



7. Multiple Regression Analysis: The Problem of Analysis

1Y & LX

Three-Variable Model: Notation and Assumptions

Let us consider the following three-variable PRF as:

$$Y_i = \beta_0 + \beta_1 X_{2i} + \beta_2 X_{3i} + u_i$$

where

Y_i is the dependent variable (regressand)

X_{2i} and X_{3i} are the regressors or the explanatory variables

u_i is the stochastic disturbance term

Remark: the subscript i is denoted the observation i from our sample data.

In case our data are time series, the subscript t will denote the t observation.

β_0 means the average value of Y when X_2 and X_3 are set equal to zero

β_1 and β_2 are called the partial regression coefficients.

We will talk about the meaning of β_1 and β_2 shortly after knowing the assumptions of the classical linear regression model (CLRM)

Y_i	X_{2i}	X_{3i}
⋮	⋮	⋮
t	Y_t	X_{2t}
	X_{3t}	

Chapter 7. Multiple Regression Analysis: The Problem of Analysis

Under the CLRM, we assume:

1. Zero mean value of u_i

$$E(u_i | X_{2i}, X_{3i}) = 0 \quad \equiv \text{Mean value of residuals} = 0$$

2. No serial correlation

$$\text{Cov}(u_i, u_j) = 0 \quad \text{for all } i \neq j$$

3. Homoscedasticity

$$\text{Var}(u_i) = \sigma^2$$

4. Zero covariance between u_i and each X variable, or

$$\text{Cov}(u_i, X_{2i}) = 0 ; \text{Cov}(u_i, X_{3i}) = 0$$

5. No specification bias or

The model is correctly specified.

6. No exact collinearity between the X variables or

Specification bias comes from

- 1) Add redundant variable(s)
- 2) Omit important variable(s) in our model ("Exclusion bias")
- 3) Use wrong functional forms

Example of an exact linear relationship between X s

$$X_{3i} = 2 \cdot X_{2i}$$

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + u_i$$

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 (2 \cdot X_{2i}) + u_i$$

$$Y_i = \beta_1 + (\beta_2 + 2\beta_3) X_{2i} + u_i$$

Observe that Partial Effect of X_{2i} on Y and of X_{3i} on Y cannot be estimated when X_{2i} and X_{3i} form an exact linear relationship

Exact linear relationship

Ex: X_{3i}	X_{2i}
2	1
4	2
6	3

Inexact linear relationship

X_{3i}	X_{2i}	$X_{3i} = 5 \cdot X_{2i} + v_i$
10	2	$10 = 5 \cdot 2 + 0$
22	4	$22 = 5 \cdot 4 + 2$
28	6	$28 = 5 \cdot 6 - 2$
40	8	$40 = 5 \cdot 8 + 0$

$$X_{3i} = 2 \cdot X_{2i}$$

$$X_{3i} - 2 \cdot X_{2i} = 0$$

With inexact linear relationship between X_{2i} and X_{3i} , OLS regression can be estimated, i.e., we can obtain β_2 and β_3 .

In general form: $\text{Lamda}2 \cdot X_{2i} + \text{Lamda}3 \cdot X_{3i} = 0$

Here, in this example, $\text{Lamda}2 = -2$, $\text{Lamda}3 = 1$

Keypoint: If $\text{Lamda}2$ and $\text{Lamda}3$ are founded to make the above equation "TRUE", then it is confirmed that X_2 and X_3 has "EXACT" linear relationship! (We call this problem as "Perfect collinearity")

If this problem occurs, regression cannot be estimated.

7.1 OLS Estimation of the Partial Regression Coefficients

By the above assumptions, we can find out the conditional expectation of Y_i :

Take $E[\cdot | X_{2i}, X_{3i}]$ to the PRF:

$$E[Y_i | X_{2i}, X_{3i}] = \beta_0 + \beta_1 X_{2i} + \beta_2 X_{3i} + E[u_i | X_{2i}, X_{3i}]$$

The meaning of partial coefficients:

$$\frac{d E[Y_i | X_{2i}, X_{3i}]}{d X_{2i}} = \beta_1$$

$$\frac{d E[Y_i | X_{2i}, X_{3i}]}{d X_{3i}} = \beta_2$$

holding X_3 constant, when X_2 changes by 1 unit, Y , on average, will change by β_1 unit.

tell us partial effect of X_2 on Y

$$\frac{d E[Y_i | X_{2i}, X_{3i}]}{d X_{2i}} = \beta_2$$

when X_2 changes by 1 unit, Y , on average, will change by β_2 unit. } partial effect of X_2 on Y

7.1 OLS Estimation of the Partial Regression Coefficients

In order to find the OLS estimators, we need to write down the sample regression function (SRF) corresponding to the PRF:

Goal: Find $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$ that $\min \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \min \sum_{i=1}^n (Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \hat{\beta}_3 X_{3i})^2$

Trick: $\min \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \min \sum_{i=1}^n (Y_i - \hat{\beta}_1 X_{1i} - \hat{\beta}_2 X_{2i} - \hat{\beta}_3 X_{3i})^2$

FOCs:

$$\frac{\partial \sum u_i^2}{\partial \hat{\beta}_1} = 0 \quad \text{--- (1)}$$

$$\frac{\partial \sum u_i^2}{\partial \hat{\beta}_2} = 0 \quad \text{--- (2)}$$

$$\frac{\partial \sum u_i^2}{\partial \hat{\beta}_3} = 0 \quad \text{--- (3)}$$

Three equations; w/ three unknowns \rightarrow we can solve for $\hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3$

holding X_2 constant, when X_3 changes by 1 unit, Y , on average, will change by β_3 unit. } partial effect of X_3 on Y

Chapter 7. Multiple Regression Analysis: The Problem of Analysis

From the FOC, we then get the normal equations:

$$\begin{aligned} Y &= \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 \\ \sum Y X_1 &= \beta_1 \sum X_1^2 + \beta_2 \sum X_1 X_2 + \beta_3 \sum X_1 X_3 \\ \sum Y X_2 &= \beta_1 \sum X_1 X_2 + \beta_2 \sum X_2^2 + \beta_3 \sum X_2 X_3 \\ \sum Y X_3 &= \beta_1 \sum X_1 X_3 + \beta_2 \sum X_2 X_3 + \beta_3 \sum X_3^2 \end{aligned}$$

We therefore get:

$$\hat{\beta}_1 = \frac{\sum X_2 X_3 \sum Y X_1 - \sum X_1 X_2 \sum Y X_3 - \sum X_1 X_3 \sum Y X_2}{\sum X_1^2 \sum X_2 X_3 - \sum X_1 X_2 \sum X_1 X_3 - \sum X_1 X_3 \sum X_2 X_3}$$

Variance and Standard Errors of OLS Estimators

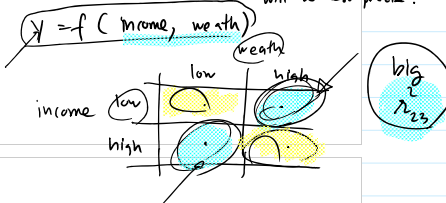
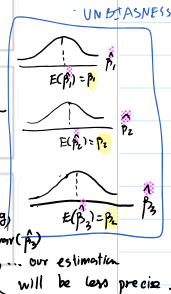
$$\text{var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum X_1^2 - \frac{(\sum X_1)^2}{n}}$$

$t = \frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)}$

$\text{se}(\hat{\beta}_1) = \frac{\sigma}{\sqrt{\sum X_1^2 - \frac{(\sum X_1)^2}{n}}}$

partial correlation coefficient bet. X_2 & X_3

if r_{23}^2 is big, $\text{var}(\hat{\beta}_2) \text{ and } \text{var}(\hat{\beta}_3)$ will be large, ... our estimation will be less precise.



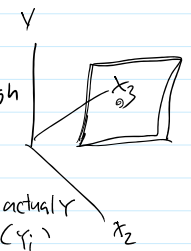
7.2 Properties of OLS Estimators

$$\text{variance of our residuals} = \frac{\sum u_i^2}{n-3} = \frac{RSS}{n-3} = \hat{\sigma}_u^2 \text{ (full description)}$$

7.2 Properties of OLS Estimators

- The three-variable regression model, the three-variable regression plane, passes through the means $\bar{Y}, \bar{X}_2, \bar{X}_3$
- The mean value of the estimated Y_i (\hat{Y}_i) is equal to the mean value of actual Y (Y_i) (see the proof in the textbook) $\hat{Y}_i = Y_i$
- with deviation forms: $y_i = \hat{y}_i + \hat{u}_i$
 $y_i = \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \hat{u}_i$

$\hat{\sigma}_u^2 \rightarrow \text{var}(\hat{\beta}_i) \rightarrow \text{se}(\hat{\beta}_i) \rightarrow t_{\hat{\beta}_i} \rightarrow \text{hypothesis testing}$



$$\hat{u}_i = y_i - \hat{\beta}_2 X_{2i} - \hat{\beta}_3 X_{3i}$$

The Multiple Coefficient of Determination R^2 and the Multiple Coefficient of Correlation R

In this section, we will study how to measure the proportion of the variation in Y explained by the variables X_2 and X_3 jointly. This is the same concept of r^2 that we have learned before.

The quantity that gives this information is known as the multiple coefficient of determination and is denoted by R^2 .

To derive R^2 , we first write down the following equation:

$$\hat{Y}_i = \beta_0 + \beta_2(X_i - \bar{X}_2) + \beta_3(X_i - \bar{X}_3) + \hat{u}_i \quad (7.1)$$

STANDARD FORM OF OUR SRF

where \hat{Y}_i is the estimated value of Y_i from the fitted regression line and is an estimator of true $E(Y_i|X_2, X_3)$.

7.1 may be written as

$$Y_i = \beta_0 + \beta_2 X_i + \beta_3 X_i + \hat{u}_i = \hat{Y}_i + \hat{u}_i \quad (7.2)$$

DEVIATION FORM OF OUR SRF

Squaring 7.2 on both sides and summing over the sample values, we obtain

$$\sum Y_i^2 = \sum \hat{Y}_i^2 + \sum \hat{u}_i^2 + 2 \sum \hat{Y}_i \hat{u}_i$$

(7.3)

TSS = ESS + RSS

Recall that $TSS = ESS + RSS$ → Variation in Y_i that is left unexplained.

Baseline or Total errors when using sample mean to estimate Y_i

Variation in Y_i that could be explained by the helps of X_{2i} and X_{3i}

$$TSS = ESS + RSS$$

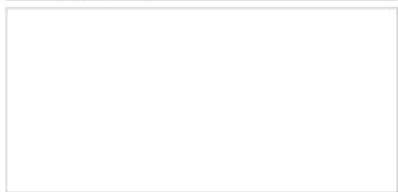
$$1 = \frac{ESS}{TSS} + \frac{RSS}{TSS}$$

$R^2 =$ Proportional variation in Y that can be explained by the regression (i.e., by using X_2 and X_3)

$$R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum \hat{u}_i^2}{\sum y_i^2} = 1 - \frac{\sum (Y_i - \hat{Y}_i)^2}{\sum (Y_i - \bar{Y})^2}$$

Residual sum of (RSS) squares

Total sum of squares (TSS)



$$R^2 = \frac{ESS}{TSS} = \frac{\beta_2 \sum X_{2i} Y_i + \beta_3 \sum X_{3i} Y_i}{\sum Y_i^2}$$

"R² in deviation form" (7.4)

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_2 X_{2i} + \hat{\beta}_3 X_{3i} + \dots + \hat{\beta}_k X_{ki}$$

Regressors = $k-1$ variables

Define R_j^2 as the R^2 in the regression of X_j on the remaining regressors, ($k-1-i = k-2$)

The three-on-one variable analogue of r is the coefficient of multiple correlation, denoted by R , and it is a measure of the degree of association between Y and all the explanatory variables jointly. Although r can be positive or negative, R is always taken to be positive.

$$\text{Var}(\hat{\beta}_j) = \frac{\sigma^2}{\sum_{i \neq j} X_i^2}$$

$\uparrow R_j^2 \rightarrow \uparrow \text{var}(\hat{\beta}_j) \rightarrow \uparrow \text{s.e.}(\hat{\beta}_j) \rightarrow \hat{t} = \frac{\hat{\beta}_j}{\text{s.e.}(\hat{\beta}_j)}$

Ex: $R_2^2 \rightarrow X_2 = \alpha_1 + \alpha_2 X_3$

$R_3^2 \rightarrow X_3 = \alpha_1 + \alpha_2 X_2$

$H_0: \beta_j = 0$

$H_a: \beta_j \neq 0$

Recall that $\sum \hat{u}_i^2 = \sum y_i^2 - \hat{\beta}_2 \sum X_{2i} y_i - \hat{\beta}_3 \sum X_{3i} y_i$

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum \hat{u}_i^2$$

↓ TSS ↓ ESS ↓ RSS

$$\sum y_i^2 = \sum \hat{y}_i^2 + \sum y_i^2 - \hat{\beta}_2 \sum X_{2i} y_i - \hat{\beta}_3 \sum X_{3i} y_i$$

↓ ESS'

$$\sum \hat{y}_i^2 = \hat{\beta}_2 \sum X_{2i} y_i + \hat{\beta}_3 \sum X_{3i} y_i$$

7.2.1 R^2 and the Adjusted R^2

It should be noted that R^2 is a nondecreasing function of the number of explanatory variables. Thus, when the number of regressors increases, R^2 almost invariably increases and never decreases. In other words, an additional X variable will not decrease R^2 .

To explain this fact, let us write down the definition of R^2 again:

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum \hat{u}_i^2}{\sum Y_i^2} \quad (7.5)$$

Therefore, in comparing two regression models with the same dependent variable but differing number of X variables, one should be very wary of choosing the model with the highest R^2 .

In light of comparing two R^2 terms, we have to take into account the number of X variables present in the model. To achieve this goal, we can consider the alternative coefficient of determination, which is as follows:

Adjusted R^2 : $R^2 = 1 - \frac{\sum \hat{u}_i^2 / (n-k)}{\sum Y_i^2 / (n-1)}$

For $k=3$

In light of comparing two R^2 terms, we have to take into account the number of X variables present in the model. To achieve this goal, we can consider the alternative coefficient of determination, which is as follows:

Adjusted R^2 :
$$\bar{R}^2 = 1 - \frac{\sum \hat{u}_i^2 / (n-k)}{\sum y_i^2 / (n-1)}$$

for $k=3$

General form:
$$\bar{R}^2 = 1 - \frac{\sum \hat{u}_i^2 / (n-k)}{\sum y_i^2 / (n-1)}$$

k = the number of parameters in the model including the intercept term.
 n = the number of observations in the sample data.

The above equation is known as the **adjusted R^2** , denoted by \bar{R}^2 . The term adjusted means adjusted for the df associated with the sums of squares entering into 7.5.

7.2 Properties of OLS Estimators 125

We can rewrite the adjusted R^2 as:

$$\bar{R}^2 = 1 - \frac{\hat{\sigma}^2}{S_Y^2}$$
 where $\hat{\sigma}^2 = \frac{\sum \hat{u}_i^2}{n-k}$ → unbiased estimator of σ^2
 $S_Y^2 = \frac{\sum y_i^2}{n-1}$ → sample variance of Y

We can also get the equation which shows the relationship between R^2 and \bar{R}^2

$$R^2 = 1 - \frac{\sum \hat{u}_i^2}{\sum y_i^2} \cdot \frac{(n-1)}{(n-k)}$$

$$\bar{R}^2 = 1 - [1 - R^2] \cdot \frac{(n-1)}{(n-k)}$$

As $k > 1$ → $\frac{(n-1)}{(n-k)} > 1$

(NOTE) \bar{R}^2 cannot be negative.

Besides R^2 and \bar{R}^2 as goodness of fit measures, other criteria are often used to judge the adequacy of a regression model. Two of these are Akaike's Information criterion and Amemiya's Prediction criteria, which are used to select between competing models. We will discuss these criteria in greater detail later.

$\frac{(n-1)}{(n-5)} = \frac{(n-1)}{(n-3)} = 1$
 $\frac{10-1}{10-3} = \frac{9}{7} > 1$

Result: $\bar{R}^2 < R^2$ ***

Remarks

When comparing two R^2 values, make sure that

- ① Dependent variables of the two models is the same
- ② Sample size must be the same

Ex:

$\ln Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i}$ ①
 $Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i}$ ②

① and ② cannot be directly compared by using R^2 b/c the dependent variables are NOT the same!



8. Multiple Regression Analysis: The Problem of Inference

In this chapter, we will extend the ideas of interval estimation and hypothesis testing developed there to models involving three or more variables.

$$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + u_i$$

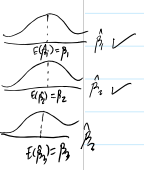
We have already known that if our objective is to do interval estimation and hypothesis testing, we need to assume that the u_i follow the normal distribution with zero mean and constant variance σ^2 .

With the normality assumption and the CLRM assumptions, we know that:

- [1] The OLS estimations of partial regression coefficients are best linear unbiased estimators (BLUE).
- [2] The estimators $\hat{\beta}_1, \hat{\beta}_2,$ and $\hat{\beta}_3$ are normally distributed with means equal to true $\beta_1, \beta_2,$ and β_3 and variances are following:

$$var(\hat{\beta}_1) = \frac{\left[\frac{1}{n} + \frac{X_2^2 \sum X_3^2 + X_3^2 \sum X_2^2 - 2X_2 X_3 \sum X_2 X_3}{\sum X_2^2 \sum X_3^2 - (\sum X_2 X_3)^2} \right] \cdot \sigma^2}{\sum X_2^2 \sum X_3^2 - (\sum X_2 X_3)^2}$$

 $se(\hat{\beta}_1) = \sqrt{var(\hat{\beta}_1)}$



$$\text{var}(\hat{\beta}_1) = \frac{\sum x_{1i}^2}{(\sum x_{1i})(\sum x_{1i}) - (\sum x_{2i}x_{1i})^2} \cdot \sigma^2$$

$$\text{var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum x_{1i}^2(1-r_{12}^2)}$$

$$\text{se}(\hat{\beta}_1) = \sqrt{\text{var}(\hat{\beta}_1)}$$

$$\text{var}(\hat{\beta}_1) = \frac{\sum x_{1i}^2}{(\sum x_{1i})(\sum x_{1i}) - (\sum x_{2i}x_{1i})^2} \cdot \sigma^2$$

$$\text{var}(\hat{\beta}_1) = \frac{\sigma^2}{\sum x_{1i}^2(1-r_{12}^2)}$$

$$\text{se}(\hat{\beta}_1) = \sqrt{\text{var}(\hat{\beta}_1)}$$

Moreover, $\frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)}$ follows the t^2 distribution with $n-3$ df. We can also show that, if we replace the true σ^2 by its unbiased estimator s^2 in the computation of the standard errors, we then get

$$t = \frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)}$$

$$t = \frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)}$$

$$t = \frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)}$$

follows the t distribution with $n-3$ df. $n-k = n-3$

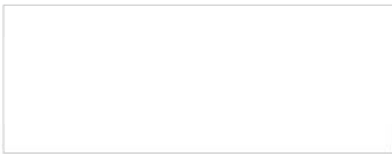
Example Consider the following regression:

$$\log(\widehat{\text{salary}}) = 4.32 + 0.280 \log(\widehat{\text{sales}}) + 0.0174 \text{ROE} + 0.00024 \text{ROS} \quad (8.1)$$

$R^2 = 0.283$

where
 salary = salary of CEO (thousand dollars)
 sales = annual firm sales
 ROE = return on equity in percent
 ROS = return on firm's stock

Interpret the partial regression coefficients



Questions What about the statistical significance of the observed results?

For the coefficient of $\log(\widehat{\text{sales}})$ of 0.280, is this coefficient statistically significant different from zero?

For the coefficient of ROE of 0.0174, is this coefficient statistically significant different from zero?

For the coefficient of ROS of 0.00024, is this coefficient statistically significant different from zero?

Are these three coefficients statistically significant?

To answer these questions, we have to learn the basics of hypothesis testing.

8.1 Hypothesis Testing About Individual Regression Coefficients

We can use the t -test to test a hypothesis about any individual partial regression coefficient.

8.1.1 Two-tail test:

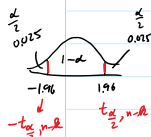
Let us assume that
 $H_0: \beta_1 = 0$ (sales do not affect CEO salary)
 $H_1: \beta_1 \neq 0$ (sales do affect CEO salary)

suppose $\alpha = 0.05$

① compute t : $t = \frac{\hat{\beta}_1 - \beta_1}{\text{se}(\hat{\beta}_1)} = \frac{0.280 - 0}{0.035} = 8.00$

② find critical t : $n = 209$ firms
 $df = n - k = 209 - 4 = 205$

$t_{\frac{\alpha}{2}, n-k} = t_{0.025, 205} = 1.96$



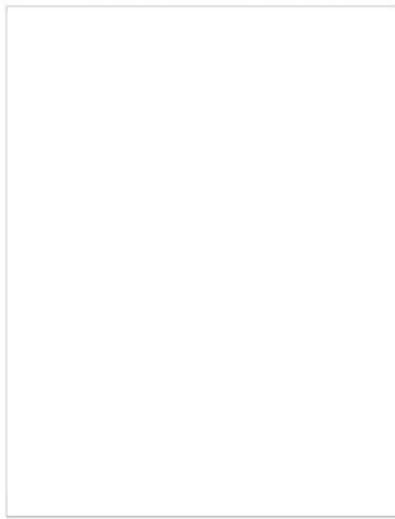
$$t_{\frac{\alpha}{2}, n-k} = t_{\frac{0.05}{2}, 205} = 1.96$$

③ Compare t w/ critical t : Since $t >$ critical t , then we can reject the null hypothesis that $H_0: \beta_2 = 0$ in favor of the alternative hypothesis that $H_a: \beta_2 \neq 0$.

In other words, $\hat{\beta}_2$ is statistically significant different from zero at 95%.

130

Chapter 8. Multiple Regression Analysis: The Problem of Inference



8.1 Hypothesis Testing About Individual Regression Coefficients

131

8.1.2 One-tail tests:

Let us consider that

$$H_0: \beta_2 \leq 0$$

$$H_1: \beta_2 > 0$$

(sales do have a positive impact on CEO salary)

$$[D-I-Y]$$



$$t_{\alpha, df} = t_{0.05, 205} = 1.645$$

8.2 Testing The Overall Significance of the Sample Regression [Joint Test]

In the previous section, we test the significance of the estimated partial regression coefficients individually, that is under the separate hypothesis that each true population partial regression coefficient was zero. But now we are interested in testing β_1, β_2 and β_3 are jointly or simultaneously equal to zero. In other words, we would like to test the following hypothesis:

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

H_A : At least one of our explanatory variables is NOT equal to zero.

In order to reach this goal, we have to learn the following test.

The Analysis of Variance Approach to Testing the Overall Significance of an Observed Multiple Regression: The F-Test

The joint hypothesis can be tested by the Analysis of Variance (ANOVA) which can be demonstrated as follows:

Table 8.1: ANOVA Table for the three-variable regression model

Source of variation	Sum of Square SS	df	Mean Sum of Square MSS
Due to regression (ESS)	$\hat{\beta}_2 \sum y_i x_{2i} + \hat{\beta}_3 \sum y_i x_{3i}$	2	$(\hat{\beta}_2 \sum y_i x_{2i} + \hat{\beta}_3 \sum y_i x_{3i}) / 2$
Due to residuals (RSS)	$\sum \hat{u}_i^2$	$n - 3$	$\sum \hat{u}_i^2 / (n - 3)$
TSS	$\sum y_i^2$	$n - 1$	$\sum y_i^2 / (n - 1)$

$TSS = ESS + RSS$
 $= \hat{\beta}_2 \sum y_i x_{2i} + \hat{\beta}_3 \sum y_i x_{3i} + \sum \hat{u}_i^2$

For Joint test

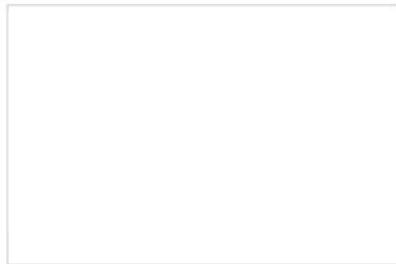
Let's say $H_0: \beta_2 = \beta_3 = 0$ Both X_2 and X_3 could not jointly explain variation in Y

$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + u_i$

Let's compute F-statistic: $F = \frac{ESS / df}{RSS / df} = \frac{(\hat{\beta}_2 \sum y_i x_{2i} + \hat{\beta}_3 \sum y_i x_{3i}) / 2}{\sum \hat{u}_i^2 / (n - 3)}$

$k - 1 = 3 - 1$

b/c, to compute $\sum \hat{u}_i^2$, we have to estimate $\hat{\beta}_1, \hat{\beta}_2$, and $\hat{\beta}_3$ which consume 3 degrees of freedom so, $df = n - 3$.



Decision Rule Given the k-variable regression model:

$Y_i = \beta_1 + \beta_2 X_{2i} + \beta_3 X_{3i} + \dots + \beta_k X_{ki} + u_i$

To test the hypothesis

$H_0: \beta_2 = \beta_3 = \dots = \beta_k = 0$

(i.e., all slope coefficients are simultaneously zero) versus

H_A : Not all slope coefficients are simultaneously zero

If $F > F_{\alpha}(k-1, n-k)$, we reject H_0 ; otherwise we cannot reject it, where $F_{\alpha}(k-1, n-k)$ is the critical F value at the α level of significance and (k-1) numerator and (n-k) denominator df.

An important relationship between R^2 and F

$$F = \frac{ESS / (k-1)}{RSS / (n-k)}$$

$$= \frac{ESS}{RSS} \cdot \frac{(n-k)}{(k-1)}$$

$$= \frac{ESS}{TSS - ESS} \cdot \frac{(n-k)}{(k-1)}$$

$$= \frac{\frac{ESS}{TSS}}{\frac{TSS - ESS}{TSS}} \cdot \frac{(n-k)}{(k-1)}$$

$$F = \frac{R^2}{1 - R^2} \cdot \frac{(n-k)}{(k-1)}$$

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

$$TSS = ESS + RSS$$

$$R^2 = \frac{ESS}{TSS}$$

- ① If $R^2 = 0$, then $F = 0$
- ② The higher the R^2 , the greater the F-value.
- ③ In limit when $R^2 \rightarrow 1$, F becomes infinite.

EX: $H_0: \beta_2 = \beta_3 = \beta_4 = 0$

H_a : Not all slope coefficients are zero simultaneously

STEP 1: Compute F-statistic:

$$F = \frac{R^2 / (k-1)}{(1 - R^2) / (n-k)} = \frac{0.28 / (4-1)}{(1 - 0.28) / (209-4)} = 26.97/2$$

STEP 2: Find critical F-stat:

$$\alpha = 0.05 \rightarrow F_{\alpha, k-1, n-k} = F_{0.05, 4-1, 209-4} = F_{0.05, 3, 205} = 2.65$$

STEP 3: Compare \hat{F} w/ critical F-stat:

Since $\hat{F} = 26.97/2 > F_{critical} = 2.65$, then we may comfortably reject $H_0: \beta_2 = \beta_3 = \beta_4 = 0$ in favor of H_a at 95% confidence level.

sales
ROE
RO9

$k = \#$ of explanatory variable including intercept.

8.2 Testing The Overall Significance of the Sample Regression

Table 8.1 ANOVA Table in Terms of $R^2 = (k-1) = (3-1)$

Source of variation	Sum of Square SS	df	Mean Sum of Square MSS
Due to regression (ESS)	$R^2 \cdot \sum y_i^2$	2	$(R^2 \cdot \sum y_i^2) / 2$
Due to residuals (RSS)	$(1 - R^2) \cdot \sum y_i^2$	$n - 3$	$[(1 - R^2) \cdot \sum y_i^2] / (n - 3)$
TSS	$\sum y_i^2$	$n - 1$	$\sum y_i^2 / (n - 1)$

Decision Rule: Testing the overall significance of a regression in terms of R^2

Given the k -variable regression model:

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + u_i$$

To test the hypothesis

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

(i.e., all slope coefficients are simultaneously zero) versus

$$H_1: \text{not all slope coefficients are simultaneously zero}$$

Compute

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

If $F > F_{\alpha}(k-1, n-k)$, we reject H_0 ; otherwise we cannot reject it, where $F_{\alpha}(k-1, n-k)$ is the critical F value at the α level of significance and $(k-1)$ numerator df and $(n-k)$ denominator df.