Statistics for Management and Economics

OPRE 6301 (Statistics and Data Analysis)
by
Ayşegül Toptal Bilhan

Chapter 17 (Multiple Regression)

Multiple Linear Regression

- The simple linear regression model was used to analyze how the dependent variable y is related to one independent variable x.
- Multiple regression allows for any number of independent variables.
- Our model is now:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$



OPRE 6301 2

Multiple Linear Regression (Coefficients)

Let us consider a regression model with two independent variables. That is,

$$\hat{y}_i = b_0 + b_1 x_{1i} + b_2 x_{2i}$$

• The coefficient estimators are:

$$b_1 = \frac{s_y(r_{x_1y} - r_{x_1x_2}r_{x_2y})}{s_{x_1}(1 - r_{x_1x_2}^2)}, b_2 = \frac{s_y(r_{x_2y} - r_{x_1x_2}r_{x_1y})}{s_{x_2}(1 - r_{x_1x_2}^2)}, b_0 = \overline{y} - b_1\overline{x}_1 - b_2\overline{x}_2$$

 r_{x_1y} is the sample correlation between X_1 and Y

 $r_{\chi_2 \gamma}$ is the sample correlation between X_2 and Y

 $r_{x_1x_2}$ is the sample correlation between X_1 and X_2

 s_{x_1} is the sample standard deviation of X_1

 s_{x_2} is the sample standard deviation of X_2

 s_y is the sample standard deviation of Y OPRE 6301



Multiple Linear Regression (Cont'd)

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.



OPRE 6301

Multiple Linear Regression (Cont'd)

Example (Cont'd): Programmer Salary Survey

Exper. (Yrs.) 4 7 1 5 8 10 0 1	78 100 86 82 86 84 75 80	Salary (\$000s) 24.0 43.0 23.7 34.3 35.8 38.0 22.2 23.1	9 2 10 5 6 8 4 6	Test Score 88 73 75 81 74 87 79 94	Salary (\$000s) 38.0 26.6 36.2 31.6 29.0 34.0 30.1 33.9
10	84	38.0	8 4	87	34.0
0	75	22.2		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 Linear Regression (Cont'd)

Example (Cont'd): Programmer Salary Survey

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 (\$000)

 x_1 = years of experience

 x_2 = score on programmer aptitude test



OPRE 6301

6

Multiple Linear Regression (Cont'd)

Example (Cont'd): Programmer Salary Survey

Excel's Regression Equation Output

	Α	В	С	D	E
38					
39			Std. Err.		
40	Intercept	3.17394	6.15607	0.5156	0.61279
41	Experience	1.4039	0.19857	7.0702	1.9E-06
42	Test Score	0.25089	0.07735	3.2433	0.00478
43					

Note: Columns F-I are not shown.

E[SALARY]= 3.174 + 1.404(EXPER) + 0.251(SCORE)

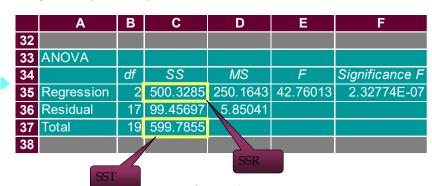


OPRE 6301

Multiple Linear Regression (Cont'd)

Example (Cont'd): Programmer Salary Survey

Excel's Regression Equation Output



 $R^2 = SSR/SST$, therefore, $R^2 = 500.3285/599.7855 = .83418$

UT DALLAS OPRE 6301

Multiple Linear Regression (Estimation of Error Variance)

Let s_{ε}^2 be an unbiased estimate of error variance σ_{ε}^2 .

$$s_{\varepsilon}^2 = \frac{SSE}{n-k-1}$$

If number of independent variables is large relative to sample size n, we use Adjusted R².

Adjusted
$$R^2 = 1 - \frac{SSE/(n-k-1)}{SST/(n-1)}$$



OPRE 6301

Multiple Linear Regression (Coefficient Standard Errors)

Square roots of the below variance estimators $(s_{b_1}, s_{b_2}, ...)$ are called coefficient standard errors.

$$s_{b_1}^2 = \frac{s_{\varepsilon}^2}{(n-1)s_{x_1}^2 \left(1 - r_{x_1 x_2}^2\right)}, \ s_{b_2}^2 = \frac{s_{\varepsilon}^2}{(n-1)s_{x_2}^2 \left(1 - r_{x_1 x_2}^2\right)}$$



OPRE 6301 10

Multiple Linear Regression (F-test for Overall Significance)

To test the validity of the regression model, we have

$$H_0$$
: $\beta_1 = \beta_2 = \beta_3 = ... = \beta_k = 0$

 H_a : At least one β_i is not equal to 0.

Test Statistic:
$$F = \frac{MSR}{MSE} = \frac{SSR/k}{SSE/(n-k-1)}$$

Rejection Region: $F > F_{\alpha,k,n-k-1}$



OPRE 6301 11

Multiple Linear Regression (Testing the Validity of the Model)

SSE	S_{ε}	R ²	F	Assessment of Model
0	0	1	8	Perfect
small	small	close to 1	large	Good
large	large	close to 0	small	Poor
$\sum (y_i - \overline{y})^2$	$\sqrt{\frac{\sum (y_i - \bar{y})^2}{n - k - 1}}$	0	0	Invalid

Once we're satisfied that the model fits the data as well as possible, and that the required conditions are satisfied, we can interpret and test the individual coefficients and use the model to predict and estimate...



Multiple Linear Regression (Tests of Hypotheses for Regression Coefficients)

To test the null hypothesis

$$H_0$$
: $\beta_j = \beta^*$

against the two-sided alternative

$$H_1: \beta_i \neq \beta^*$$

the decision rule is as follows:

$$\text{Reject H}_0 \text{ if } \frac{b_j - \beta^*}{s_{b_j}} > t_{n-k-1,\alpha/2} \text{ or } \frac{b_j - \beta^*}{s_{b_j}} < -t_{n-k-1,\alpha/2}$$



OPRE 6301 13