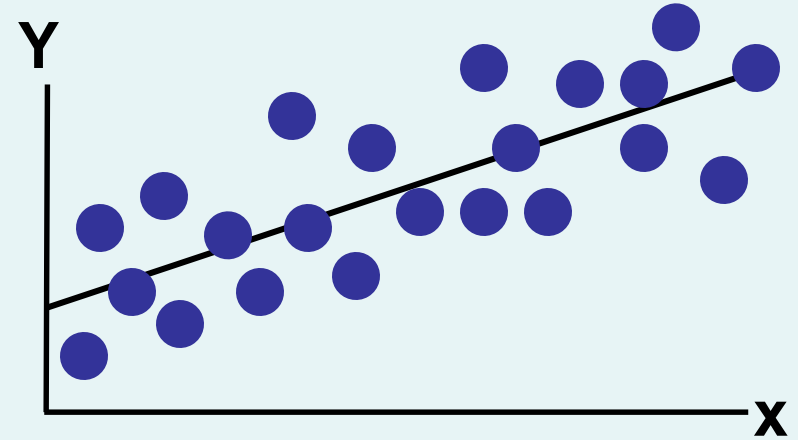
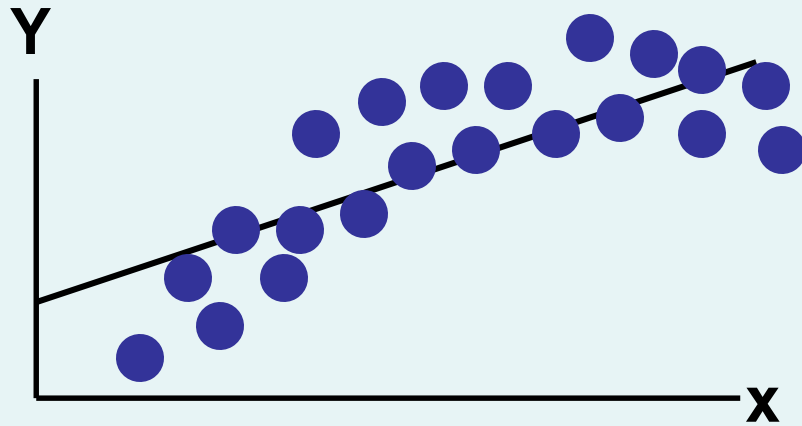
DCOVA

$$e_i = Y_i - \hat{Y}_i$$

- The residual for observation i , e_i , is the difference between its observed and predicted value
- Check the assumptions of regression by examining the residuals
 - Examine for linearity assumption
 - Evaluate independence assumption
 - Evaluate normal distribution assumption
 - Examine for constant variance for all levels of X (homoscedasticity)
- Graphical Analysis of Residuals
 - Can plot residuals vs. X

Residual Analysis for Linearity

DCOVA



Not Linear

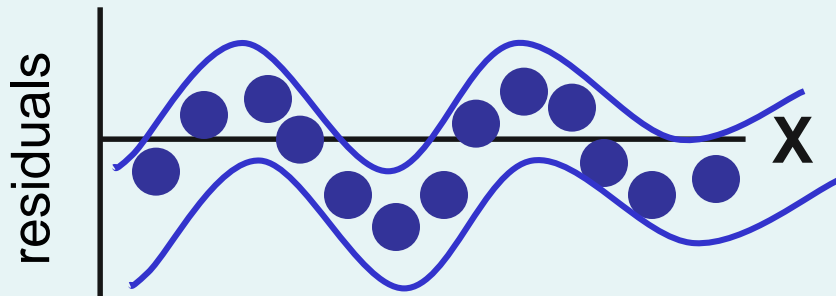
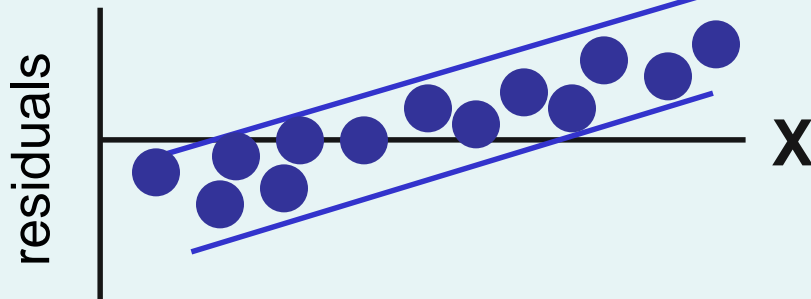


Linear

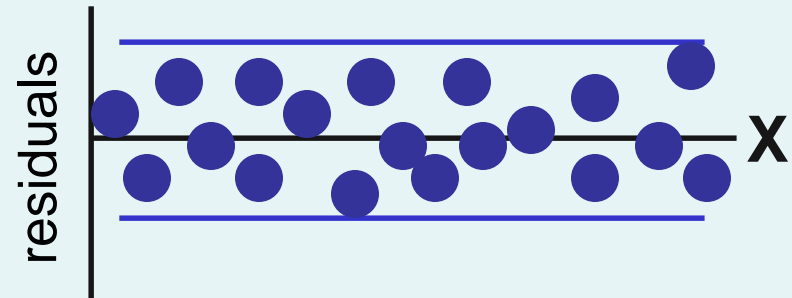
Residual Analysis for Independence



Not Independent



Independent





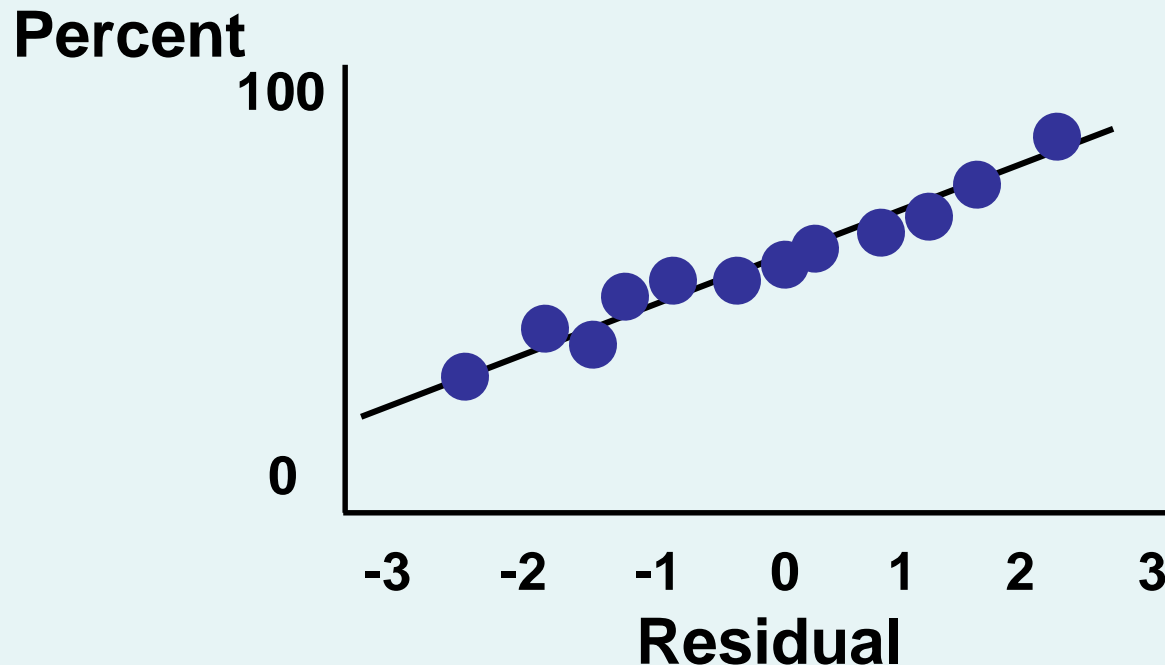
Checking for Normality

DCOVA

- Examine the Stem-and-Leaf Display of the Residuals
- Examine the Boxplot of the Residuals
- Examine the Histogram of the Residuals
- Construct a Normal Probability Plot of the Residuals

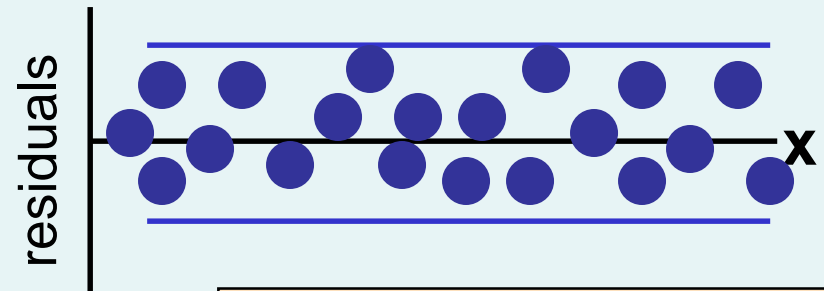
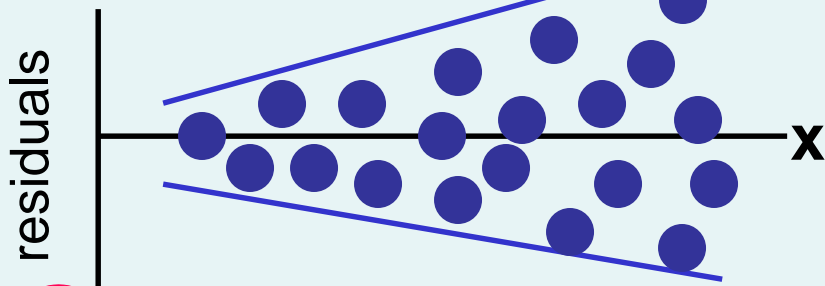
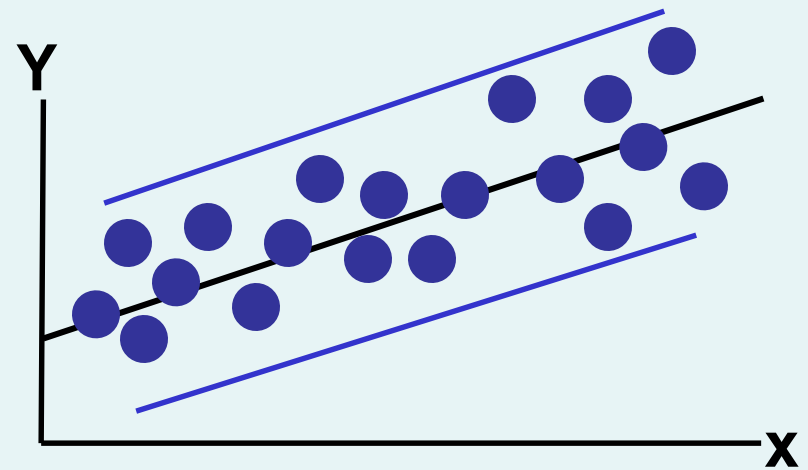
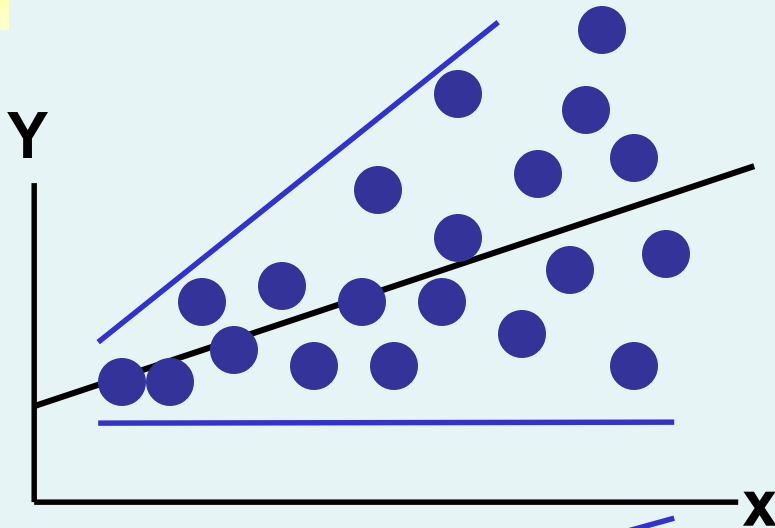
Residual Analysis for Normality

When using a normal probability plot, normal errors will approximately display in a straight line



Residual Analysis for Equal Variance

DCOVA



Non-constant variance



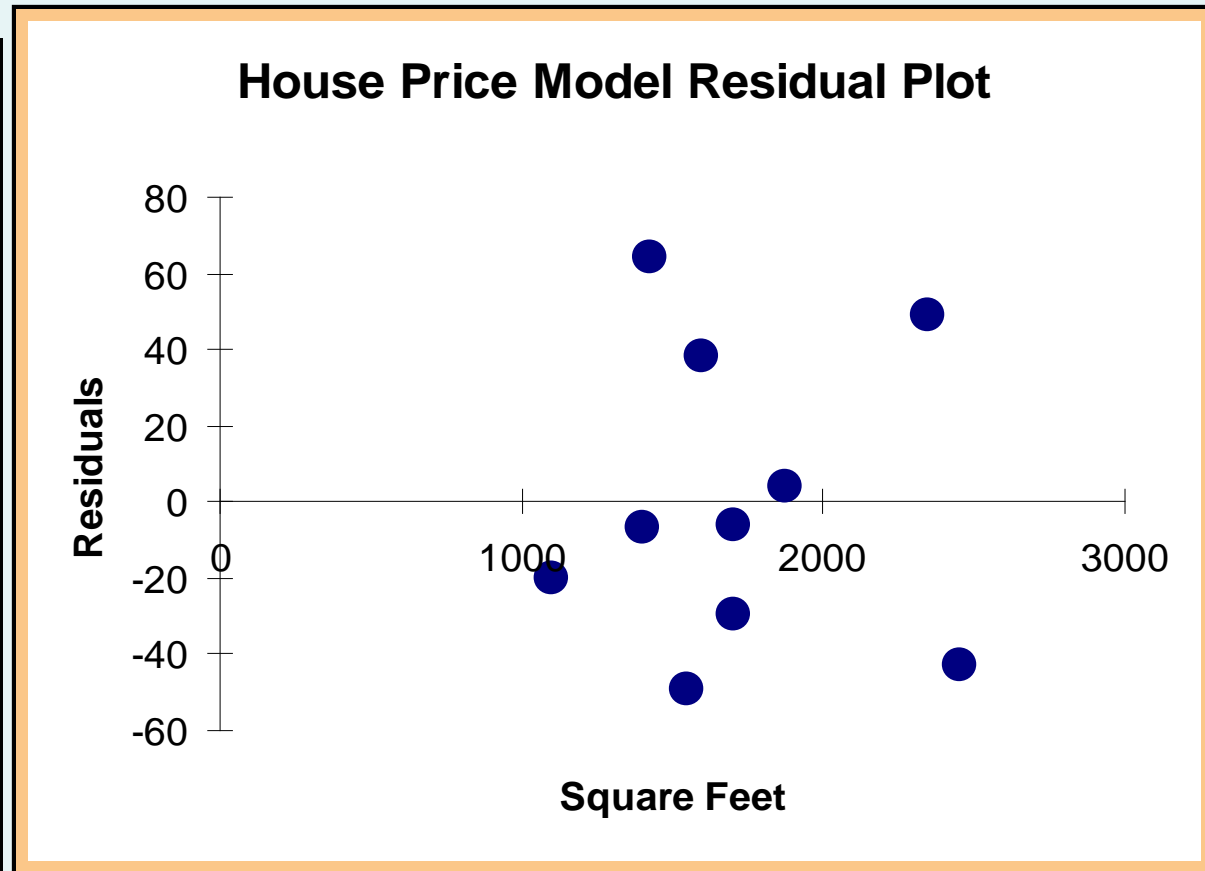
Constant variance

Simple Linear Regression

Example: Excel Residual Output

DCOVA

RESIDUAL OUTPUT		
	<i>Predicted House Price</i>	<i>Residuals</i>
1	251.92316	-6.923162
2	273.87671	38.12329
3	284.85348	-5.853484
4	304.06284	3.937162
5	218.99284	-19.99284
6	268.38832	-49.38832
7	356.20251	48.79749
8	367.17929	-43.17929
9	254.6674	64.33264
10	284.85348	-29.85348



Does not appear to violate
any regression assumptions



Measuring Autocorrelation: The Durbin-Watson Statistic

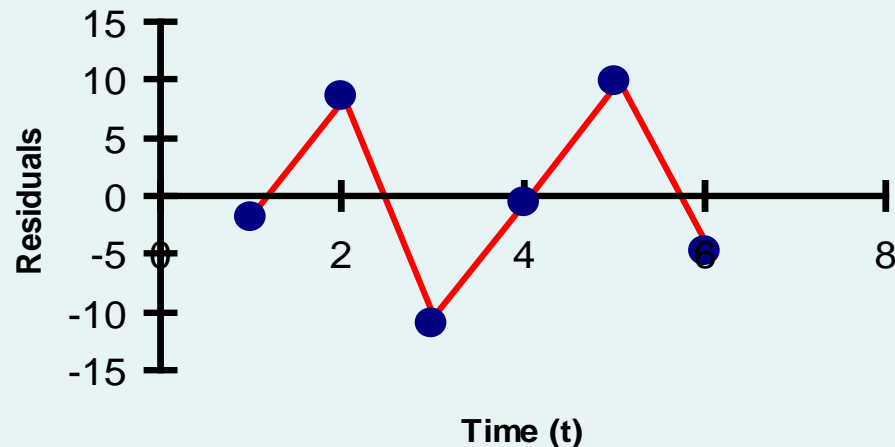
DCOVA A

- Used when data are **collected over time** to detect if autocorrelation is present
- Autocorrelation exists if residuals in one time period are related to residuals in another period

Autocorrelation

- Autocorrelation is correlation of the errors (residuals) over time

Time (t) Residual Plot



- Here, residuals show a cyclic pattern, not random. Cyclical patterns are a sign of positive autocorrelation

- Violates the regression assumption that residuals are random and independent

The Durbin-Watson Statistic

- The Durbin-Watson statistic is used to test for autocorrelation

H_0 : residuals are not correlated

H_1 : positive autocorrelation is present

$$D = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}$$

- The possible range is $0 \leq D \leq 4$
- D should be close to 2 if H_0 is true
- D less than 2 may signal positive autocorrelation, D greater than 2 may signal negative autocorrelation

Testing for Positive Autocorrelation

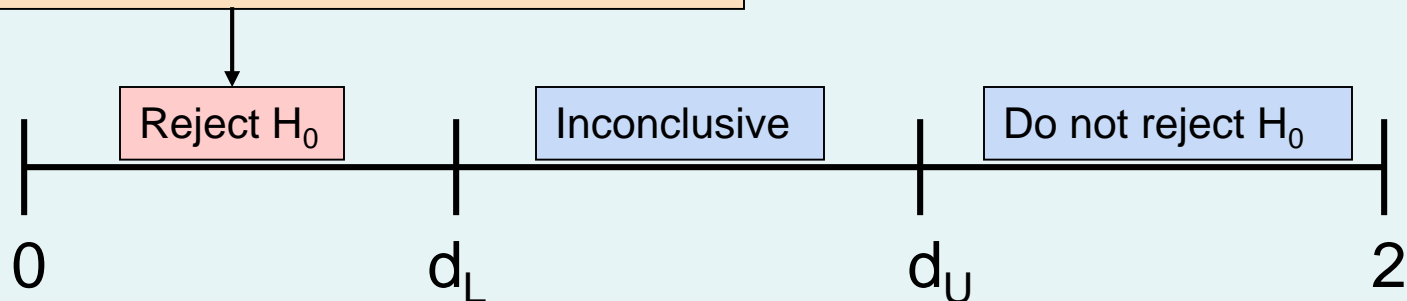
DCOVA

H_0 : positive autocorrelation does not exist

H_1 : positive autocorrelation is present

- Calculate the Durbin-Watson test statistic = D
(The Durbin-Watson Statistic can be found using Excel or Minitab)
- Find the values d_L and d_U from the Durbin-Watson table
(for sample size n and number of independent variables k)

Decision rule: reject H_0 if $D < d_L$

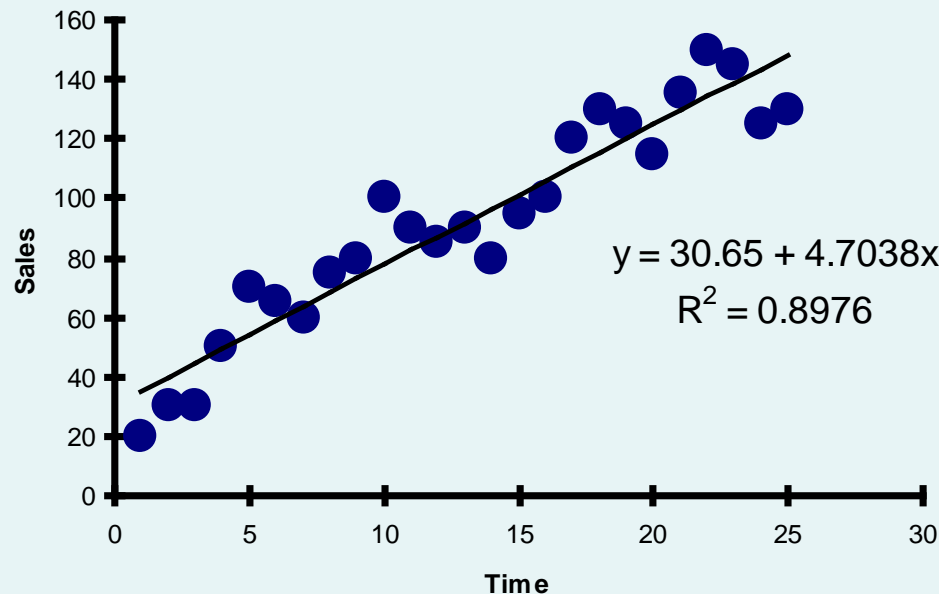


Testing for Positive Autocorrelation

(continued)

DCOVA

- Suppose we have the following time series data:



- Is there autocorrelation?

Testing for Positive Autocorrelation

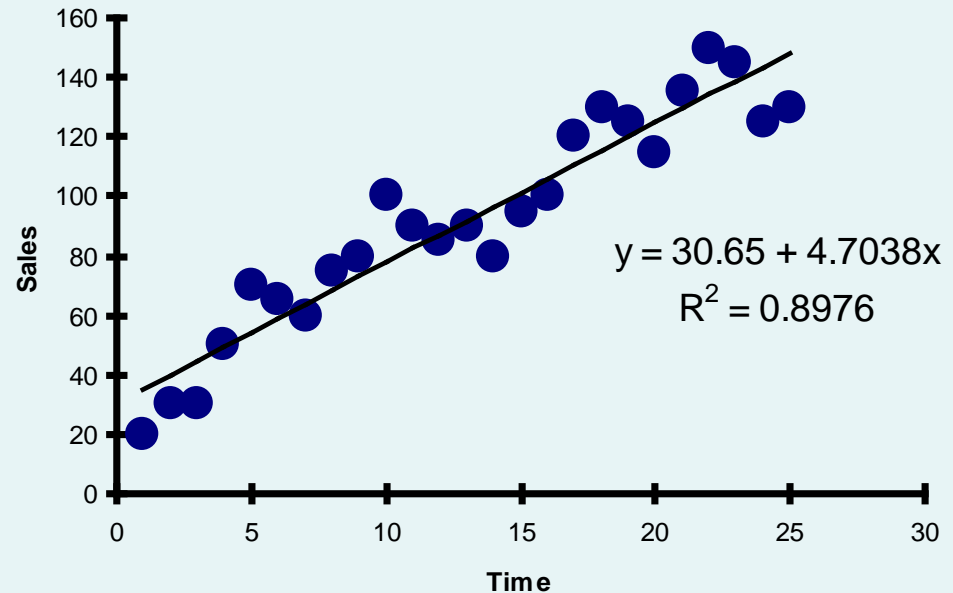
(continued)

DCOVA A

- Example with $n = 25$:

Excel/PHStat output:

Durbin-Watson Calculations	
Sum of Squared Difference of Residuals	3296.18
Sum of Squared Residuals	3279.98
Durbin-Watson Statistic	1.00494



$$D = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2} = \frac{3296.18}{3279.98} = 1.00494$$

Testing for Positive Autocorrelation

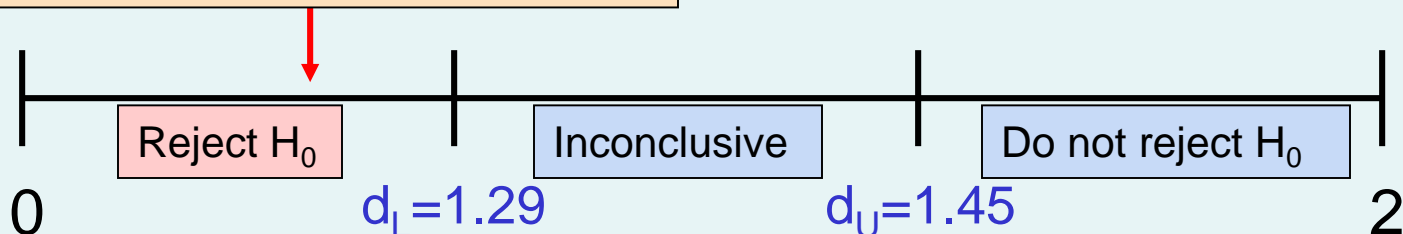
(continued)

DCOVA

- Here, $n = 25$ and there is $k = 1$ one independent variable
- Using the Durbin-Watson table, $d_L = 1.29$ and $d_U = 1.45$
- $D = 1.00494 < d_L = 1.29$, so reject H_0 and conclude that significant positive autocorrelation exists

Decision: reject H_0 since

$$D = 1.00494 < d_L$$



Inferences About the Slope

- The standard error of the regression slope coefficient (b_1) is estimated by

$$S_{b_1} = \frac{S_{YX}}{\sqrt{SSX}} = \frac{S_{YX}}{\sqrt{\sum (X_i - \bar{X})^2}}$$

where:

S_{b_1} = Estimate of the standard error of the slope

$S_{YX} = \sqrt{\frac{SSE}{n-2}}$ = Standard error of the estimate

Inferences About the Slope: t Test

- t test for a population slope
 - Is there a linear relationship between X and Y?
- Null and alternative hypotheses
 - $H_0: \beta_1 = 0$ (no linear relationship)
 - $H_1: \beta_1 \neq 0$ (linear relationship does exist)
- Test statistic

$$t_{\text{STAT}} = \frac{b_1 - \beta_1}{S_{b_1}}$$

$$\text{d.f.} = n - 2$$

where:

b_1 = regression slope
coefficient

β_1 = hypothesized slope

S_{b_1} = standard
error of the slope

Inferences About the Slope: t Test Example

DCOVA

House Price in \$1000s (y)	Square Feet (x)
245	1400
312	1600
279	1700
308	1875
199	1100
219	1550
405	2350
324	2450
319	1425
255	1700

Estimated Regression Equation:

$$\text{house price} = 98.25 + 0.1098 (\text{sq.ft.})$$

The slope of this model is 0.1098

Is there a relationship between the square footage of the house and its sales price?

Inferences About the Slope: t Test Example

DCOVA

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

From Excel output:

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>
Intercept	98.24833	58.03348	1.69296	0.12892
Square Feet	0.10977	0.03297	3.32938	0.01039

b_1

S_{b_1}

$$t_{\text{STAT}} = \frac{b_1 - \beta_1}{S_{b_1}} = \frac{0.10977 - 0}{0.03297} = 3.32938$$

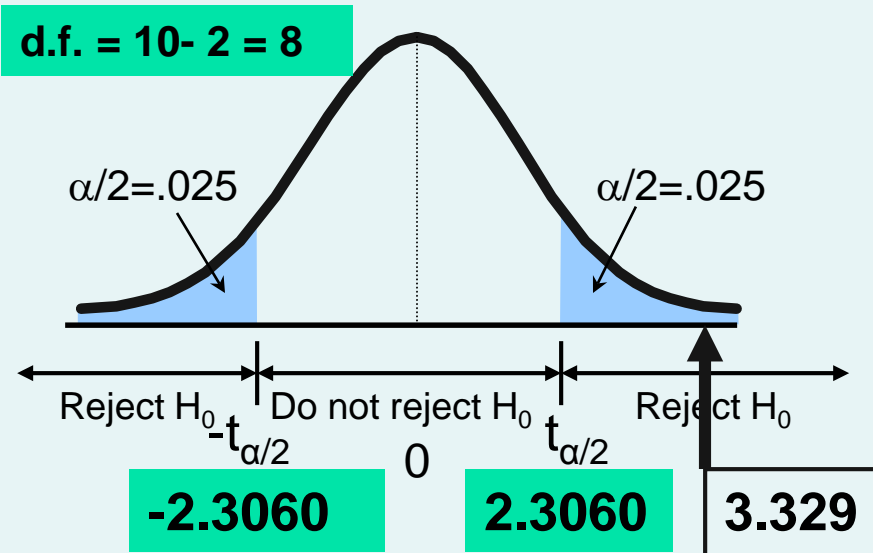
Inferences About the Slope: t Test Example

DCOVA

Test Statistic: $t_{\text{STAT}} = 3.329$

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$



Decision: Reject H_0

There is sufficient evidence
that square footage affects
house price

Inferences About the Slope: t Test Example

DCOVA

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

From Excel output:

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>
Intercept	98.24833	58.03348	1.69296	0.12892
Square Feet	0.10977	0.03297	3.32938	0.01039

p-value

Decision: Reject H_0 , since p-value $< \alpha$

There is sufficient evidence that square footage affects house price.

F Test for Significance

- F Test statistic:

$$F_{STAT} = \frac{MSR}{MSE}$$

where

$$MSR = \frac{SSR}{k}$$

$$MSE = \frac{SSE}{n - k - 1}$$

where F_{STAT} follows an F distribution with k numerator and $(n - k - 1)$ denominator **degrees of freedom**

(k = the number of independent variables in the regression model)

F-Test for Significance

Excel Output

Regression Statistics

Multiple R	0.76211
R Square	0.58082
Adjusted R Square	0.52842
Standard Error	41.33032
Observations	10

$$F_{\text{STAT}} = \frac{\text{MSR}}{\text{MSE}} = \frac{18934.9348}{1708.1957} = 11.0848$$

With 1 and 8 degrees of freedom

p-value for the F-Test

ANOVA	df	SS	MS	F	Significance F
Regression	1	18934.9348	18934.9348	11.0848	0.01039
Residual	8	13665.5652	1708.1957		
Total	9	32600.5000			

F Test for Significance

(continued)

DCOVA

$$H_0: \beta_1 = 0$$

$$H_1: \beta_1 \neq 0$$

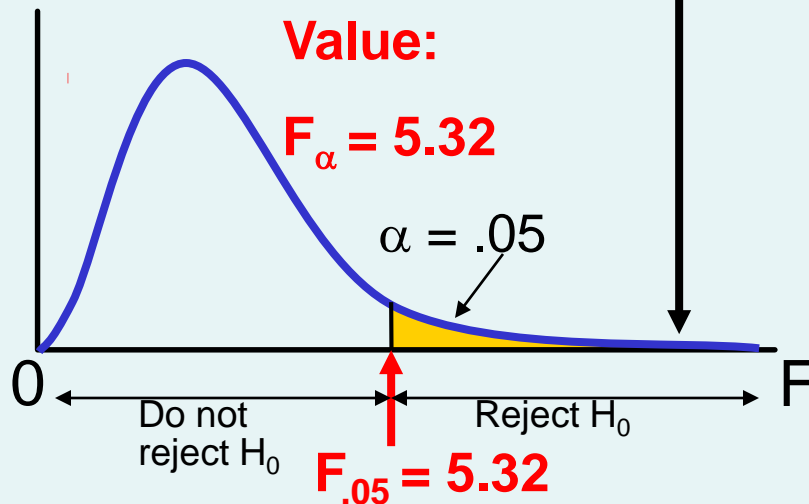
$$\alpha = .05$$

$$df_1 = 1 \quad df_2 = 8$$

Critical Value:

$$F_{\alpha} = 5.32$$

$$\alpha = .05$$



Test Statistic:

$$F_{\text{STAT}} = \frac{MSR}{MSE} = 11.08$$

Decision:

Reject H_0 at $\alpha = 0.05$

Conclusion:

There is sufficient evidence that house size affects selling price

Confidence Interval Estimate for the Slope

DCOVA

Confidence Interval Estimate of the Slope:

$$b_1 \pm t_{\alpha/2} S_{b_1}$$

d.f. = n - 2

Excel Printout for House Prices:

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	98.24833	58.03348	1.69296	0.12892	-35.57720	232.07386
Square Feet	0.10977	0.03297	3.32938	0.01039	0.03374	0.18580

At 95% level of confidence, the confidence interval for the slope is (0.0337, 0.1858)

Confidence Interval Estimate for the Slope

(continued)

DCOVA

	<i>Coefficients</i>	<i>Standard Error</i>	<i>t Stat</i>	<i>P-value</i>	<i>Lower 95%</i>	<i>Upper 95%</i>
Intercept	98.24833	58.03348	1.69296	0.12892	-35.57720	232.07386
Square Feet	0.10977	0.03297	3.32938	0.01039	0.03374	0.18580

Since the units of the house price variable is \$1000s, we are 95% confident that the average impact on sales price is between \$33.74 and \$185.80 per square foot of house size

This 95% confidence interval **does not include 0**.

Conclusion: There is a significant relationship between house price and square feet at the .05 level of significance



t Test for a Correlation Coefficient

DCOVA

- Hypotheses

$H_0: \rho = 0$ (no correlation between X and Y)

$H_1: \rho \neq 0$ (correlation exists)

- Test statistic

$$t_{\text{STAT}} = \frac{r - \rho}{\sqrt{\frac{1 - r^2}{n - 2}}}$$

(with $n - 2$ degrees of freedom)

where

$$r = +\sqrt{r^2} \text{ if } b_1 > 0$$

$$r = -\sqrt{r^2} \text{ if } b_1 < 0$$

t-test For A Correlation Coefficient

(continued)

DCOVA

Is there evidence of a linear relationship between square feet and house price at the .05 level of significance?

$H_0: \rho = 0$ (No correlation)

$H_1: \rho \neq 0$ (correlation exists)

$\alpha = .05$, $df = 10 - 2 = 8$

$$t_{\text{STAT}} = \frac{r - \rho}{\sqrt{\frac{1 - r^2}{n - 2}}} = \frac{.762 - 0}{\sqrt{\frac{1 - .762^2}{10 - 2}}} = 3.329$$

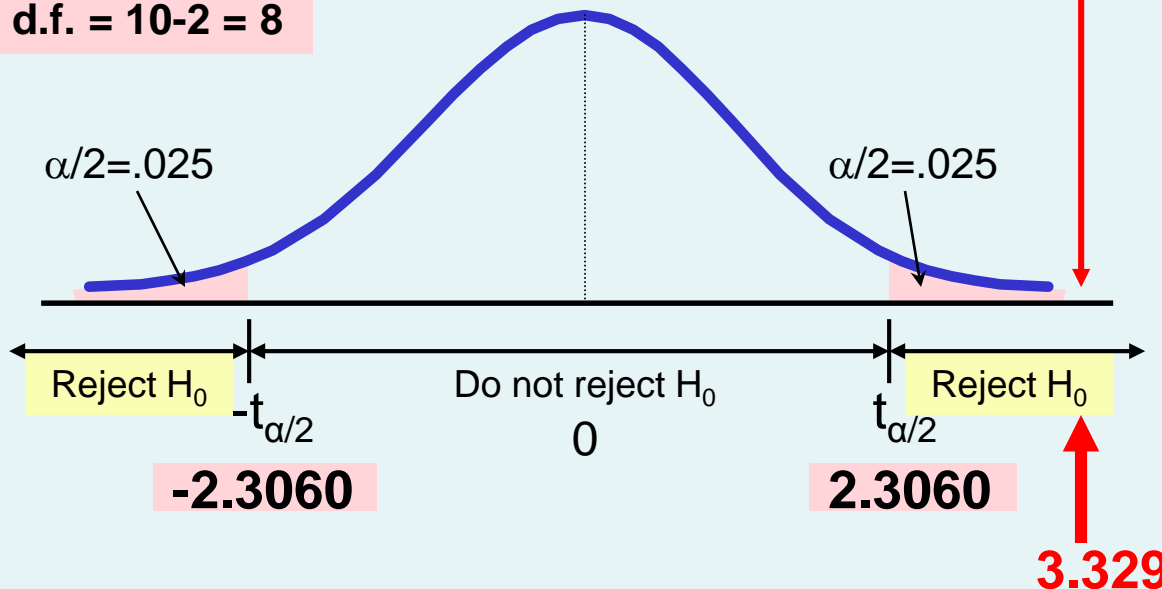
t-test For A Correlation Coefficient

(continued)

DCOVA

$$t_{\text{STAT}} = \frac{r - \rho}{\sqrt{\frac{1 - r^2}{n - 2}}} = \frac{.762 - 0}{\sqrt{\frac{1 - .762^2}{10 - 2}}} = 3.329$$

d.f. = 10 - 2 = 8



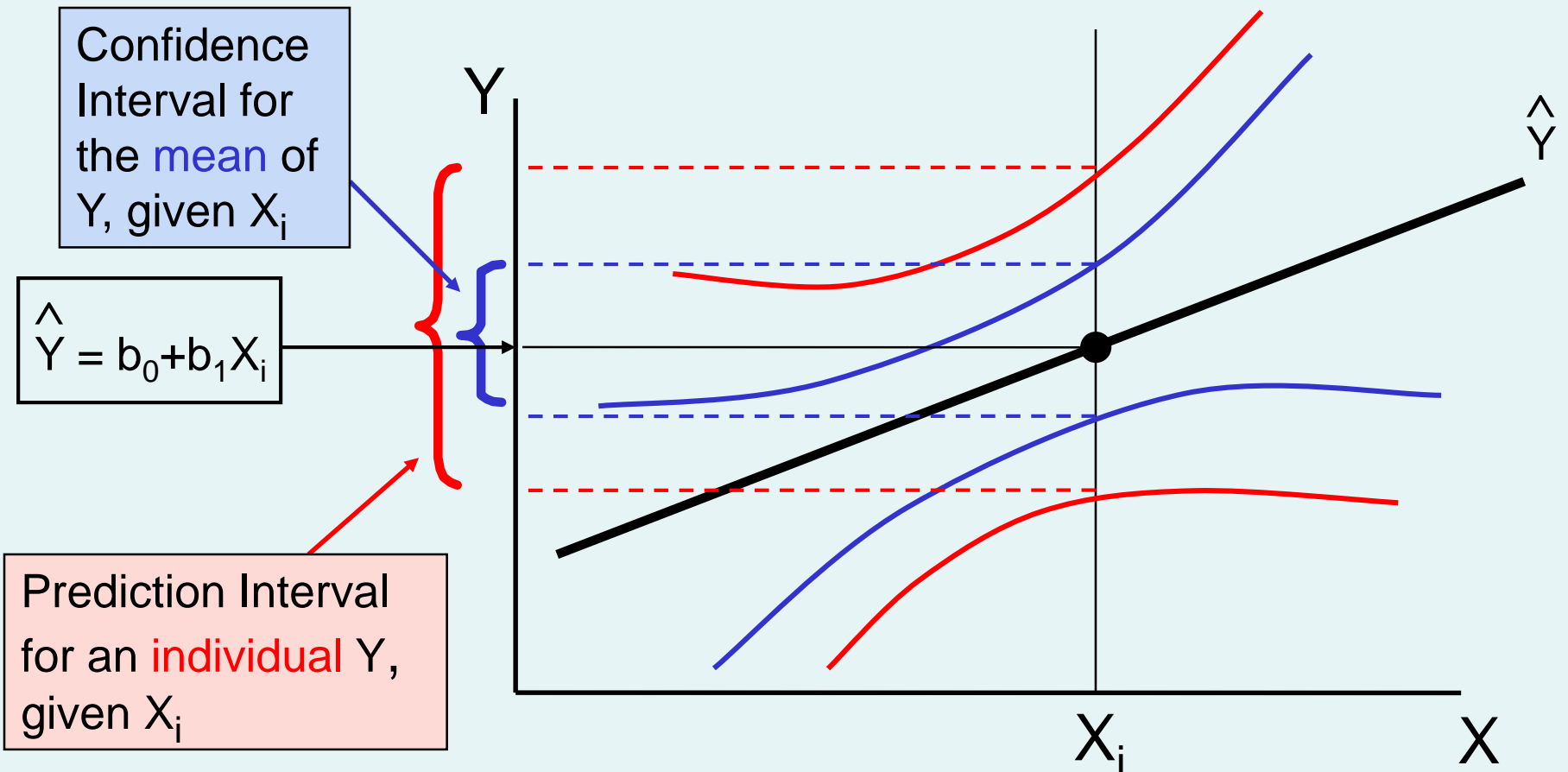
Decision:
Reject H_0

Conclusion:
There is **evidence** of a linear association at the 5% level of significance

Estimating Mean Values and Predicting Individual Values

DCOVA A

Goal: Form intervals around Y to express uncertainty about the value of Y for a given X_i



Confidence Interval for the Average Y, Given X

DCOVA

Confidence interval estimate for the **mean value of Y** given a particular X_i

Confidence interval for $\mu_{Y|X=X_i}$:

$$\hat{Y} \pm t_{\alpha/2} S_{YX} \sqrt{h_i}$$

Size of interval varies according to distance away from mean, \bar{X}

$$h_i = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{SSX} = \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum (X_i - \bar{X})^2}$$

Prediction Interval for an Individual Y, Given X

DCOVA

Confidence interval estimate for an
Individual value of Y given a particular X_i

Confidence interval for $Y_{X=X_i}$:

$$\hat{Y} \pm t_{\alpha/2} S_{YX} \sqrt{1 + h_i}$$

This extra term adds to the interval width to reflect
the added uncertainty for an individual case

Estimation of Mean Values: Example

DCOVA

Confidence Interval Estimate for $\mu_{Y|X=X_i}$

Find the 95% confidence interval for the mean price of 2,000 square-foot houses

Predicted Price $\hat{Y}_i = 317.85$ (\$1,000s)

$$\hat{Y} \pm t_{0.025} S_{YX} \sqrt{\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum (X_i - \bar{X})^2}} = 317.85 \pm 37.12$$

The confidence interval endpoints (from Excel) are 280.66 and 354.90, or from \$280,660 to \$354,900

Estimation of Individual Values: Example

DCOVA

Prediction Interval Estimate for $Y_{X=X_i}$

Find the 95% prediction interval for an individual house with 2,000 square feet

Predicted Price $\hat{Y}_i = 317.85$ (\$1,000s)

$$\hat{Y} \pm t_{0.025} S_{YX} \sqrt{1 + \frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum (X_i - \bar{X})^2}} = 317.85 \pm 102.28$$

The prediction interval endpoints from Excel are 215.50 and 420.07, or from \$215,500 to \$420,070



Finding Confidence and Prediction Intervals in Excel

DCOVA A

- From Excel, use
PHStat | regression | simple linear regression ...
- Check the
“confidence and prediction interval for $X=$ ”
box and enter the X -value and confidence level
desired

Finding Confidence and Prediction Intervals in Excel

(continued)

DCOVA

	A	B
1	Confidence Interval Estimate	
2		
3	Data	
4	X Value	2000
5	Confidence Level	95%
6		
7	Intermediate Calculations	
8	Sample Size	10
9	Degrees of Freedom	8
10	t Value	2.306006
11	Sample Mean	1715
12	Sum of Squared Difference	1571500
13	Standard Error of the Estimate	41.33032
14	h Statistic	0.151686
15	Average Predicted Y (YHat)	317.7838
16		
17	For Average Predicted Y (YHat)	
18	Interval Half Width	37.11952
19	Confidence Interval Lower Limit	280.6643
20	Confidence Interval Upper Limit	354.9033
21		
22	For Individual Response Y	
23	Interval Half Width	102.2813
24	Prediction Interval Lower Limit	215.5025
25	Prediction Interval Upper Limit	420.0651

Input values

\hat{Y}

Confidence Interval Estimate for $\mu_{Y|X=X_i}$

Prediction Interval Estimate for $Y_{X=X_i}$



Pitfalls of Regression Analysis

- Lacking an awareness of the assumptions underlying least-squares regression
- Not knowing how to evaluate the assumptions
- Not knowing the alternatives to least-squares regression if a particular assumption is violated
- Using a regression model without knowledge of the subject matter
- Extrapolating outside the relevant range



Strategies for Avoiding the Pitfalls of Regression

- Start with a scatter plot of X vs. Y to observe possible relationship
- Perform residual analysis to check the assumptions
 - Plot the residuals vs. X to check for violations of assumptions such as homoscedasticity
 - Use a histogram, stem-and-leaf display, boxplot, or normal probability plot of the residuals to uncover possible non-normality



Strategies for Avoiding the Pitfalls of Regression

(continued)

- If there is violation of any assumption, use alternative methods or models
- If there is no evidence of assumption violation, then test for the significance of the regression coefficients and construct confidence intervals and prediction intervals
- Avoid making predictions or forecasts outside the relevant range



Chapter Summary

In this chapter we discussed

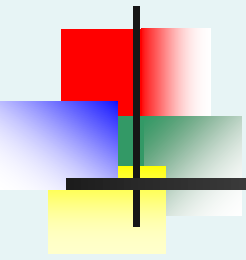
- Types of regression models
- The assumptions of regression and correlation
- Determining the simple linear regression equation
- Measures of variation
- Residual analysis
- Measuring autocorrelation



Chapter Summary

(continued)

- Making inferences about the slope
- Correlation -- measuring the strength of the association
- The estimation of mean values and prediction of individual values
- Possible pitfalls in regression and recommended strategies to avoid them



This work is protected by United States copyright laws and is provided solely for the use of instructors in teaching their courses and assessing student learning. Dissemination or sale of any part of this work (including on the World Wide Web) will destroy the integrity of the work and is not permitted. The work and materials from it should never be made available to students except by instructors using the accompanying text in their classes. All recipients of this work are expected to abide by these restrictions and to honor the intended pedagogical purposes and the needs of other instructors who rely on these materials.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording, or otherwise, without the prior written permission of the publisher.

Printed in the United States of America.