### Statistics for Business and Economics (14e) Metric Version

Chapters 14~15 (迴歸分析)

Anderson Sweeney Williams Camm Cochran Fry Ohlmann **Statistics for Business** ERSION & Economics Metric Version, 14th Edition

### Chapter 14 - Simple Linear Regression

- 14.1 Simple Linear Regression Model
- 14.2 Least Squares Method
- 14.3 Coefficient of Determination
- 14.4 Model Assumptions
- 14.5 Testing for Significance
- 14.6 Using the Estimated Regression Equation for Estimation and Prediction
- 14.7 Computer Solution
- 14.8 Residual Analysis: Validating Model Assumptions
- 14.9 Residual Analysis: Outliers and Influential Observations
- 14.10 Practical Advice: Big Data and Hypothesis Testing in Simple Linear Regression

量化模型

或是

透過量化模型描述觀察結果:

*y* = *f*(*x*) + error; 觀察值 = 訊號 + 雜訊

■ 數量化模型的關鍵:

→量化目標值y:定義問題!

→選取關鍵變數: $x_1, x_2, \ldots, x_p$ 

→建立量化模型:統計學習、機器學習。



http://www.ljagilamplighter.com/wp-content/uploads/noise.png

http://www.ljagilamplighter.com/wp-content/uploads/noise\_signal-mlab1\_png\_pagespeed\_ce\_b\_GTiE6tAg.png

## What we observe can be divided into:

## signal

MMMMM

noise

### what we see

https://blog.solidsignal.com/wp-content/uploads/2018/03/expclas1.gif

### Simple Linear Regression (1 of 2)

- Managerial decisions often are based on the relationship between two or more variables.
- <u>Regression analysis</u> can be used to develop an equation showing how the variables are related.
- The variable being predicted is called the <u>dependent variable</u> and is denoted by y.
- The variables being used to predict the value of the dependent variable are called the <u>independent variables</u> and are denoted by *x*.
- 註:因變數(應變數、相依變數) vs. 自變數(獨立變數、解釋變數)

### Simple Linear Regression (2 of 2)

- <u>Simple linear regression</u> involves one independent variable and one dependent variable.
- The relationship between the two variables is approximated by a straight line.
- Regression analysis involving two or more independent variables is called <u>multiple</u> <u>regression</u>.



© 2020 Cengage. May not be scanned, copied or duplicated, or posted to a publicly accessible website, in whole or i password-protected website or school-approved learning management system for classroom use.

什麼是迴歸分析?

**Regression** analysis

FITS A STRAIGHT LINE TO THIS MESSY SCATTERPLOT. 2 15 CALLED THE INDEPENDENT OR PREDICTOR VARIABLE, AND Y IS THE DEPENDENT OR RESPONSE VARIABLE. THE REGRESSION OR PREDICTION LINE HAS THE FORM

$$y = a + bx$$



### Simple Linear Regression Model

- The equation that describes how y is related to x and an error term is called the <u>regression model</u>.
- The <u>simple linear regression model</u> is

~

where:  $\beta_0$  and  $\beta_1$  are the parameters of the model.  $\varepsilon$  is a random variable called the <u>error term</u>.

### Simple Linear Regression Equation (1 of 4)

The simple linear regression equation is

 $E(y) = \beta_0 + \beta_1 x$ 

where:  $\beta_0$  is the y intercept of the regression line.  $\beta_1$  are the slope of the regression line. E(y) is the expected value of y for a given x value.

The graph of the regression equation is a straight line.

### Simple Linear Regression Equation (2 of 4)

**Positive Linear Relationship** 



### Simple Linear Regression Equation (3 of 4)

Negative Linear Relationship



### Simple Linear Regression Equation (4 of 4)

No Relationship



# 相關係數與迴歸係數



### **Estimated Simple Linear Regression Equation**

The <u>estimated simple linear regression equation</u>:

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

where:  $b_0$  is the y intercept of the regression line.  $b_1$  are the slope of the regression line.  $\hat{y}$  is the estimated value of y for a given x value.

### **Estimation Process**



Least Squares Method (1 of 3)

Least Squares Criterion

$$min\sum(y_i - \hat{y}_i)^2$$

where:

 $y_i$  = the observed value of the dependent variable for the  $i^{th}$  observation  $\hat{y}_i$  = the estimated value of the dependent variable for the  $i^{th}$  observation

Least Squares Method (2 of 3)

Slope for the Estimated Regression Equation

$$b_{1} = \frac{\sum (x_{i} - \bar{x})(y_{i} - \bar{y})}{\sum (x_{i} - \bar{x})^{2}}$$

where:  $x_i$  = the value of the independent variable for the  $i^{th}$  observation  $y_i$  = the value of the dependent variable for the  $i^{th}$  observation  $\bar{x}$  = the mean value for the independent variable  $\bar{y}$  = the mean value for the dependent variable Least Squares Method (3 of 3)

*y*-intercept for the estimated regression equation

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

### Simple Linear Regression

Reed Auto periodically has a special week-long sale. As part of the advertising campaign, Reed runs one or more television commercials during the weekend preceding the sale. Here are the data from a sample of 5 previous sales:

Number of
<u>Cars Sold (y)</u>
14
24
18
17
27
$\Sigma y = 100$
$\overline{y} = 20$

### **Estimated Regression Equation**

• Slope for the estimated regression equation

$$b_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \frac{20}{4} = 5$$

• *y*-intercept for the estimated regression equation

$$b_0 = \bar{y} - b_1 \bar{x} = 20 - 5(2) = 10$$

• Estimated Regression Equation:  $\hat{y} = 10 + 5x$ 

# Using Excel's Chart Tools for Scatter Diagram & Estimated Regression Equation



Coefficient of Determination (1 of 3)

Relationship Among SST, SSR, SSE

SST = SSR + SSE $\sum (y_i - \bar{y})^2 = \sum (\hat{y}_i - \bar{y})^2 + \sum (y_i - \hat{y}_i)^2$ 

where:

SST = total sum of squares SSR = sum of squares due to regression SSE = sum of squares due to error Coefficient of Determination (2 of 3)

The coefficient of determination is:

$$r^2 = \frac{SSR}{SST}$$

### where:

SSR = sum of squares due to regression SST = total sum of squares

### Coefficient of Determination (3 of 3)

$$r^2 = \frac{SSR}{SST} = \frac{100}{114} = 0.8772$$

The regression relationship is very strong; 87.72% of the variability in the number of cars sold can be explained by the linear relationship between the number of TV ads and the number of cars sold.

### Using Excel to Compute the Coefficient of Determination

Adding  $r^2$  Value to Scatter Diagram



### Sample Correlation Coefficient (1 of 2)

 $r_{xy} = (\text{sign of } b_1) \sqrt{\text{Coefficient of Determination}}$ 

 $r_{xy} = (\text{sign of } b_1) \sqrt{r^2}$ 

#### where:

 $b_1$  = the slope of the estimated regression equation

### Sample Correlation Coefficient (2 of 2)

$$r_{xy} = (\text{sign of } b_1) \sqrt{r^2}$$

The sign of  $b_1$  in the equation  $\hat{y} = 10 + 5x$  is positive, so

$$r_{xy} = \sqrt{0.8772} = 0.9366$$

### Assumptions About the Error Term $\varepsilon$

- 1. The error  $\varepsilon$  is a random variable with mean of zero.
- 2. The variance of  $\varepsilon$ , denoted by  $\sigma^2$ , is the same for all values of the independent variable.
- 3. The values of  $\varepsilon$  are independent.
- 4. The error  $\varepsilon$  is a normally distributed random variable.

### Testing for Significance (1 of 3)

- To test for a significant regression relationship, we must conduct a hypothesis test to determine whether the value of  $\beta_1$  is zero.
- Two tests commonly used are the *t* test and *F* test.
- Both the t test and F test require an estimate of  $\sigma^2$ , the variance of  $\varepsilon$  in the regression model.

### Testing for Significance (2 of 3)

An Estimate of  $\sigma^2$ 

The mean square error (MSE) provides the estimate of  $\sigma^2$ , and the notation  $s^2$  is also used.

$$s^2 = MSE = \frac{SSE}{n-2}$$

where:

$$SSE = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - b_0 - b_1 x_i)^2$$

### Testing for Significance (3 of 3)

To estimate  $\sigma$ , we take the square root of  $s^2$ .

The resulting *s* is called the <u>standard error of the estimate</u>.

$$s^{2} = MSE = \frac{SSE}{n-2}$$
  
 $s = \sqrt{MSE} = \sqrt{\frac{SSE}{n-2}}$ 

Testing for Significance: *t* Test (1 of 4)

Hypotheses:

$$\begin{array}{l} H_0: \ \beta_1 = 0 \\ H_a: \ \beta_1 \neq 0 \end{array}$$

**Test Statistic:** 

$$t = \frac{b_1}{s_{b_1}}$$
 where  $s_{b_1} = \frac{s}{\sqrt{\sum (x_i - \bar{x})^2}}$ 

Testing for Significance: t Test (2 of 4)

**Rejection Rule:** 

Reject  $H_0$  if p –value  $\leq \alpha$ Reject  $H_0$  if  $t \leq -t_{\alpha/2}$  or if  $t \geq t_{\alpha/2}$ 

where:

 $t_{\alpha/2}$  is based on a *t* distribution with n-2 degrees of freedom

### Testing for Significance: *t* Test (3 of 4)

1. Determine the hypotheses.

$$H_0: \ \beta_1 = 0$$
$$H_a: \ \beta_1 \neq 0$$

 $H \cdot R = 0$ 

2. Specify the level of significance.  $\alpha = 0.05$ 

3. Select the test statistic.

$$t = \frac{b_1}{s_{b_1}}$$

4. State the rejection rule.

Reject  $H_0$  if the *p*-value  $\leq 0.05$  or |t| > 3.182 with 3 degrees of freedom.

### Testing for Significance: *t* Test (4 of 4)

5. Compute the value of the test statistic.

$$t = \frac{b_1}{s_{b_1}} = \frac{5}{1.08} = 4.63$$

6. Determine whether to reject  $H_0$ .

t = 4.4541 provides an area of 0.01 in the upper tail. Hence the *p*-value < 0.02. Also, t = 4.63 > 3.182, so we can reject  $H_0$ .

## Confidence Interval for $\beta_{1}$ (1 of 3)

- We can use a 95% confidence interval for  $\beta_1$  to test the hypotheses just used in the *t* test.
- $H_0$  is rejected if the hypothesized value of  $\beta_1$  is not included in the confidence interval for  $\beta_1$ .
#### Confidence Interval for $\beta_{1}$ (2 of 3)

## $b_1 \pm t_{\alpha/2}(s_{b_1})$

#### where

 $b_1$  = the point estimator,

 $t_{\alpha/2}(s_{b_1})$  = the margin of error, and

 $t_{\alpha/2}$  = the *t* value providing an area of  $\alpha/2$  in the upper tail of a *t* distribution with n - 2 degrees of freedom.

## Confidence Interval for $\beta_{1}$ (3 of 3)

• Rejection Rule

Reject  $H_0$  if 0 is not included in the confidence interval for  $\beta_1$ .

• 95% Confidence Interval for  $\beta_1$ 

$$b_1 \pm t_{\alpha/2}(s_{b_1})$$
  
 $5 \pm 3.182(1.08)$   
 $5 \pm 3.44$   
1.56 to 8.44

Conclusion

#### 0 is not included in the confidence interval. Reject $H_0$ .

#### Testing for Significance: F Test (1 of 3)

• Hypotheses:  $H_0: \ \beta_1 = 0$  $H_a: \ \beta_1 \neq 0$ 

• Test Statistic: 
$$F = \frac{MSR}{MSE}$$

• Rejection Rule:

Reject  $H_0$  if the *p*-value  $\leq 0.05$  or if  $F \geq F_{\alpha}$ where  $F_{\alpha}$  is based on an *F* distribution with 1 degree of freedom in the numerator and n - 2 degrees of freedom in the denominator.

#### Testing for Significance: F Test (2 of 3)

1. Determine the hypotheses.

$$H_0: \beta_1 = 0$$
$$H_a: \beta_1 \neq 0$$

2. Specify the level of significance.  $\alpha = 0.05$ 

3. Select the test statistic.  $F = \frac{MSR}{MSE}$ 

4. State the rejection rule.

Reject  $H_0$  if the *p*-value  $\leq 0.05$  or  $F \geq 10.13$  with 1 *df* in the numerator and 3 *df* in the denominator.

#### Testing for Significance: F Test (3 of 3)

5. Compute the value of the test statistic.

$$F = \frac{MSR}{MSE} = \frac{100}{4.667} = 21.43$$

6. Determine whether to reject  $H_0$ .

F = 17.44 provides an area of 0.025 in the upper tail Thus, the *p*-value corresponding to F = 21.43 is less than 0.025. Hence, we reject  $H_0$ .

The statistical evidence is sufficient to conclude that we have a significant relationship between the number of TV ads aired and the number of cars sold.

## Some Cautions about the Interpretation of Significance Tests

- Rejecting H<sub>0</sub>: β<sub>1</sub> = 0 and concluding that the relationship between x and y is significant does not enable us to conclude that a <u>cause-and-effect relationship</u> is present between x and y.
- Just because we are able to reject  $H_0$ :  $\beta_1 = 0$  and demonstrate statistical significance does not enable us to conclude that there is a <u>linear relationship</u> between x and y.

0

Using the Estimated Regression Equation for Estimation and Prediction (1 of 2)

- A <u>confidence interval</u> is an interval estimate of the *mean value of* y for a given value of x.
- A <u>prediction interval</u> is used whenever we want to predict an individual value of y for a new observation corresponding to a given value of x.
- The margin of error is larger for a prediction interval.

假設需要預測以下兩個問題:

1.對於那些積雪深度為7米的年份,預測灌溉面積

2.今年最大積雪深度為7米,預測今年灌溉面積

→第一個問題「<u>預測平均值</u>」,第2個問題「<u>預測個別值</u>」

變異數

• 由(2.29b)可以注意到,  $\hat{Y}_h$  抽樣分配的變異數受到  $X_h$  與.  $\overline{X}$ 距離多遠的影響, 即受到  $(X_h - \overline{X})^2$  的影響。  $X_h$  離.  $\overline{X}$  越 遠,則  $(X_h - \overline{X})^2$  就越大,  $\hat{Y}_h$  的變異數也越大。圖 2.3 可 以對這個現象做一個直觀的解釋。在同一組X值之下重複 抽取兩個樣本所得到的兩條樣本迴歸線。

$$\hat{Y}_{h} = b_{0} + b_{1}X_{h} \qquad \sigma^{2}\left\{\hat{Y}_{h}\right\} = \sigma^{2}\left[\frac{1}{n} + \frac{\left(X_{h} - \overline{X}\right)^{2}}{\sum(X_{i} - \overline{X})^{2}}\right] \qquad (2.29b)$$

圖 2.3

具有相同平均數 7

及 X 的 兩 樣 本 其 不

同樣本之 b1 的變化

對 $\hat{Y}_{\mu}$ 的影響。



Using the Estimated Regression Equation for Estimation and Prediction (2 of 2)

• Confidence Interval Estimate of  $E(y^*)$ 

$$\hat{y}^* \pm t_{\alpha/2}(s_{\hat{y}^*})$$

• Prediction Interval Estimate of  $E(y^*)$  $\hat{y}^* \pm t_{\alpha/2}(s_{pred})$ 

where  $1 - \alpha$  is the confidence coefficient and  $t_{\alpha/2}$  is based on a t distribution with n - 2 degrees of freedom.

#### **Point Estimation**

If 3 TV ads are run prior to a sale, we expect the mean number of cars sold to be:

$$\hat{y} = 10 + 5(3) = 25$$
 cars

## Confidence Interval for $E(y^*)_{(1 \text{ of } 2)}$

Estimate of the Standard Deviation of  $\hat{y}^*$ :

$$s_{\hat{y}^*} = s \sqrt{\frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum (x_i - \bar{x})^2}}$$

$$s_{\hat{y}^*} = 2.16025 \sqrt{\frac{1}{5} + \frac{(3-2)^2}{(1-2)^2 + (3-2)^2 + \dots + (3-2)^2}}$$
$$s_{\hat{y}^*} = 2.16025 \sqrt{\frac{1}{5} + \frac{1}{4}} = \mathbf{1.4491}$$

## Confidence Interval for $E(y^*)$ (2 of 2)

The 95% confidence interval estimate of the mean number of cars sold when 3 TV ads are run is:

$$\hat{y}^* \pm t_{\alpha/2}(s_{\hat{y}^*})$$

#### $25 \pm 3.1824(1.4491)$

 $25 \pm 4.61$ 

#### 20.39 to 29.61 cars

#### Prediction Interval for $y^*$ (1 of 2)

Estimate of the Standard Deviation of an Individual Value of  $y^*$ 

$$s_{\text{pred}} = s \sqrt{1 + \frac{1}{n} + \frac{(x^* - \bar{x})^2}{\sum(x_i - \bar{x})^2}}$$
$$s_{\text{pred}} = 2.16025 \sqrt{1 + \frac{1}{5} + \frac{1}{4}}$$
$$s_{\text{pred}} = 2.16025(1.20416) = 2.6013$$

Prediction Interval for  $y^*$  (2 of 2)

The 95% prediction interval estimate of the number of cars sold in one particular week when 3 TV ads are run is:

 $\hat{y}^* \pm t_{\alpha/2}(s_{\rm pred})$ 

25 ± 3.1824(2.6013)

 $25\pm8.28$ 

#### 16.72 to 33.28 cars

#### **Computer Solution**

- Up to this point, you have seen how Excel can be used for various parts of a regression analysis.
- Excel also has a <u>comprehensive tool</u> in its Data Analysis package called <u>Regression</u>.
- The Regression tool can be used to perform a <u>complete regression analysis</u>.
- Performing the regression analysis computations without the help of a computer can be quite time consuming.
- On the next slide we show Minitab output for the Reed Auto Sales example.
- Recall that the independent variable was named Ads and the dependent variable was named Cars in the example.

#### Minitab Output (1 of 2)



#### Minitab Output (2 of 2)

- Minitab prints the estimated regression equation as Cars = 10 + 5 Ads.
- For each of the coefficients b<sub>0</sub> and b<sub>1</sub>, the output shows its value, standard deviation, t value, and p-value.
- Minitab prints the standard error of the estimate, *s*, as well as information about the goodness of fit.
- The standard ANOVA table is printed.
- Also provided are the 95% confidence interval estimate of the expected number of cars sold and the 95% prediction interval estimate of the number of cars sold for an individual weekend with 3 ads.

## Using Excel's Regression Tool (1 of 4)

#### Excel Output (top portion)

	А	В	С
9			
10	Regression Statistics	Regression Statistics	
11	Multiple R	0.936585812	
12	R Square	0.877192982	
13	Adjusted R Square	0.83625731	
14	Standard Error	2.160246899	
15	Observations	5	
16			

#### Using Excel's Regression Tool (2 of 4)

#### Excel Output (middle portion)

	Α	В	С	D	E	F
16						
17	ANOVA					
18		df	SS	MS	F	Significance F
19	Regression	1	100	100	21.4286	0.018986231
20	Residual	3	14	4.66667		
21	Total	4	114			
22						

## Using Excel's Regression Tool (3 of 4)

Excel Output (bottom-left portion)

	Α	В	С	D	E
22					
23		Coeffic.	Std. Err.	t Stat	P-values
24	Intercept	10	2.36643	4.2258	0.02424
25	TV Ads	5	1.08012	4.6291	0.01899
26					

#### Note: Columns F – I are not shown

#### Using Excel's Regression Tool (4 of 4)

Excel Output (bottom-right portion)

	Α	В	F	G	Н	I
22						
23		Coeffic.	Low. 95%	Up. 95%	Low. 95.0%	Up. 95.0%
24	Intercept	10	2.46895	17.53105	2.46895044	17.5310496
25	TV Ads	5	1.562562	8.437438	1.56256189	8.43743811
	А	В	F	G	н	I

Note: Columns C – E are hidden

#### **Residual Analysis**

- If the assumptions about the error term ε appear questionable, the hypothesis tests about the significance of the regression relationship and the interval estimation results may not be valid.
- The residuals provide the best information about  $\varepsilon$ .
- Residual for observation *i*:

$$y_i - \hat{y}_i$$

• Much of the residual analysis is based on an examination of graphical plots.

#### Residual Plot Against X (1 of 4)

If the assumption that the variance of  $\varepsilon$  is the same for all values of x is valid, and the assumed regression model is an adequate representation of the relationship between the variables, then the residual plot should give an overall impression of a horizontal band of points.



#### Residual Plot Against X (2 of 4)



#### Residual Plot Against X (3 of 4)

Residuals

<u>Observation</u>	Predicted Cars Sold	<u>Residuals</u>
1	15	-1
2	25	-1
3	20	-2
4	15	2
5	25	2

#### Residual Plot Against X (4 of 4)



## Standardized Residuals

Standardized Residual for Observation *i* 

$$\frac{y_i - \hat{y}_i}{s_{y_i - \hat{y}_i}}$$

$$s_{y_i-\hat{y}_i}=s\sqrt{1-h_i}$$

Where:

$$h_i = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{\sum (x_i - \bar{x})^2}$$

Standardized Residual Plot (1 of 3)

- The standardized residual plot can provide insight about the assumption that the error term  $\epsilon$  has a normal distribution.
- If this assumption is satisfied, the distribution of the standardized residuals should appear to come from a standard normal probability distribution.

#### Standardized Residual Plot (2 of 3)

#### **Standardized Residuals**

<u>Observation</u>	Predicted y	<u>Residual</u>	Standardized Residual
1	15	-1	-0.5345
2	25	-1	-0.5345
3	20	-2	-1.0690
4	15	2	1.0690
5	25	2	1.0690

#### Standardized Residual Plot (3 of 3)

#### **Standardized Residual Plot**

All of the standardized residuals are between -1.5 and +1.5, indicating that there is no reason to question the assumption that  $\varepsilon$  has a normal distribution.



#### **Outliers and Influential Observations**

**Detecting Outliers** 

- An <u>outlier</u> is an observation that is unusual in comparison with the other data.
- Minitab classifies an observation as an outlier if its standardized residual value is

< -2 or > +2.

- This standardized residual rule sometimes fails to identify an unusually large observation as being an outlier.
- This rule's shortcoming can be circumvented by using studentized deleted residuals.

• The  $|i^{th}$  studentized deleted residual will be larger than the  $|i^{th}$  studentized residual.

## 為什麼要檢查殘差?



#### Statistics for Business and Economics (14e) Metric Version

Anderson, Sweeney, Williams, Camm, Cochran, Fry, Ohlmann © 2020 Cengage Learning Anderson Sweeney Williams Camm Cochran Fry Ohlmann

# Statistics for Business & Economics



#### Chapter 15 - Multiple Regression

- 15.1 Multiple Regression Model
- 15.2 Least Squares Method
- 15.3 Multiple Coefficient of Determination
- 15.4 Model Assumptions
- 15.5 Testing for Significance
- 15.6 Using the Estimated Regression Equation for Estimation and Prediction
- 15.7 Categorical Independent Variables
- 15.8 Residual Analysis
- 15.9 Logistic Regression

#### MULTIPLE REGRESSION ANALYSIS PROCEDURE



## Multiple Regression ( 複迴歸 )

- In this chapter we continue our study of regression analysis by considering situations involving two or more independent variables.
- This subject area, called <u>multiple regression analysis</u>, enables us to consider more factors and thus obtain better estimates than are possible with simple linear regression.
#### Multiple Regression Model

The equation that describes how the dependent variable y is related to the independent variables  $x_1, x_2, ..., x_p$  and an error term is:

$$y=\beta_0+\beta_1x_1+\beta_2x_2+\cdots+\beta_px_p+\varepsilon$$

where  $\beta_0$ ,  $\beta_1$ ,  $\beta_2$ ,...,  $\beta_p$  are the parameters, and  $\varepsilon$  is a random variable called the error term.

#### **Multiple Regression Equation**

# The equation that describes how the mean value of y is related to $x_1, x_2, \dots, x_p$ is:

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p$$

#### Estimated Multiple Regression Equation

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots + b_p x_p$$

A simple random sample is used to compute the sample statistics  $b_0$ ,  $b_1$ ,  $b_2$ ,...,  $b_p$  that are used as the point estimators of the parameters  $\beta_0$ ,  $\beta_1$ ,  $\beta_2$ ,...,  $\beta_p$ .

#### Multiple Regression Model (1 of 3)

Example: Programmer Salary Survey

A software firm collected data for a sample of 20 computer programmers. A suggestion was made that regression analysis could be used to determine if salary was related to the years of experience and the score on the firm's programmer aptitude test.

The years of experience, score on the aptitude test, and corresponding annual salary (\$1000s) for a sample of 20 programmers is shown on the next slide.

## Multiple Regression Model (2 of 3)

Exper. (Yrs.)	Test score	Salary (\$1000s)	Exper. (Yrs.)	Test score	Salary (\$1000s)
4	78	24.0	9	88	38.0
7	100	43.0	2	73	26.6
1	86	23.7	10	75	36.2
5	82	34.3	6	74	29.0
10	84	38.0	8	87	34.0
0	75	22.2	4	79	30.1
1	80	23.1	6	94	33.9
6	83	30.0	3	70	28.2
6	91	33.0	3	89	30.0

#### Multiple Regression Model (3 of 3)

Suppose we believe that salary (y) is related to the years of experience  $(x_1)$  and the score on the programmer aptitude test  $(x_2)$  by the following regression model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \varepsilon$$

#### where

y =annual salary (\$1000s)

 $x_1$  = years of experience

 $x_2$  = score on programmer aptitude test

## Solving for the Estimates of $\beta_0$ , $\beta_1$ , $\beta_2$ (1 of 2)



## Solving for the Estimates of $\beta_0$ , $\beta_1$ , $\beta_2$ (2 of 2)

#### **Regression Equation Output:**

Predictor	Coef	SE Coef	т	р
Constant	3.17394	6.15607	0.5156	0.61279
Experience	1.4039	0.19857	7.0702	1.9E-06
Test Score	0.25089	0.07735	3.2433	0.00478
		-		

The estimated regression equation is:

SALARY = 3.174 + 1.404(EXPER) + 0.251(SCORE) (Note: Predicted salary will be in thousands of dollars.)

## **Testing for Significance**

- In simple linear regression, the F and t tests provide the same conclusion.
- In multiple regression, the F and t tests have different purposes.
- The *F* test is referred to as the <u>test for overall significance</u>.
- If the F test shows an overall significance, the t test is used to determine whether each of the individual independent variables is significant.
- A separate *t* test is conducted for each of the independent variables in the model.
- We refer to each of these t tests as a <u>test for individual significance</u>.

#### Testing for Significance: F Test

 $H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0$ • Hypotheses:

 $H_a$ : One or more of the parameters is not equal to zero.

• Test Statistic: 
$$F = \frac{MSR}{MSE}$$

$$' = \frac{MSR}{MSE}$$

Reject  $H_0$  if the *p*-value  $\leq \alpha$  or if  $F \geq F_{\alpha}$ • Rejection Rule: where  $F_{\alpha}$  is based on an F distribution with p degree of freedom in the numerator and n - p - 1 degrees of freedom in the denominator.

#### F Test for Overall Significance (1 of 3)

• Hypotheses:  $H_0: \beta_1 = \beta_2 = 0$ 

 $H_a$ : One or more of the parameters is not equal to zero.

• Rejection Rule: For  $\alpha = 0.05$  and df = 2, 17;  $F_{0.05} = 3.59$ 

Reject  $H_0$  if the *p*-value  $\leq 0.05$  or  $F \geq 3.59$ 

F Test for Overall Significance (2 of 3)

#### ANOVA Output



#### F Test for Overall Significance (3 of 3)

• Test Statistics: 
$$F = \frac{MSR}{MSE} = \frac{250.16}{5.85} = 42.76$$

• Conclusion: p-value  $\leq 0.05$  so we can reject  $H_0$ . Also,  $F = 42.76 \geq 3.59$ 

#### Testing for Significance: t Test

• Hypotheses:  $H_0: \beta_i = 0$  $H_a: \beta_i \neq 0$ 

• Test Statistic: t =

$$=\frac{b_i}{s_{b_i}}$$

• Rejection Rule: Reject  $H_0$  if the *p*-value  $\leq \alpha$  or if  $t \leq -t_{\alpha/2}$  or if  $t \geq t_{\alpha/2}$ where  $t_{\alpha/2}$  is based on a *t* distribution with n - p - 1 degrees of freedom.

### *t* Test for Significance of Individual Parameters (1 of 4)

- Hypotheses:  $H_0: b_i = 0$  $H_a: b_i \neq 0$
- $\begin{array}{ll} \bullet \mbox{ Rejection Rule:} & \mbox{ For } \alpha = 0.05 \mbox{ and } df = 17; \ t_{0.025} = 2.11 \\ & \mbox{ Reject } H_0 \mbox{ if the } p\mbox{-value} \leq 0.05 \mbox{ or } t \leq -2.11 \mbox{ or if } t \geq 2.11 \\ \end{array}$

## t Test for Significance of Individual Parameters (2 of 4)

**Regression Equation Output** 

Predictor	Coef	SE Coef	т	р
Constant	3.17394	6.15607	0.5156	0.61279
Experience	1.4039	0.19857	7.0702	1.9E-06
Test Score	0.25089	0.07735	3.2433	0.00478

## t Test for Significance of Individual Parameters (3 of 4)

#### **Regression Equation Output**

Predictor	Coef	SE Coef	т	р				
Constant	3.17394	6.15607	0.5156	0.61279				
Experience	1.4039	0.19857	7.0702	1.9E-06				
Test Score	0.25089	0.07735	3.2433	0.00478				
<i>t</i> statistic and <i>p</i> -value used to test for the individual significance of "Test Score"								

*t* Test for Significance of Individual Parameters (4 of 4)

• Test Statistics:

$$t = \frac{b_1}{s_{b_1}} = \frac{1.4039}{0.1986} = 7.07$$
$$t = \frac{b_2}{s_{b_2}} = \frac{0.25089}{0.07735} = 3.24$$

Conclusions

Reject both  $H_0$ :  $\beta_1 = 0$  and  $H_0$ :  $\beta_2 = 0$ . Both independent variables are significant.

## Testing for Significance: Multicollinearity (共線性)

- The term <u>multicollinearity</u> refers to the correlation among the independent variables.
- When the independent variables are highly correlated (say, |r| > 0.7), it is not possible to determine the separate effect of any particular independent variable on the dependent variable.
- If the estimated regression equation is to be used only for predictive purposes, multicollinearity is usually not a serious problem.
- Every attempt should be made to avoid including independent variables that are highly correlated.

## **Chapter 16 - Regression Analysis: Model Building**

- 16.1 General Linear Model
- 16.2 Determining When to Add or Delete Variables
- 16.3 Analysis of a Larger Problem
- 16.4 Variable Selection Procedures
- 16.5 Multiple Regression Approach to Experimental Design
- 16.6 Autocorrelation and the Durbin-Watson Test

## Variable Selection Procedures

- 1. Stepwise Regression
- 2. Forward Selection
- 3. Backward Elimination
- 4. Best-Subsets Regression
- The first three procedures are Iterative; one independent variable at a time is added or deleted based on the *F* statistic.
- The first 3 procedures are heuristics and therefore offer no guarantee that the best model will be found.
- In the fourth procedure different subsets of the independent variables are evaluated.

#### Variable Selection: Stepwise Regression (1 of 2)

- At each iteration, the first consideration is to see whether the least significant variable currently in the model can be removed because its *F* value is less than the user-specified or default *Alpha to remove*.
- If no variable can be removed, the procedure checks to see whether the most significant variable not in the model can be added because its *F* value is greater than the user-specified or default *Alpha to enter*.
- If no variable can be removed and no variable can be added, the procedure stops.

#### Variable Selection: Stepwise Regression (2 of 2)



#### Variable Selection: Forward Selection (1 of 2)

- This procedure is similar to stepwise regression, but does not permit a variable to be deleted.
- It adds variables one at a time as long as a significant reduction in the error sum of squares (SSE) can be achieved.

#### Variable Selection: Forward Selection (2 of 2)



## Variable Selection: Backward Elimination (1 of 2)

- This procedure begins with a model that includes all the independent variables the modeler wants considered.
- It then attempts to delete one variable at a time by determining whether the least significant variable currently in the model can be removed because its *p*-value is less than the userspecified or default value.
- Once a variable has been removed from the model, it cannot reenter at a subsequent step.

#### Variable Selection: Backward Elimination (2 of 2)



#### Variable Selection: Best-Subsets Regression (1 of 2)

- The three preceding procedures are one-variable-at-a-time methods offering no guarantee that the best model for a given number of variables will be found.
- Minitab output identifies the two best one-variable estimated regression equations, the two best two-variable equations, and so on.
- Some software packages include <u>best-subsets regression</u> that enables the user to find, given a specified number of independent variables, the best regression model.

## Variable Selection: Best-Subsets Regression (2 of 2)

Example: PGA Tour Data

The Professional Golfers Association keeps a variety of statistics regarding performance measures. Data include the average driving distance, percentage of drives that land in the fairway, percentage of greens hit in regulation, average number of putts, percentage of sand saves, and average score.

#### Variable-Selection Procedures (1 of 8)

Variable Names and Definitions

- **Drive**: average length of a drive in yards
- **Fair**: percentage of drives that land in the fairway
- **Green**: percentage of greens hit in regulation (a par-3 green is "hit in regulation" if the player's first shot lands on the green)
- **Putt**: average number of putts for greens that have been hit in regulation
- **Sand**: percentage of sand saves (landing in a sand trap and still scoring par or better)
- Score: average score for an 18-hole round

#### Variable-Selection Procedures (2 of 8)

#### Sample Data (Part 1)

<u>Drive</u>	<u>Fair</u>	<u>Green</u>	<u>Putt</u>	<u>Sand</u>	<u>Score</u>
277.6	.681	.667	1.768	.550	69.10
259.6	.691	.665	1.810	.536	71.09
269.1	.657	.649	1.747	.472	70.12
267.0	.689	.673	1.763	.672	69.88
267.3	.581	.637	1.781	.521	70.71
255.6	.778	.674	1.791	.455	69.76
272.9	.615	.667	1.780	.476	70.19
265.4	.718	.699	1.790	.551	69.73

#### Variable-Selection Procedures (3 of 8)

#### Sample Data (Part 2)

<u>Drive</u>	<u>Fair</u>	<u>Green</u>	<u>Putt</u>	<u>Sand</u>	<u>Score</u>
272.6	.660	.672	1.803	.431	69.97
263.9	.668	.669	1.774	.493	70.33
267.0	.686	.687	1.809	.492	70.32
266.0	.681	.670	1.765	.599	70.09
258.1	.695	.641	1.784	.500	70.46
255.6	.792	.672	1.752	.603	69.49
261.3	.740	.702	1.813	.529	69.88
262.2	.721	.662	1.754	.576	70.27

#### Variable-Selection Procedures (4 of 8)

#### Sample Data (Part 3)

<u>Drive</u>	<u>Fair</u>	<u>Green</u>	<u>Putt</u>	<u>Sand</u>	<u>Score</u>
260.5	.703	.623	1.782	.567	70.72
271.3	.671	.666	1.783	.492	70.30
263.3	.714	.687	1.796	.468	69.91
276.6	.634	.643	1.776	.541	70.69
252.1	.726	.639	1.788	.493	70.59
263.0	.639	.647	1.760	.374	70.81
253.5	.732	.693	1.797	.518	70.26
266.2	.681	.657	1.812	.472	70.96

#### Variable-Selection Procedures (5 of 8)

#### Sample Correlation Coefficients

	Score	Drive	Fair	Green	Putt
Drive	154				
Fair	427	679			
Green	556	045	.421		
Putt	.258	139	.101	.354	
Sand	278	024	.265	.083	296

#### Variable-Selection Procedures (6 of 8)

#### Best Subsets Regression of SCORE

Varsx	R-sq	R-sq(a)	С-р	S	D	F	G	Р	S
1	30.9	27.9	26.9	.39685			Х		
1	18.2	14.6	35.7	.43183		Х			
2	54.7	50.5	12.4	.32872	Х	Х			
2	54.6	50.5	12.5	.32891			Х	Х	
3	60.7	55.1	10.2	.31318	Х	Х		Х	
3	59.1	53.3	11.4	.31957	Х	Х	Х		
4	72.2	66.8	4.2	.26913	Х	Х	Х	Х	
4	60.9	53.1	12.1	.32011	Х	Х		Х	Х
5	72.6	65.4	6.0	.27499	Х	Х	Х	Х	Х

## Variable-Selection Procedures (7 of 8)

Minitab Output

#### The regression equation

Score = 74.678 - .0398(Drive) - 6.686(Fair) - 10.342(Green) + 9.858(Putt)

Predictor	Coef	Stdev	t-ratio	р
Constant	74.678	6.952	10.74	.000
Drive	0398	.01235	-3.22	.004
Fair	-6.686	1.939	-3.45	.003
Green	-10.342	3.561	-2.90	.009
Putt	9.858	3.180	3.10	.006
<b>s</b> = .2691	R-sq =	72.4%	R-sq(adj)	= 66.8%
### Variable-Selection Procedures (8 of 8)

### Minitab Output

#### **Analysis of Variance**

SOURCE	DF	SS	MS	F	Р
Regression	4	3.79469	.94867	13.10	.000
Error	20	1.44865	.07243		
Total	24	5.24334			

## Logistic regression

- Most important model for *categorical response*  $(y_i)$  data
- Categorical response with 2 levels (*binary*: 0 and 1)
- Categorical response with  $\geq$  3 levels (nominal or ordinal)
- Predictor variables (x<sub>i</sub>) can take on *any* form: binary, categorical, and/or continuous

# Logistic Regression

- Models relationship between set of variables  $X_i$ 
  - dichotomous (yes/no, smoker/nonsmoker,...)
  - categorical (social class, race, ... )
  - continuous (age, weight, gestational age, ...)

### and

- dichotomous categorical response variable *Y*
- e.g. <u>Success/Failure</u>, <u>Remission/No Remission</u>, <u>Survived/Died</u>, <u>CHD/No CHD</u>, <u>Low Birth Weight/Normal Birth Weight</u>...



### Logistic Regression Curve



## Logit Transformation

Logistic regression models transform probabilities called *logits*.

where 
$$logit(p_i) = log\left(\frac{p_i}{1-p_i}\right)$$

- *i* indexes all cases (observations).
- $p_i$  is the probability the event (a sale, for example) occurs in the *i*<sup>th</sup> case.
- log is the natural log (to the base e).

# Logistic regression model with a single continuous predictor

logit 
$$(p_i) = \log (\text{odds}) = \beta_0 + \beta_1 X_1$$

where

logit( $p_i$ )logit transformation of the probability of the event $\beta_0$ intercept of the regression line $\beta_1$ slope of the regression line

## Assumption



# Interpretation of a single *continuous* parameter

- The sign ( $\pm$ ) of  $\beta$  determines whether the **log odds** of y is increasing or decreasing *for every 1-unit increase* in x.
- If  $\beta > 0$ , there is an increase in the **log odds** of y for every 1-unit increase in x.
- If  $\beta < 0$ , there is a decrease in the **log odds** of y for every 1-unit increase in x.
- If β = 0 there is *no linear relationship* between the log odds and x.

### Parameter Interpretation (ctd).

• Exponentiating both sides of the logit link function we get the following:

$$\left(\frac{p_i}{1-p_i}\right) = \text{odds} = \exp(\beta_0 + \beta_1 X_1) = e^{\beta 0} e^{\beta 1 X_1}$$

- The odds increase **multiplicatively** by  $e^{\beta}$  for every 1-unit increase in *x*.
- Whether the increase is greater than 1 or less than one depends on whether  $\beta > 0$  or  $\beta < 0$ .
- The odds at X = x+1 are  $e^{\beta}$  times the odds at X = x. Therefore,  $e^{\beta}$  is an odds ratio!

Logistic regression model with a single *categorical* ( $\geq 2$  *levels*) predictor

$$logit (p_i) = log (odds) = \beta_0 + \beta_k X_k$$

where

 $\begin{array}{ll} logit(p_i) & \mbox{logit transformation of the probability of the event} \\ \beta_0 & \mbox{intercept of the regression line} \\ \beta_k & \mbox{difference between the logits for category k vs. the} \\ & \mbox{reference category} \end{array}$ 

## Logistic Regression

**Example:** Coronary Heart Disease (CD) and Age In this study sampled individuals were examined for signs of CD (present = 1/absent = 0) and its potential relationship with the age (yrs.) was considered.

•											
		Agegrp	Age	CD	Agegrp	Age	CD		Agegrp	Age	CD
•	1	1	20	0	2	30	0		8	60	0
•	2	1	23	0	2	30	0	. –	8	60	1
•	3	1	24	0	2	30	0	-	8	61	1
٠	4	1	25	0	2	30	0	-	8	62	1
٠	5	1	25	1	2	30	1	-	8	62	1
٠	6	1	26	0	2	32	0		8	63	1
٠	7	1	26	0	2	32	0		8	64	
٠	8	1	28	0	2	33	0	-		04	
٠	9	1	28	0	2	33	0	-	8	64	1
٠	10	1	29	0	2	34	0	1	8	65	1
٠	11	2	30	0	2	34	0	_	8	69	1

Note: There are 100 subjects participating in the study.

## Logistic Regression • How can we analyze these data?



Means and Std Deviations								
Level	Number	М	ean	Std Dev	Std			
0	57	39.1	754	10.2018				
1	43	51.2	2791	9.9793				
t Test								
1-0								
Assuming	unequal v	ariand	es		[			
Differenc	e 12.	1036	t Ratio	5.94	727			
Std Err Di	if 2.	0352	DF	91.61	987			
Upper CL	. Dif 16.	1459	Prob >	t  <.00	001*			
Lower Cl	L Dif 8.	0614	Prob >	t <.00	001*			
Confiden	се	0.95	Prob <	t 1.00	000			

#### **Non-pooled t-test**

The mean age of the individuals with some signs of coronary heart disease is 51.28 years vs. 39.18 years for individuals without signs (t = 5.95, p < .0001).

© 2020 Cengage. password-protected