

統計學 (二)

Anderson's Statistics for Business & Economics (14/E)

Chapter 15(2): Multiple Regression

上課時間地點: 二 D56, 資訊 140306

授課教師: 吳漢銘 (國立政治大學統計學系副教授)

教學網站: <http://www.hmwu.idv.tw>

系級: _____ 學號: _____ 姓名: _____

15.1 Multiple Regression Model

- (Recall) that the variable being predicted or explained is called the _____ variable and the variable being used to predict or explain the dependent variable is called the _____ variable.
- Multiple regression analysis is the study of how a dependent variable y is related to _____ variables. In the general case, we will use _____ to denote the number of independent variables.
- The concepts of a regression model and a regression equation introduced in the preceding chapter are _____ in the multiple regression case.
- Multiple regression model:** The equation that describes how the dependent variable y is related to the independent variables x_1, x_2, \dots, x_p and an error term is called the multiple regression model.

(15.1)

- In the multiple regression model, $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ are the _____ and the error term ϵ is a _____. y is a linear function of x_1, x_2, \dots, x_p plus the error term ϵ .
- The error term accounts for the _____ in y that _____ by the linear effect of the p independent variables.

7. **(Multiple regression equation):**The equation that describes how the mean value of y is related to x_1, x_2, \dots, x_p is called the multiple regression equation.

$$\text{_____} \tag{15.2}$$

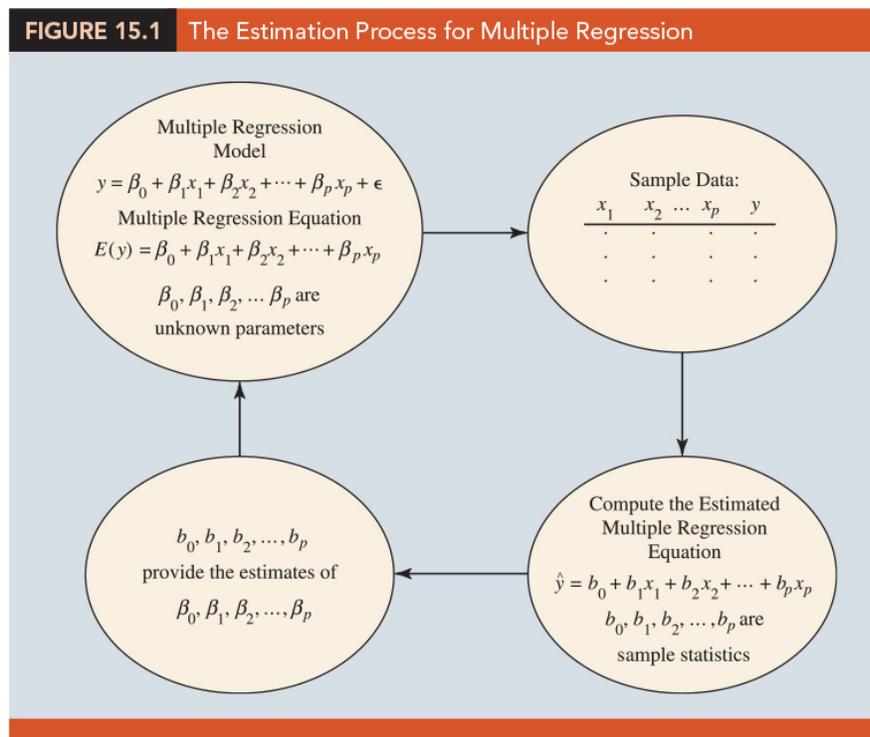
under the assumption that the mean or expected value of ϵ is zero.

8. **The estimated multiple regression equation:**

$$\text{_____} \tag{15.3}$$

where $b_0, b_1, b_2, \dots, b_p$ are the estimates of $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ and \hat{y} is the predicted value of the dependent variable

9. (Figure 15.1)



15.2 Least Squares Method

1. The least squares method is used to develop the estimated multiple regression equation:

$$\text{_____} \quad (15.4)$$

where y_i is observed value of the dependent variable for the i th observation, \hat{y}_i is predicted value of the dependent variable for the i th observation

2. In multiple regression, however, the presentation of the formulas for the regression coefficients $\beta_0, \beta_1, \beta_2, \dots, \beta_p$ involves the use of _____ and is beyond the scope of this text.
3. Therefore, in presenting multiple regression, we focus on how statistical software can be used to obtain the estimated regression equation and other information. The emphasis will be on how to _____ the computer output rather than on how to make the multiple regression computations.

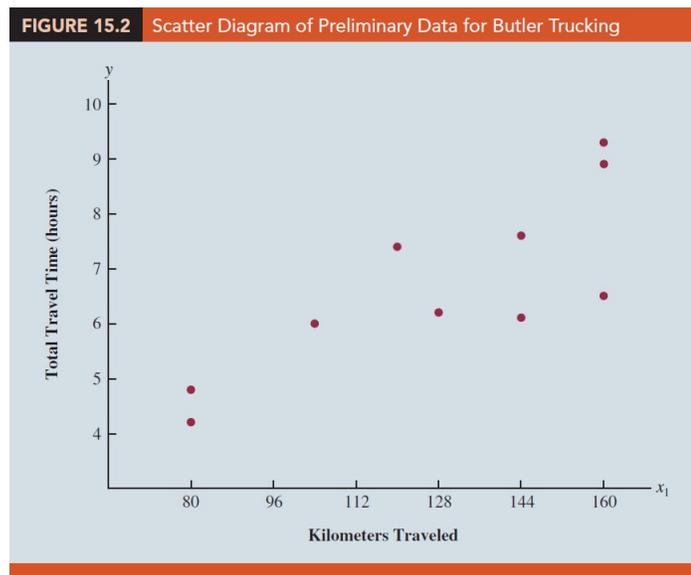
An Example: Butler Trucking Company

1. The Butler Trucking Company, an independent trucking company in southern California.
2. A major portion of Butler's business involves deliveries throughout its local area. To develop better work schedules, the managers want to predict the total daily travel time for their drivers.
 - (a) Initially the managers believed that the total daily travel time would be closely related to the number of miles traveled in making the daily deliveries.
 - (b) (Table 15.1)(Figure 15.2) A simple random sample of 10 driving assignments provided the data shown in Table 15.1 and the scatter diagram shown in Figure 15.2.

TABLE 15.1 Preliminary Data for Butler Trucking

Driving Assignment	x_1 = Kilometers Traveled	y = Travel Time (hours)
1	160	9.3
2	80	4.8
3	160	8.9
4	160	6.5
5	80	4.2
6	128	6.2
7	120	7.4
8	104	6.0
9	144	7.6
10	144	6.1

Source: PC Magazine website, April, 2015. (<https://www.pcmag.com/reviews/monitors>)



- (c) After reviewing this scatter diagram, the managers hypothesized that the simple linear regression model $y = \beta_0 + \beta_1 x_1 + \epsilon$ could be used to describe the relationship between the total travel time (y) and the number of miles traveled (x_1).
- (d) (Figure 15.3) we show statistical software output from applying simple linear regression to the data in Table 15.1. The estimated regression equation is _____
- At the 0.05 level of significance, the F value of _____ and its corresponding p -value of _____ indicate that the relationship is significant; that is, we can reject $H_0 : \beta_1 = 0$ because the p -value is less than $\alpha = 0.05$.
 - Note that the same conclusion is obtained from the t value of _____ and its associated p -value of _____.

- iii. Thus, we can conclude that the relationship between the total travel time and the number of miles traveled is _____; longer travel times are associated with more miles traveled.

FIGURE 15.3 Output for Butler Trucking with One Independent Variable

Analysis of Variance					
Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	1	15.871	15.8713	15.81	.004
Error	8	8.029	1.0036		
Total	9	23.900			

Model Summary		
S	R-sq	R-sq (adj)
1.00179	66.41%	62.21%

Coefficients				
Term	Coef	SE Coef	T-Value	P-Value
Constant	1.27	1.40	.91	.390
Kilometers	.0424	.0107	3.98	.004

Regression Equation

Time = 1.27 + .0424 Kilometers

- iv. With a coefficient of determination (expressed as a percentage) of _____, we see that _____ in travel time can be explained by the linear effect of the number of miles traveled.

3. (Table 15.2) The managers might want to consider adding a second independent variable (number of deliveries) to explain some of the remaining variability in the dependent variable.

TABLE 15.2 Data for Butler Trucking with Kilometers Traveled (x_1) and Number of Deliveries (x_2) as the Independent Variables			
Driving Assignment	x_1 = Kilometers Traveled	x_2 = Number of Deliveries	y = Travel Time (hours)
1	160	4	9.3
2	80	3	4.8
3	160	4	8.9
4	160	2	6.5
5	80	2	4.2
6	128	2	6.2
7	120	3	7.4
8	104	4	6.0
9	144	3	7.6
10	144	2	6.1

4. (Figure 15.4) Computer output with both miles traveled (x_1) and number of deliveries (x_2) as independent variables is shown in Figure 15.4. The estimated regression equation is

$$\hat{y} = \underline{\hspace{2cm}} \quad (15.6)$$

FIGURE 15.4 Output for Butler Trucking with Two Independent Variables

Analysis of Variance					
Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	2	21.6006	10.8003	32.88	.000
Error	7	2.2994	.3285		
Total	9	23.900			

Model Summary		
S	R-sq	R-sq (adj)
.573142	90.38%	87.63%

Coefficients				
Term	Coef	SE Coef	T-Value	P-Value
Constant	-.869	.952	-.91	.392
Kilometers	.03821	.00618	6.18	.000
Deliveries	.923	.221	4.18	.004

Regression Equation

Time = -.869 + .03821 Kilometers + 0.923 Deliveries

Note on Interpretation of Coefficients

- One observation can be made at this point about the relationship between the estimated regression equation with only the miles traveled as an independent variable and the equation that includes the _____ as a second independent variable.
- The value of _____ is not the same in both cases. In simple linear regression, we interpret β_1 as an estimate of the change in y for a _____ in the independent variable.
- In multiple regression analysis, we interpret each regression coefficient as follows: b_i represents an estimate of the _____ corresponding to a _____ when all other independent variables are _____.

4. Butler Trucking example

- (a) $\beta_1 = 0.06113$, an estimate of the expected increase in travel time corresponding to an increase of one mile in the distance traveled when the number of deliveries is held constant is 0.06113 hours.
- (b) $\beta_2 = 0.923$, an estimate of the expected increase in travel time corresponding to an increase of one delivery when the number of miles traveled is held constant is 0.923 hours.

☺ EXERCISES 15.2: 1, 5, 6

15.3 Multiple Coefficient of Determination

1. In simple linear regression, we showed that the total sum of squares can be partitioned into two components: the sum of squares due to regression and the sum of squares due to error. The same procedure applies to the sum of squares in multiple regression.

$$\text{SST} = \text{SSR} + \text{SSE} \quad (15.7)$$

where

SST: total sum of squares = _____.

SSR: sum of squares due to regression = _____.

SSE: sum of squares due to error = _____.

2. Example Butler Trucking problem (Figure 15.4)

FIGURE 15.4 Output for Butler Trucking with Two Independent Variables

Analysis of Variance

Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	2	21.6006	10.8003	32.88	.000
Error	7	2.2994	.3285		
Total	9	23.900			

Model Summary

S	R-sq	R-sq (adj)
.573142	90.38%	87.63%

Coefficients

Term	Coef	SE Coef	T-Value	P-Value
Constant	-.869	.952	-.91	.392
Kilometers	.03821	.00618	6.18	.000
Deliveries	.923	.221	4.18	.004

Regression Equation

Time = $-.869 + .03821$ Kilometers + 0.923 Deliveries

$SST = 23.900$, $SSR = 21.6006$, and $SSE = 2.2994$.

- With only one independent variable (number of miles traveled), the output in Figure 15.3 shows that $SST = 23.900$, $SSR = 15.871$, and $SSE = 8.029$. The value of SST is the same in both cases because it does not depend on \hat{y} , but SSR increases and SSE decreases when a second independent variable (number of deliveries) is added.
- The multiple coefficient of determination, denoted R^2 , measures the goodness of fit for the estimated multiple regression equation.

(15.8)

- The multiple coefficient of determination can be interpreted as the _____ in the dependent variable that can be explained by the estimated multiple regression equation.
- Hence, when multiplied by 100, it can be interpreted as the percentage of the variability in y that can be explained _____.

7. **Example** In the two-independent-variable Butler Trucking example, with $SSR = 21.6006$ and $SST = 23.900$, we have $R^2 = 21.6006/23.900 = 0.9038$.
8. Therefore, 90.38% of the variability in travel time y is explained by the estimated multiple regression equation with miles traveled and number of deliveries as the independent variables.
9. (Figure 15.3) the R-sq value for the estimated regression equation with only one independent variable, number of miles traveled (x_1), is 66.41%. Thus, the percentage of the variability in travel times that is explained by the estimated regression equation increases from _____ when number of deliveries is added as a second independent variable.
10. In general, R^2 always increases as independent variables are added to the model.
11. Many analysts prefer adjusting R^2 for the number of independent variables to avoid _____ the impact of adding an independent variable on the amount of variability explained by the estimated regression equation.
12. With n denoting the number of observations and p denoting the number of independent variables, the adjusted multiple coefficient of determination is computed as follows:

$$(15.9)$$

13. **Example** With $n = 10$ and $p = 2$, we have

$$R^2 = 1 - (1 - 0.9038) \frac{10 - 1}{10 - 2 - 1}$$

14. Thus, after adjusting for the two independent variables, we have an adjusted multiple coefficient of determination of 0.8763. This value (expressed as a percentage) is provided in the output in Figure 15.4 as _____.
15. If the value of R^2 is small and the model contains a large number of independent variables, the adjusted coefficient of determination can take a _____; in such cases, statistical software usually sets the adjusted coefficient of determination to _____.

☺ EXERCISES 15.3: 11, 14, 15

15.4 Model Assumptions

1. The multiple regression model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon \quad (15.10)$$

2. The assumptions about the _____ in the multiple regression model:

- (1) The error term ϵ is a random variable with mean or expected value of zero; that is, _____.

Implication: For given values of x_1, x_2, \dots, x_p , the expected, or average, value of y is given by

$$E(y) = \underline{\hspace{10em}} \quad (15.11)$$

Equation (15.11) is the _____. $E(y)$ represents the average of all possible values of y that might occur for the given values of x_1, x_2, \dots, x_p .

- (2) The variance of ϵ is denoted by σ^2 and is the same for all values of the independent variables x_1, x_2, \dots, x_p ; that is, _____.

Implication: The variance of y about the regression line equals _____ and is the same for all values of x_1, x_2, \dots, x_p .

- (3) The values of ϵ are _____.

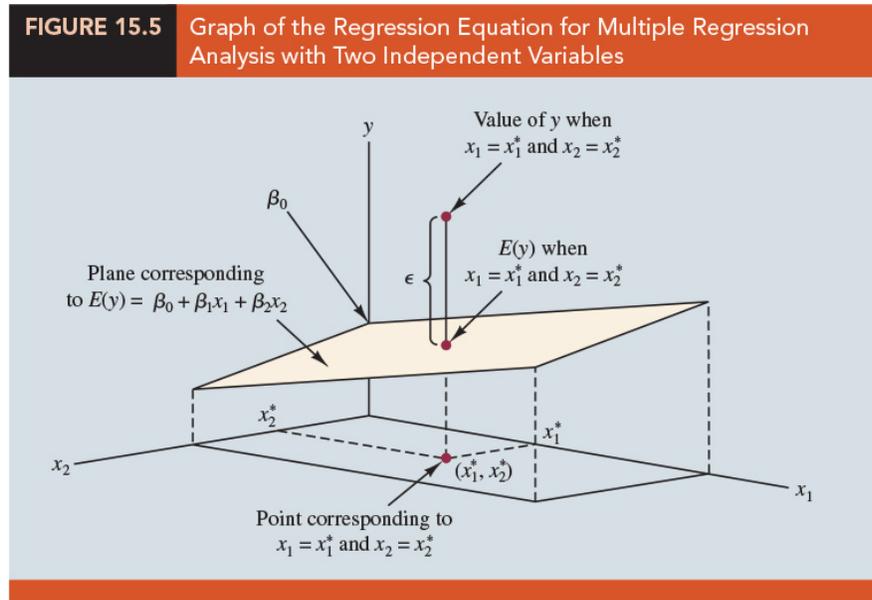
Implication: The value of ϵ for a particular set of values for the independent variables is not related to the value of ϵ for any other set of values.

- (4) The error term ϵ is a _____ random variable reflecting the deviation between the _____ value and the _____ given by $\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p$.

Implication: Because $\beta_0, \beta_1, \dots, \beta_p$ are _____ for the given values of x_1, x_2, \dots, x_p , the dependent variable y is also a _____ distributed random variable.

3. (Figure 15.5) Consider the following two-independent-variable multiple regression equation.

$$E(y) = \beta_0 + \beta_1x_1 + \beta_2x_2$$



4. Note that the value of ϵ shown is the _____ between the actual y value and the expected value of y , $E(y)$, when $x_1 = x_1^*$ and $x_2 = x_2^*$.
5. In regression analysis, the term response variable is often used in place of the term _____. Furthermore, since the multiple regression equation generates a plane or surface, its graph is called a _____.

15.5 Testing for Significance

1. In simple linear regression, both _____ and an _____ provide the same conclusion; that is, if the null hypothesis is rejected, we conclude that _____.

2. In multiple regression, the t test and the F test have different purposes.
 - (a) The F test is used to determine whether a significant relationship exists between the dependent variable and the set of _____ the independent variables; we will refer to the F test as the test for _____.
 - (b) If the F test shows an overall significance, the _____ 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 test for _____.
3. In the material that follows, we will explain the F test and the t test and apply each to the Butler Trucking Company example.

F Test

1. The hypotheses for the F test involve the parameters of the multiple regression model.

$$H_0 : \underline{\hspace{10em}}$$

$$H_a : \text{One or more of the parameters are not equal to zero}$$

2. If H_0 is rejected, the test gives us _____ to conclude that one or more of the parameters are not equal to zero and that the _____ between y and the set of independent variables x_1, x_2, \dots, x_p is _____.
3. However, if H_0 cannot be rejected, we do not have _____ to conclude that a significant relationship is present.
4. (Review)(Chapter 14)
 - (a) A mean square is a _____ divided by its corresponding degrees of freedom.
 - (b) In the multiple regression case, the total sum of squares (SST) has _____ degrees of freedom, the sum of squares due to regression (SSR) has _____ degrees of freedom, and the sum of squares due to error (SSE) has _____ degrees of freedom.

- (c) Hence, the mean square due to regression (MSR) is _____ and the mean square due to error (MSE) is _____.
 - (d) MSE provides an unbiased estimate of _____, the variance of the error term ϵ .
 - (e) If _____ is true, _____ also provides an unbiased estimate of σ^2 , and the value of MSR/MSE should be close to _____.
 - (f) However, if H_0 is false, MSR _____ σ^2 and the value of MSR/MSE becomes _____.
5. To determine how large the value of _____ must be to reject H_0 , we make use of the fact that if _____ and the _____ about the multiple regression model are _____, the sampling distribution of MSR/MSE is an _____ distribution with _____ degrees of freedom in the numerator and _____ in the denominator.

6. F test for overall significance

(a) Hypothesis:

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

H_a : One or more of the parameters are not equal to zero

(b) Test statistic:

$$\text{_____} \quad (15.14)$$

(c) Rejection rule:

- i. p -value approach: Reject H_0 if _____.
- ii. Critical value approach: Reject H_0 if _____

TABLE 15.3 ANOVA Table for a Multiple Regression Model with p Independent Variables				
Source	Sum of Squares	Degrees of Freedom	Mean Square	F
Regression	SSR	p	$MSR = \frac{SSR}{p}$	$F = \frac{MSR}{MSE}$
Error	SSE	$n - p - 1$	$MSE = \frac{SSE}{n - p - 1}$	
Total	SST	$n - 1$		

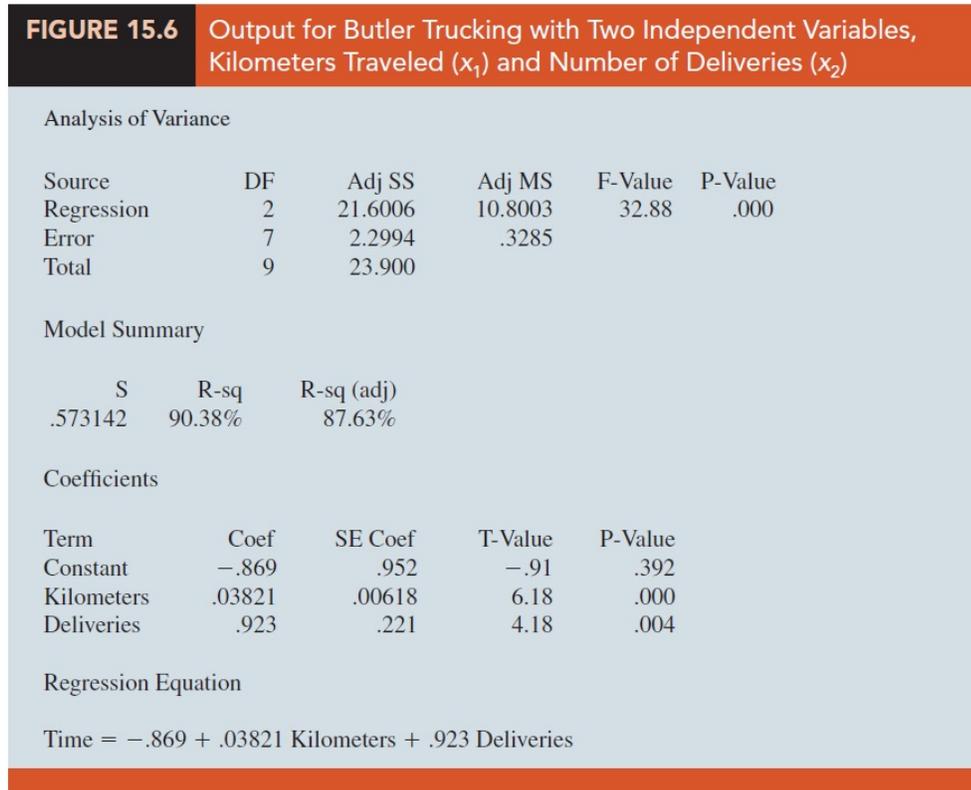
7. **Example** Butler Trucking Company

(a) Hypotheses:

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

$$H_a : \beta_1 \text{ and/or } \beta_2 \text{ is not equal to zero}$$

(b) (Figure 15.6)



- (c) $MSR = 10.8003$ and $MSE = 0.3285$, $F = \underline{\hspace{2cm}}$. Using $\alpha = 0.01$, $\underline{\hspace{2cm}}$. With $F = 32.88 > 9.55$, we reject $H_0 : \beta_1 = \beta_2 = 0$.
- (d) Using $\alpha = 0.01$, the p -value = 0.000 indicates that we can reject $H_0 : \beta_1 = \beta_2 = 0$ because the p -value is less than $\alpha = 0.01$.
- (e) Conclude that a $\underline{\hspace{2cm}}$ is present between travel time y and the two independent variables, miles traveled and number of deliveries.

***t* Test**

1. If the F test shows that the multiple regression relationship is significant, a t test can be conducted to determine the significance of each of the _____ parameters.

2. The t test for individual significance

- (a) Hypothesis: For any parameter β_i

$$H_0 : \underline{\hspace{2cm}}$$

$$H_a : \beta_i \neq 0$$

- (b) Test statistic:

$$\underline{\hspace{2cm}} \quad (15.15)$$

- (c) Rejection rule: _____

i. p -value approach: Reject H_0 if $p\text{-value} \leq \alpha$.

ii. Critical value approach: Reject H_0 if _____ or if _____.

3. In the test statistic, s_{b_i} is the estimate of the standard deviation of b_i . The value of s_{b_i} will be provided by the computer software package.

補充:

The multiple regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_{p-1} x_{p-1} + \epsilon,$$

or

$$y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_{p-1} x_{i,p-1} + \epsilon_i, \quad i = 1, \cdots, n.$$

- (a) In the matrix notation:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \quad \text{or} \quad \mathbf{y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\epsilon}_{n \times 1}.$$

- (b) Use Least-squares to fit a regression line to the data $\{\mathbf{x}_i, y_i\}_{i=1}^n$, where $\mathbf{x}_i = \{x_{i,1}, x_{i,2}, \cdots, x_{i,p-1}\}$

$$Q(\boldsymbol{\beta}) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \sum_{i=1}^n (y_i - \mathbf{x}_i \boldsymbol{\beta})^2.$$

$$\begin{aligned}\frac{\partial Q}{\partial \boldsymbol{\beta}} &= -2\mathbf{X}'(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) = \mathbf{0} \\ \Rightarrow & (\mathbf{X}'\mathbf{X})\boldsymbol{\beta} = \mathbf{X}'\mathbf{y} \\ \Rightarrow & \hat{\boldsymbol{\beta}} = \mathbf{b} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}\end{aligned}$$

(c) Variance of the sampling distribution of $b_i, i = 1, 2, \dots, p$.

$$Var(b_i) = \frac{\sigma^2}{(n-1)S_{x_i}^2(1-R_i^2)},$$

where $S_{x_i}^2$ is the sample variance of variable x_i and R_i^2 is R -square of the regression of x_i on the rest of the explanatory variables of the models (including the constant term). Note that the variance should be conditional on the observed values of the explanatory variables.

4. Example Butler Trucking Company

(a) (Figure 15.6) that shows the output for the t -ratio calculations:

$$b_1 = 0.06113, b_2 = 0.923, s_{b_1} = 0.00989, s_{b_2} = 0.221$$

(b) The test statistic for the hypotheses involving parameters β_1 and β_2 :

$$t = 0.06113/0.00989 = 6.18, \quad t = 0.923/0.221 = 4.18$$

(c) Using $\alpha = 0.01$, the p -values of _____ and _____ in the output indicate that we can reject $H_0 : \beta_1 = 0$ and $H_0 : \beta_2 = 0$. Hence, both parameters are statistically significant.

(d) Alternatively, _____. With $6.18 > 3.499$, we reject $H_0 : \beta_1 = 0$. Similarly, with $4.18 > 3.499$, we reject $H_0 : \beta_2 = 0$.

Multicollinearity

1. We use the term _____ in regression analysis to refer to any variable being used to predict or explain the value of the dependent variable.

2. The term does not mean, however, that the independent variables _____ are independent in any statistical sense. On the contrary, most independent variables in a multiple regression problem are _____ to some degree with one another.
3. **Example** Butler Trucking Example
 - (a) Butler Trucking example involves the two independent variables x_1 (miles traveled) and x_2 (number of deliveries), we could treat the miles traveled as the dependent variable and the number of deliveries as the independent variable to determine whether those two variables are themselves related.
 - (b) Compute the sample correlation coefficient $r(x_1, x_2) = 0.16$ and find that some degree of linear association between the two independent variables.
4. In multiple regression analysis, _____ refers to the correlation among the independent variables.
5. **Example** Modified Butler Trucking Example, the potential problems of multicollinearity.
 - (a) Consider a modification of the Butler Trucking example. Instead of x_2 being the number of deliveries, let x_2 denote the number of gallons of gasoline consumed. Clearly, x_1 (the miles traveled) and x_2 are related; that is, we know that the number of gallons of gasoline used depends on the number of miles traveled.
 - (b) We would conclude logically that x_1 and x_2 are highly correlated independent variables.
 - (c) Assume that we obtain the equation $\hat{y} = b_0 + b_1x_1 + b_2x_2$ and find that the F test shows the relationship to be significant. Then suppose we conduct a t test on β_1 to determine whether $\beta_1 \neq 0$, and we cannot reject $H_0 : \beta_1 = 0$. Does this result mean that travel time is not related to miles traveled? Not necessarily.
 - (d) What it probably means is that with _____, x_1 does not make a significant contribution to determining the value of y .

- (e) This interpretation makes sense in our example; if we know the amount of gasoline consumed (x_2), we do not gain much additional information useful in predicting y by knowing the miles traveled (x_1).
- (f) Similarly, a t test might lead us to conclude $\beta_2 = 0$ on the grounds that, with x_1 in the model, knowledge of the amount of gasoline consumed does not add much.
6. To summarize, in _____ for the significance of individual parameters, the difficulty caused by multicollinearity is that it is possible to conclude that _____ of the individual parameters is significantly different from zero when an _____ on the _____ multiple regression equation indicates a significant relationship.
7. Statisticians have developed several _____ for determining whether multicollinearity is high enough to cause problems.
8. According to the rule of thumb test, multicollinearity is a potential problem if the absolute value of the _____ exceeds _____ for any two of the independent variables.
9. The other types of tests are more advanced and beyond the scope of this text. If possible, every attempt should be made to avoid including independent variables that are highly correlated.
10. When multicollinearity is severe,
- (a) it is not possible to determine the separate effect of any particular independent variable on the dependent variable.
 - (b) we can have difficulty interpreting the results of t tests on the individual parameters.
 - (c) Least squares estimates may have the wrong sign.
11. 補充:
- (a) Multicollinearity in Regression Analysis: Problems, Detection, and Solutions
<https://statisticsbyjim.com/regression/multicollinearity-in-regression-analysis/>
 - (b) Multicollinearity in Regression: Why it is a problem? How to check and fix it

<https://towardsdatascience.com/multi-collinearity-in-regression-fe7a2c1467ea>

(c) Eight Ways to Detect Multicollinearity

<https://www.theanalysisfactor.com/eight-ways-to-detect-multicollinearity/>

(d) Multicollinearity (Wikipedia)

<https://en.wikipedia.org/wiki/Multicollinearity>

☺ EXERCISES 15.5: 19, 23, 24

15.6 Using the Estimated Regression Equation for Estimation and Prediction

1. The procedures for estimating the mean value of y and predicting an individual value of y in multiple regression are similar to those in regression analysis involving one independent variable.
2. We substitute the given values of x_1, x_2, \dots, x_p into the estimated regression equation and use the corresponding value of \hat{y} as the _____.
3. **Example** Butler Trucking example
 - (a) We want to use the estimated regression equation involving x_1 (miles traveled) and x_2 (number of deliveries) to develop two interval estimates:
 - i. A _____ of the mean travel time for all trucks that travel 100 miles and make two deliveries.
 - ii. A _____ of the travel time for one specific truck that travels 100 miles and makes two deliveries
 - (b) Using the estimated regression equation $\hat{y} = -0.869 + 0.06113x_1 + 0.923x_2$ with $x_1 = 100$ and $x_2 = 2$, we obtain

$$\hat{y} = \underline{\hspace{10em}}$$

Hence, the point estimate of travel time in both cases is approximately seven hours.

- (c) To develop interval estimates for the mean value of y and for an individual value of y , we use a procedure similar to that for regression analysis involving one independent variable. The formulas required are beyond the scope of the text, but statistical _____ for multiple regression analysis will often provide confidence intervals once the values of x_1, x_2, \dots, x_p are specified by the user.
- (d) (Table 15.4)

Value of x_1	Value of x_2	95% Confidence Interval		95% Prediction Interval	
		Lower Limit	Upper Limit	Lower Limit	Upper Limit
160	4	8.135	9.742	7.363	10.514
80	3	4.127	5.789	3.369	6.548
160	4	8.135	9.742	7.363	10.514
160	2	6.258	7.925	5.500	8.683
80	2	3.146	4.924	2.414	5.656
128	2	5.232	6.505	4.372	7.366
120	3	6.037	6.936	5.059	7.915
104	4	5.960	7.637	5.205	8.392
144	3	6.917	7.891	5.964	8.844
144	2	5.776	7.184	4.953	8.007
120	4	6.669	8.152	5.865	8.955

- (e) Note that the interval estimate for an individual value of y is _____ the interval estimate for the expected value of y . This difference simply reflects the fact that for given values of x_1 and x_2 we can estimate the mean travel time for all trucks with _____ than we can predict the travel time for one specific truck.

😊 **EXERCISES 15.6:** 27, 29

15.7 Categorical Independent Variables

- (a) Thus far, the examples we have considered involved _____ independent variables such as student population, distance traveled, and number of deliveries.
- (b) In many situations, however, we must work with _____ independent variables such as gender (male, female), method of payment (cash, credit card, check), and so on.

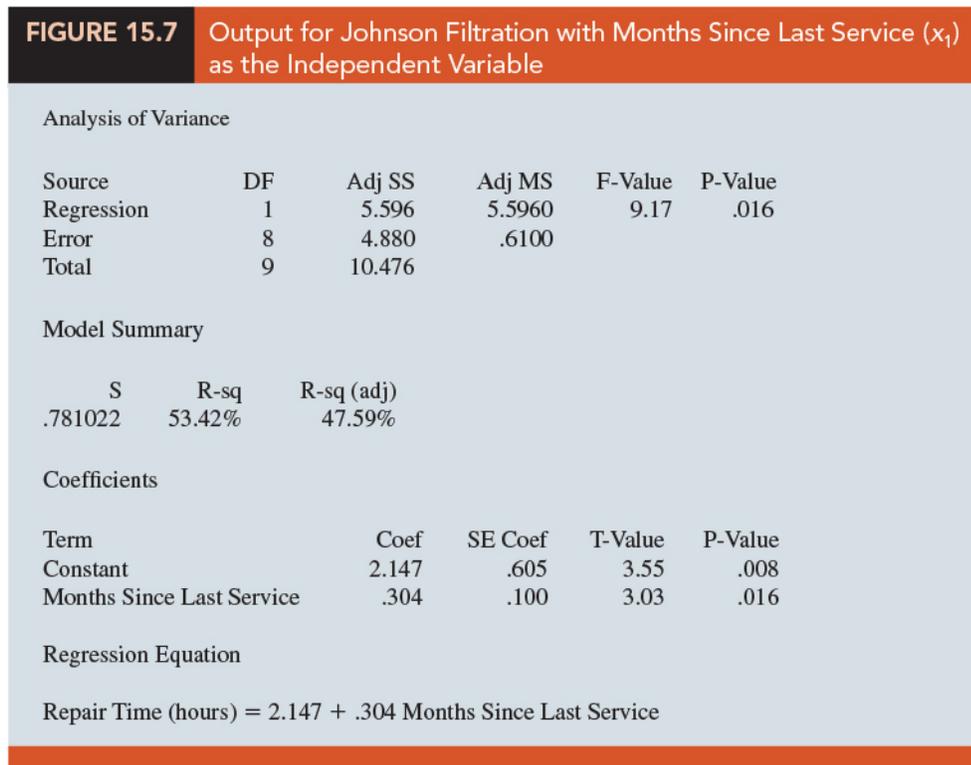
An Example: Johnson Filtration, Inc.

- (a) (Background) Johnson Filtration, Inc., provides maintenance service for water-filtration systems throughout southern Florida. Customers contact Johnson with requests for maintenance service on their water-filtration systems. To estimate the service time and the service cost, Johnson's managers want to predict the repair time necessary for each maintenance request.
- (b) (Dependent variable/Independent variables) Hence, repair time in hours is the dependent variable. Repair time is believed to be related to two factors, the number of months since the last maintenance service and the type of repair problem (mechanical or electrical).
- (c) (Data)(Table 15.5)

Service Call	Months Since Last Service	Type of Repair	Repair Time in Hours
1	2	Electrical	2.9
2	6	Mechanical	3.0
3	8	Electrical	4.8
4	3	Mechanical	1.8
5	2	Electrical	2.9
6	7	Electrical	4.9
7	9	Mechanical	4.2
8	8	Mechanical	4.8
9	4	Electrical	4.4
10	6	Electrical	4.5

- (d) (SLR) Let y denote the repair time in hours and x_1 denote the number of months since the last maintenance service. The regression model that uses only x_1 to predict y is $y = \beta_0 + \beta_1 x_1 + \epsilon$

(e) (Figure 15.7)



- i. The estimated regression equation is _____.
 - ii. At the 0.05 level of significance, the p -value of _____ for the t (or F) test indicates that the number of months since the last service is significantly related to repair time.
 - iii. R -sq = _____ indicates that x_1 alone explains _____ of the _____ in repair time.
4. To incorporate the type of repair into the regression model, we define
- $$x_2 = \begin{cases} \text{_____,} & \text{if the type of repair is mechanical} \\ \text{_____,} & \text{if the type of repair is electrical} \end{cases}$$
5. In regression analysis x_2 is called a _____ or _____.
 6. Using this dummy variable, we can write the multiple regression model as

$$y = \text{_____}$$

7. (Table 15.6) Data for the Johnson Filtration Example with Type of Repair Indicated by a Dummy Variable ($x_2 = 0$ for Mechanical; $x_2 = 1$ for Electrical)

Customer	Months Since Last Service (x_1)	Type of Repair (x_2)	Repair Time in Hours (y)
1	2	1	2.9
2	6	0	3.0
3	8	1	4.8
4	3	0	1.8
5	2	1	2.9
6	7	1	4.9
7	9	0	4.2
8	8	0	4.8
9	4	1	4.4
10	6	1	4.5

8. (Figure 15.7) Output for Johnson Filtration with Months Since Last Service (x_1) as the Independent Variable

Analysis of Variance					
Source	DF	Adj SS	Adj MS	F-Value	P-Value
Regression	2	9.0009	4.50046	21.36	.001
Error	7	1.4751	.21073		
Total	9	10.4760			
Model Summary					
S	R-sq	R-sq (adj)			
.459048	85.92%	81.90%			
Coefficients					
Term	Coef	SE Coef	T-Value	P-Value	
Constant	.930	.467	1.99	.087	
Months Since Last Service	.3876	.0626	6.20	.000	
Type of Repair	1.263	.314	4.02	.005	
Regression Equation					
Repair Time (hours) = .930 + .3876 Months Since Last Service + 1.263 Type of Repair					

- (a) The estimated multiple regression equation is

$$\text{Repair Time (hours)} = .930 + .3876 \text{ Months Since Last Service} + 1.263 \text{ Type of Repair} \quad (15.17)$$

- (b) At the 0.05 level of significance, the p -value of _____ associated with the F test (_____) indicates that the regression relationship is significant.

- (c) The t test shows that both months since last service (p -value = _____) and type of repair (p -value = _____) are statistically significant.
- (d) In addition, R -Sq = _____ and R -Sq (adj) = _____ indicate that the estimated regression equation does a good job of explaining the variability in repair times.
- (e) Thus, equation (15.17) should prove helpful in predicting the repair time necessary for the various service calls.

Interpreting the Parameters

1. The multiple regression equation for the Johnson Filtration example is

$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \quad (15.18)$$

2. Consider the case when $x_2 = 0$ (mechanical repair). Using _____ to denote the mean or expected value of repair time given a mechanical repair, we have

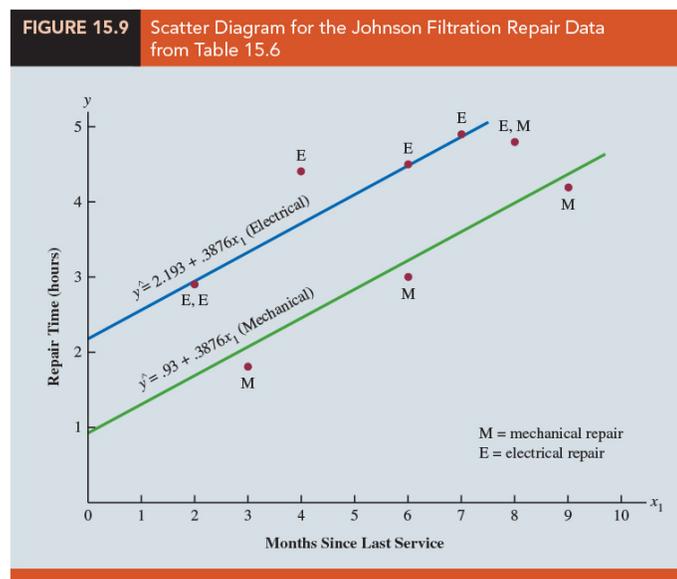
$$E(y|\text{mechanical}) = \underline{\hspace{2cm}} = \underline{\hspace{2cm}} \quad (15.19)$$

3. Similarly, for an electrical repair ($x_2 = 1$), we have

$$E(y|\text{electrical}) = \underline{\hspace{2cm}} = \underline{\hspace{2cm}} \quad (15.20)$$

4. Comparing equations (15.19) and (15.20), we see that the mean repair time is a linear function of _____ for both mechanical and electrical repairs. The slope of both equations is _____, but the _____ differs.
5. The y -intercept is _____ in equation (15.19) for mechanical repairs and _____ in equation (15.20) for electrical repairs.
6. The interpretation of β_2 is that it indicates the _____ between the _____ for an electrical repair and the mean repair time for a mechanical repair.
- (a) If _____, the mean repair time for an electrical repair will be _____ that for a mechanical repair;

- (b) if _____, the mean repair time for an electrical repair will be _____ that for a mechanical repair.
- (c) if _____, there is _____ in the mean repair time between electrical and mechanical repairs and the type of repair is _____ to the repair time.
7. Using the estimated multiple regression equation $\hat{y} = 0.93 + 0.3876x_1 + 1.263x_2$, we see that 0.93 is the estimate of β_0 and 1.263 is the estimate of β_2 .
8. Thus, when $x_2 = 0$ (mechanical repair)
- $$\hat{y} = 0.93 + 0.3876x_1 \quad (15.21)$$
- and when $x_2 = 1$ (electrical repair)
- $$\hat{y} = 0.93 + 0.3876x_1 + 1.263(1) = 2.193 + 0.3876x_1 \quad (15.22)$$
9. In effect, the use of a dummy variable for type of repair provides _____ that can be used to predict the repair time, one corresponding to mechanical repairs and one corresponding to electrical repairs.
10. In addition, with $\beta_2 = 1.263$, we learn that, on average, electrical repairs require _____ than mechanical repairs.
11. (Figure 15.9) Scatter Diagram for the Johnson Filtration Repair Data



More Complex Categorical Variables

1. If a categorical variable has k levels, $k-1$ dummy variables are required, with each dummy variable being coded as _____.
2. **Example** Suppose a manufacturer of copy machines organized the sales territories for a particular state into three regions: A, B, and C. The managers want to use regression analysis to help predict the number of copiers sold per week.
3. With the number of units sold as the dependent variable, they are considering several independent variables (the number of sales personnel, advertising expenditures, and so on).
4. Suppose the managers believe sales region is also an important factor in predicting the number of copiers sold. Because sales region is a categorical variable with three levels, A, B and C, we will need _____ dummy variables to represent the sales region. Each variable can be coded 0 or 1:

$$x_1 = \begin{cases} 1, & \text{if sales region B} \\ 0, & \text{otherwise} \end{cases}$$

$$x_2 = \begin{cases} 1, & \text{if sales region C} \\ 0, & \text{otherwise} \end{cases}$$

5. We have the following values of x_1 and x_2 :

Region	x_1	x_2
A	0	0
B	1	0
C	0	1

6. Observations corresponding to region A would be coded _____; observations corresponding to region B would be coded _____; and observations corresponding to region C would be coded _____.
7. The regression equation relating the expected value of the number of units sold, $E(y)$, to the dummy variables would be written as

$$E(y) = \underline{\hspace{2cm}}$$

8. To help us interpret the parameters β_0 , β_1 , and β_2 , consider the following three variations of the regression equation.

$$E(y|\text{region A}) = \beta_0 + \beta_1(0) + \beta_2(0) = \beta_0$$

$$E(y|\text{region B}) = \beta_0 + \beta_1(1) + \beta_2(0) = \beta_0 + \beta_1$$

$$E(y|\text{region C}) = \beta_0 + \beta_1(0) + \beta_2(1) = \beta_0 + \beta_2$$

- (a) Thus, β_0 is the mean or expected value of sales for _____;
- (b) β_1 is the _____ between the mean number of units sold in _____ and the mean number of units sold in _____;
- (c) and β_2 is the _____ between the mean number of units sold in _____ and the mean number of units sold in _____.
9. Two dummy variables were required because sales region is a categorical variable with three levels.
10. The assignment was _____. For example, we could have chosen $x_1 = 1, x_2 = 0$ to indicate region A, $x_1 = 0, x_2 = 0$ to indicate region B, and $x_1 = 0, x_2 = 1$ to indicate region C.

Region	x_1	x_2
A	1	0
B	0	0
C	0	1

In that case, β_1 would have been interpreted as the mean difference between regions A and B and β_2 as the mean difference between regions C and B.

11. The important point to remember is that when a categorical variable has k levels, $k-1$ dummy variables are required in the multiple regression analysis. Thus, if the sales region example had a fourth region, labeled D, three dummy variables would be necessary. For example, the three dummy variables can be coded as follows.

$$x_1 = \begin{cases} 1, & \text{if sales region B} \\ 0, & \text{otherwise} \end{cases} \quad x_2 = \begin{cases} 1, & \text{if sales region C} \\ 0, & \text{otherwise} \end{cases} \quad x_3 = \begin{cases} 1, & \text{if sales region D} \\ 0, & \text{otherwise} \end{cases}$$

☺ EXERCISES 15.7: 32, 34, 35