

# INTRODUCTION TO ECONOMETRICS

Econometrics can be defined as measurement in economics. It therefore deals with Modelling of the economy factors or elements e.g. Modelling demand and supply. Econometrics models can be classified into 2 depending on the nature of the data.

1. Regression models
2. Time series models

In this course we shall deal with regression models.

## Regression model

It is a mathematical equation that is used in forecasting such that the sum of the squared deviations of the estimated values from the actual values is minimised.

## Classification of Regression Models

- a) Simple regression models - This is a regression model that contains only one independent variable e.g.  $y = a + bx$ ,  $y = e^x$ .
- b) Multiple regression model - This is a regression model with more than one independent variable e.g.  $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$ ,  $y = e^{B_0 + B_1 x_1 + B_2 x_2}$ .
- c) Linear regression models - This is a regression model where independent variable(s) power is one i.e. the degree of the independent variable(s) is one.  
e.g.  $y = a + bx$ , - simple linear regression  
 $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_n x_n$  - multiple linear regression

d) Non-linear regression models - these includes all other regression models that are not linear i.e. logarithmic, exponential, quadratic etc.

In regression analysis we deal with;

- 1) Estimation of regression parameters.

The parameters are estimated such that the sum of the squared error term is minimised.

- 2) Evaluation of the regression model

This involves analysing how well the estimated regression model fits the data from which it was obtained.

### Simple linear regression model

This is a regression model of the form  $y = \alpha + \beta x$

$$y = \alpha + \beta x$$

$$\alpha = \frac{\sum y - b \sum x}{n}$$

$$b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

Estimation of  $\alpha$

$$\frac{\partial \sum e^2}{\partial \alpha} = 2 \sum (y - \alpha - \beta x)(-1)$$

$$\frac{\partial}{\partial \alpha} = -2 \sum (y - \alpha - \beta x)$$

$$\sum (y - \alpha - \beta x) = 0$$

$$\sum e^2 = \sum (y - \hat{y})^2$$

$$\sum e^2 = \sum (y - (\alpha + \beta x))^2$$

$$\sum e^2 = \sum (y - \alpha - \beta x)^2$$

$$\sum y - n\alpha - \beta \sum x = 0$$

$$n\alpha = \sum y - \beta \sum x$$

$$\alpha = \frac{\sum y - \beta \sum x}{n}$$

We can estimate the parameters  $\alpha$  and  $\beta$  such that we minimize the sum of the square error where  $\sum e^2 = \sum (y - \hat{y})^2$

We differentiate partially the sum of the squared term with respect to each of the parameters

Estimation of  $\beta$  is as shown above.

Estimation of  $\beta$ .

$$\sum e^2 = \sum (y - \alpha - \beta x)^2$$

$$\frac{\partial \sum e^2}{\partial \beta} = 2 \sum (y - \alpha - \beta x)(-x)$$

$$\frac{\partial}{\partial \beta} = -2 \sum (y - \alpha - \beta x)(-x)$$

$$\frac{\partial}{\partial \beta} = \sum (y - \alpha - \beta x)(-x)$$

$$\frac{\partial}{\partial \beta} = (\sum y - n\alpha - \beta \sum x)(-x)$$

$$\frac{\partial}{\partial \beta} = -\sum xy \sum (-xy + x\alpha + \beta x^2)$$

$$\frac{\partial}{\partial \beta} = -\sum xy + \alpha \sum x + \beta \sum x^2$$

$$\text{but } \alpha = \frac{\sum y - \beta \sum x}{n}$$

$$\frac{\partial}{\partial \beta} = -\sum xy + \left( \frac{\sum y - \beta \sum x}{n} \right) \sum x + \beta \sum x^2$$

$$\frac{\partial}{\partial \beta} = -n \sum xy + \sum x \sum y - \beta (\sum x)^2 + \beta n \sum x^2$$

$$\beta (\sum x)^2 - \beta n \sum x^2 = \sum x \sum y - n \sum xy$$

$$\beta (\sum x^2 -$$

$$\begin{array}{r} + \\ - \\ - \\ + \end{array}$$

$$0 = \sum xy - A \sum x - B \sum x^2$$

but

$$A = \frac{\sum y - B \sum x}{n}$$

$$0 = \sum xy - \left( \frac{\sum y - B \sum x}{n} \right) \sum x - B \sum x^2$$

$$n \sum xy - \sum x \sum y + B (\sum x)^2 - n B \sum x^2 = 0.$$

$$n B \sum x^2 - B (\sum x)^2 = n \sum xy - \sum x \sum y$$

$$B (n \sum x^2 - (\sum x)^2) = n \sum xy - \sum x \sum y$$

$$B = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

### Illustration.

The following data has been provided by the mgt accountant of ABC Ltd

Production Units	Pdctn	Cost.	
x	sh'000's.	xy	$x^2$
14	376	5264	196
18	398	7164	324
22	409	8988	484
27	458	12366	729
28	437	12288	784
32	459	14658	1024
36	471	16956	1296
39	482	18798	1521
42	491	20622	1764
45	519	23355	2025
303	4500	140447	10147

Required

- A regression equation of the pdctn cost on production units.
- Estimate the production cost for producing 65 units.

$$\begin{aligned} \text{Ans} &= \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2} \\ &= \frac{10 \times 140447 - 303 \times 4500}{10 \times 10147 - (303)^2} \\ &= -\frac{1259030}{10} \end{aligned}$$

$$\begin{aligned} B &= \frac{40970}{9661} & A &= 321.5 \\ &= 4.24 & \\ A &= \frac{\sum y - B \sum x}{n} & \\ &= \frac{4500 - 4.24(303)}{10} \end{aligned}$$

$$(i) \hat{y} = 321.5 + 4.24x$$

(ii) For 65 units

$$\begin{aligned}\hat{y} &= 321.5 + 4.24(65) \\ &= 321.5 + 275.6 \\ &= \underline{\underline{\text{ksh } 597,100}}\end{aligned}$$

### Multiple Linear Regression

This is a regression model of the form

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n.$$

In this course we shall start by analysing a multiple linear regression model with two independent variables of the form

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

#### Estimation of parameters

The parameters estimation should be such that the sum of the squared error term is minimized.

$$\begin{aligned}\sum e^2 &= \sum (\hat{Y} - Y)^2 \\ &= \sum (Y - (\beta_0 + \beta_1 X_1 + \beta_2 X_2))^2\end{aligned}$$

We need to differentiate the sum of the squared error term partially with respect to each of the independent variables. We shall have three eqns with three unknowns, the eqns are then solved simultaneously to obtain the estimate of each parameters.

$$\hat{Y} = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$\begin{aligned}\sum e^2 &= \sum (Y - \hat{Y})^2 \\ &= \sum (Y - (\beta_0 + \beta_1 X_1 + \beta_2 X_2))^2 \\ &= \sum (Y - \beta_0 - \beta_1 X_1 - \beta_2 X_2)^2\end{aligned}$$

$$\frac{\partial \sum e^2}{\partial \beta_0} = 2 \sum (Y - \beta_0 - \beta_1 X_1 - \beta_2 X_2)(-1)$$

$$\frac{\partial}{\partial \beta_1} = -2 \sum (Y - \beta_0 - \beta_1 X_1 - \beta_2 X_2)$$

$$\sum (Y - \beta_0 - \beta_1 X_1 - \beta_2 X_2) = 0$$

$$\sum Y - \beta_0 \sum 1 - \beta_1 \sum X_1 - \beta_2 \sum X_2 = 0$$

$$\sum Y = n \beta_0 + \beta_1 \sum X_1 + \beta_2 \sum X_2 \dots \dots \dots (1)$$

$$\frac{\partial \sum e^2}{\partial B_1} = 2 \sum (y - B_0 - B_1 x_1 - B_2 x_2) (-x_1)$$

$$\frac{\partial}{\partial B_2} = -2 \sum x_1 (y - B_0 - B_1 x_1 - B_2 x_2)$$

$$0 = \sum x_1 y - B_0 \sum x_1 - B_1 \sum x_1^2 - B_2 \sum x_1 x_2$$

$$\sum x_1 y = B_0 \sum x_1 + B_1 \sum x_1^2 + B_2 \sum x_1 x_2 \dots \text{(ii)}$$

$$\frac{\partial \sum e^2}{\partial B_2} = 2 \sum (y - B_0 - B_1 x_1 - B_2 x_2) (-x_2)$$

$$\frac{\partial}{\partial B_1} = -2 \sum x_2 (y - B_0 - B_1 x_1 - B_2 x_2)$$

$$0 = \sum x_2 y - B_0 \sum x_2 - B_1 \sum x_1 x_2 - B_2 \sum x_2^2$$

$$\sum x_2 y = B_0 \sum x_2 + B_1 \sum x_1 x_2 + B_2 \sum x_2^2 \dots \text{(iii)}$$

Writing the above 3 eqns in matrix form.

$$\begin{bmatrix} n & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 & \sum x_1 x_2 & \sum x_2^2 \end{bmatrix} \begin{bmatrix} B_0 \\ B_1 \\ B_2 \end{bmatrix} = \begin{bmatrix} \sum y \\ \sum x_1 y \\ \sum x_2 y \end{bmatrix}$$

The matrix can be solved using either ~~reverse~~

- i) inverse method
- ii) Grammer's rule

Cramers rule

$$B_0 = \frac{\begin{vmatrix} \sum y & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 & \sum x_1 x_2 & \sum x_2^2 \end{vmatrix}}{\begin{vmatrix} n & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 & \sum x_1 x_2 & \sum x_2^2 \end{vmatrix}}$$

$$B_2 = \frac{\begin{vmatrix} n & \sum y & \sum x_2 \\ \sum x_1 & \sum x_1 y & \sum x_1 x_2 \\ \sum x_2 & \sum x_2 y & \sum x_2^2 \end{vmatrix}}{\begin{vmatrix} n & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 & \sum x_1 x_2 & \sum x_2^2 \end{vmatrix}}$$

$$\beta_2 = \frac{\begin{vmatrix} n & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 & \sum x_2 x_1 & \sum x_2^2 \end{vmatrix}}{\begin{vmatrix} n & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_2 x_1 \\ \sum x_2 & \sum x_2 x_1 & \sum x_2^2 \end{vmatrix}}$$

Solve either using  
 \* Artistic method  
 \* Co factor method.

### Example 2.

$x_1$	$x_2$	$y$	$x_1 x_2$	$x_1^2$	$x_2^2$	$x_1 y$	$x_2 y$
2	3	18	6	4	9	36	54
4	8	12	32	16	64	48	96
6	12	15	72	36	144	90	180
8	4	24	32	64	16	192	96
10	1	32	10	100	1	320	32
12	5	25	60	144	25	300	125
14	6	29	84	196	36	406	174
16	7	41	112	256	49	656	287
18	9	48	162	324	81	864	432
20	11	55	220	400	121	1100	605
Rqd.		66	299	790	1540	546	2081

i) Estimated regression eqn of  $y$  on  $x_1$  and  $x_2$

a) Estimated value of  $y$  if  $x_1 = 30$  and  $x_2 = 17$

$$\Sigma y = n \beta_0 + \beta_1 \sum x_1 + \beta_2 \sum x_2$$

$$\sum x_1 y = \beta_0 \sum x_1 + \beta_1 \sum x_1^2 + \beta_2 \sum x_1 x_2$$

$$\sum x_2 y = \beta_0 \sum x_2 + \beta_1 \sum x_1 x_2 + \beta_2 \sum x_2^2$$

$$299 = 10 \beta_0 + 110 \beta_1 + 66 \beta_2$$

$$4012 = 110 \beta_0 + 1540 \beta_1 + 790 \beta_2$$

$$2081 = 66 \beta_0 + 790 \beta_1 + 546 \beta_2$$

$$299 \begin{bmatrix} 10 & 110 & 66 \\ 110 & 1540 & 790 \\ 66 & 790 & 546 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix} = \begin{bmatrix} 299 \\ 4012 \\ 2081 \end{bmatrix}$$

$$\beta_0 = \frac{\begin{vmatrix} 299 & 110 & 66 \\ 4012 & 1540 & 790 \\ 2081 & 790 & 546 \end{vmatrix}}{\begin{vmatrix} 10 & 110 & 66 \\ 110 & 1540 & 790 \\ 66 & 790 & 546 \end{vmatrix}}$$

Solving the numerator using artifice method.

$$\begin{vmatrix} 299 & 110 & 66 \\ 4012 & 1540 & 790 \\ 2081 & 790 & 546 \end{vmatrix} = \begin{matrix} \cancel{299} & \cancel{110} & \cancel{66} & \cancel{299} & \cancel{110} \\ \cancel{4012} & \cancel{1540} & \cancel{790} & \cancel{4012} & \cancel{1540} \\ \cancel{2081} & \cancel{790} & \cancel{546} & \cancel{2081} & \cancel{790} \end{matrix}$$

$$= 15144$$

$$(251,411,160 + 180,838,900 + 209,185,680) - (211,512,840 + 186,605,900 + 240,960,720)$$

$$= 641,435,740 - 639,079,460$$

$$= 2,356,280.$$

$$\begin{vmatrix} 10 & 110 & 66 \\ 110 & 1540 & 790 \\ 66 & 790 & 546 \end{vmatrix} = \begin{matrix} \cancel{10} & \cancel{110} & \cancel{66} & \cancel{10} & \cancel{110} \\ \cancel{110} & \cancel{1540} & \cancel{790} & \cancel{110} & \cancel{1540} \\ \cancel{66} & \cancel{790} & \cancel{546} & \cancel{66} & \cancel{790} \end{matrix}$$

$$= (8408400 + 5735400 + 5735400) - (6708240 + 6241000 + 6606600)$$

$$= (19879200 - 19555840)$$

$$= 323,360$$

$$\beta_0 = \frac{2356280}{323360}$$

$$= \underline{\underline{7.29}}.$$

$$\beta_1 = \frac{\begin{vmatrix} 10 & 299 & 66 \\ 110 & 4012 & 790 \\ 66 & 2081 & 546 \end{vmatrix}}{\begin{vmatrix} 10 & 110 & 66 \\ 110 & 1540 & 790 \\ 66 & 790 & 546 \end{vmatrix}} = \frac{10(546562) - 299(7920) + 66(-35882)}{-110(7920) + 1540(1104) - 790(640)}$$

Using the co factor method, above.

$$= \frac{729328}{323360}$$

$$= 2.26$$

$$\beta_2 = \frac{\begin{vmatrix} 10 & 110 & 299 \\ 110 & 1540 & 4012 \\ 66 & 790 & 2081 \end{vmatrix}}{\begin{vmatrix} 10 & 110 & 66 \\ 110 & 540 & 790 \\ 66 & 790 & 546 \end{vmatrix}} = \frac{10(35260) - 110(-35882) + 299(-14740)}{66(-51260) - 790(640) + 546(-6700)} = \frac{-107640}{323360} = -0.33.$$

$$\therefore \hat{y} = 7.29 + 2.26x_1 - 0.33x_2.$$

ii) When  $x_1 = 30$   $x_2 = 17$

$$\begin{aligned} y &= 7.29 + 2.26(30) - 0.33(17) \\ &= 7.29 + 67.8 - 5.61 \\ &= \underline{\underline{69.48}} \end{aligned}$$

### Evaluation of Regression Model

This involves analysing how well the estimated regression model fits the data from which it was obtained.

Evaluation of regression model can be classified into

- i) Evaluation of the whole regression model
- ii) Evaluation of b-coefficient / Slope.

### Evaluation of regression model as a whole

This involves determination of how well all the predictor variables in a regression model taken together predicts the dependent variable.

The methods used include:-

- i) Standard error of estimate ( $s_e$ )
- ii) Coefficient of determination ( $r^2$ )
- iii) F-statistic

### Standard error of estimate ( $s_e$ )

This is the average deviation of the estimated value from the actual value.

$$s_e = \sqrt{\frac{\sum (Y - \hat{Y})^2}{n - k}}$$

n - number of observations

k - number of parameters being estimated

$n - k$  - d.f of error term

- The smaller the standard error the better the regression
- The std. error of estimate is used in construction of confidence interval for the actual value of the predicted variable.

$$C.I = \hat{y} \pm Z_{\alpha/2} S_e \quad \text{for } n \geq 30.$$

$$\hat{y} - Z_{\alpha/2} S_e \leq y \leq \hat{y} + Z_{\alpha/2} S_e.$$

or

$$CI = \hat{y} \pm t_{\alpha/2} S_e$$

$$\hat{y} - t_{\alpha/2} S_e \leq y \leq \hat{y} + t_{\alpha/2} S_e \quad \text{for } n < 30.$$

### Example 3.

Using the data in example 1

- Compute the std. error of estimate
- Construct a 95% Confidence Interval (CI) for the production cost of 50 units

x	y	$\hat{y} = 321.5 + 4.24x$	$(y - \hat{y})^2$
14	376	380.86	23.6196
18	398	397.82	0.0324
22	409	414.78	33.4054
27	458	435.98	484.8804
28	437	440.22	10.2400
32	459	457.18	3.3124
36	471	474.14	9.8596
39	482	486.86	23.6196
42	491	499.58	73.6164
45	519	512.3	44.89
			<u>707.4788</u>

i)

$$\begin{aligned}
 S_e &= \sqrt{\frac{707.4788}{10-2}} \\
 &= \sqrt{\frac{707.4788}{8}} = \sqrt{88.43485} \\
 &= \underline{\underline{8.87}} \quad = \underline{\underline{9.40}} \quad V = df = n - k = 8
 \end{aligned}$$

$$\begin{aligned}
 ii) \quad \hat{y} &= 321.5 + 4.24(50) \\
 &= 533.5
 \end{aligned}$$

$$t_{0.975, 8} = 2.306$$

$$\begin{aligned}
 CI &= \hat{y} \pm t_{\alpha/2} S_e \\
 &= 533.5 \pm 2.306 \times 9.4 \\
 &= 533.5 \pm 21.6764 \\
 &= 511.82 \leq y \leq 555.19
 \end{aligned}$$

Example 4.

Using the data in example 2

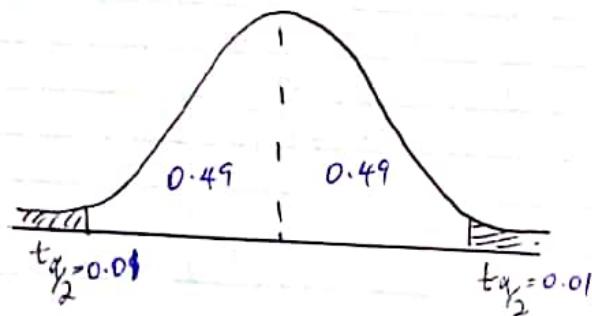
i) Compute the std error

ii) Construct a 98% CI for  $\bar{y}$  if  $x_1 = 20$   $x_2 = 25$

$x_1$	$x_2$	$y$	$\hat{Y} = 7.287 + 2.255x_1 - 0.333x_2$	$(\bar{Y} - \hat{Y})^2$
2	3	18	10.798	51.869
4	8	12	13.643	2.699
6	12	15	16.821	3.316
8	4	24	23.995	0.000
10	1	32	29.504	6.230
12	5	25	32.682	59.013
14	6	29	36.859	61.764
16	7	41	41.036	0.001
18	9	48	44.880	9.734
20	11	55	48.724	39.388
				234.014

$$\begin{aligned}
 \text{(i)} \quad S_e &= \sqrt{\frac{\sum (\bar{Y} - \hat{Y})^2}{n-k}} = \sqrt{\frac{234.014}{10-3}} = \sqrt{\frac{234.014}{7}} \\
 &= \underline{\underline{5.782}}
 \end{aligned}$$

$$\begin{aligned}
 \text{(ii)} \quad \hat{Y} &= 7.287 + 2.255x_1 - 0.333x_2 \quad \text{but } x_1 = 20, x_2 = 25 \\
 &= 7.287 + 2.255(20) - 0.333(25) \\
 &= 44.062
 \end{aligned}$$



$$\begin{aligned}
 V &= d.f = n-k \\
 &= 10-3 \\
 &= 7
 \end{aligned}$$

$$t_{0.99, 7} = 2.998$$

$$\begin{aligned}
 CI &= \bar{y} \pm t_{0.99, 7} S_e \\
 &= 44.062 \pm 2.998 \times 5.782 \\
 &= 44.062 \pm 17.334 \\
 &= 26.728 \leq \bar{y} \leq 61.396
 \end{aligned}$$

## Coefficient of determination

Coefficient of determination is defined as the proportion of variations that is explained by the regression model.  
It can be computed in two different approaches.

$$r^2 = \text{coefficient of determination}$$

$y$  - Actual  
 $\hat{y}$  - estimated  
 $\bar{y}$  - expected

$$\text{Sum of squares of Total Variation (SST)} \\ SST = \sum (y - \bar{y})^2$$

$$\text{Sum of squares of Error / Residual (SSE)} \\ SSE = \sum (\hat{y} - y)^2$$

$$\text{Sum of squares of regression model (SSR)} \\ SSR = \sum (\hat{y} - \bar{y})^2$$

$$SST = SSR + SSE$$

$$r^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

### Example 5

Using the data The variations are summarised in a table known as analysis of variance (ANOVA) table.

The general format of the ANOVA table is as follows

Source of Variation	sum of squares	degree of freedom	mean sum of squares	F-value
Regression	SSE	k-1	$MSR = \frac{SSR}{k-1}$	$\frac{SSE}{k-1} = \frac{MSR}{MSE}$
Error	SSE	n-k	$MSE = \frac{SSE}{n-k}$	
Total	SST	n-1		

### Example 5

From the data in example 1,

(i) Construct ANOVA Table

(ii) Compute the coefficient of determination and comment on your results.

x	y	$\bar{y}$	$SSE$	$SST$
			$(y - \bar{y})^2$	$(\bar{y} - \bar{\bar{y}})^2$
14	376	380.82	23.62	4780.34
18	398	397.82	0.03	2722.75
22	409	414.78	33.41	1240.45
27	458	435.98	484.88	196.56
28	437	440.22	10.24	95.65
32	459	457.18	3.31	51.55
36	471	474.14	9.82	582.74
39	482	486.86	23.62	1358.66
42	491	499.58	73.62	2458.18
45	519	512.3	<u>44.89</u>	<u>3881.29</u>
		$\bar{y} = 450$	<u>707.48</u>	<u>17368.17</u>
				<u>18,082</u>

$$\text{Prove that } \sum (y - \bar{y})^2 + \sum (\bar{y} - \bar{\bar{y}})^2 = \sum (y - \bar{\bar{y}})^2.$$

i)

Source of Variation	Sum of Squares	degree of freedom	Mean of squares	F - value
Regression	17368.16	1	17368.16	196.36
Error	707.60	8	88.45	
Total	18082*	9		

$$ii) r^2 = \frac{SSR}{SST} = \frac{17368.16}{18082} = 0.9605.$$

$$= 96.05\%.$$

96.05% of the variations in production cost is explained by variations in production units, the remaining proportion is explained by other factors which include the error term.

### Example 6

Using the data in example 2 Compute the coefficient of determination and comment on your results

$x_1$	$x_2$	$y$	$\bar{y}$	$SSE$	$SSR$	$SST$
2	3	18	10.80	51.87	364.81	141.61
4	8	12	13.64	2.70	264.39	320.41
6	12	15	16.82	3.32	171.09	222.01
8	4	24	24.	0.00	34.81	34.81
10	1	32	29.50	6.23	0.16	4.41
12	5	25	32.68	59.01	7.73	24.01
14	6	29	36.86	61.76	48.44	0.81
16	7	41	41.04	0.00	124.10	123.21
18	9	48	44.88	9.73	224.40	327.61
20	11	55	48.72	39.39	354.19	630.01
			$\bar{y} = 29.9$	<u>234.01</u>	<u>1594.12</u>	<u>1828.9</u>

$$r^2 = \frac{SSR}{SST} = \frac{1594.12}{1828.9} = 0.8720 = 87.20\%$$

~~87.22~~ ~~87.20%~~ of the variations in  $y$  is explained by variations in  $x_1$  and  $x_2$ , the remaining ~~87.20%~~ is explained by other factors which include the error term.

### F - statistics

This statistic is used to carry out hypothesis testing on the significance of the regression model.

The hypothesis is tested as follows: —

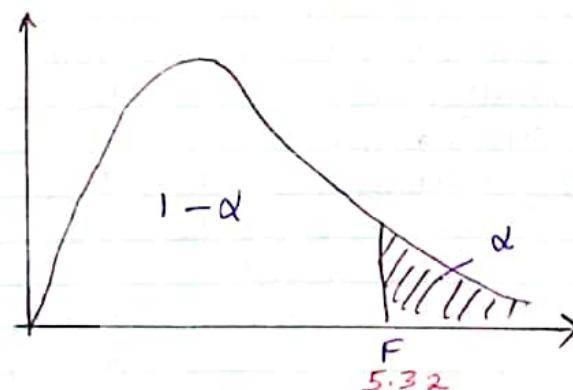
$H_0 : r^2 = 0$  Regression model is not statistically significant.  
 $H_a : r^2 \neq 0$ . Regression model is statistically significant

F - statistic

$$F = \frac{SSR/k-1}{SSE/(n-k)}$$

F - Distribution

$$F_{k-1, n-k}$$



### Example 7.

Test on the significance of the regression model obtained in example 1 at 5% significance level.

$$H_0 : r^2 = 0$$

$$F_{k-1, n-k}$$

$$H_a : r^2 \neq 0.$$

$$F_{1, 8} = 5.32$$

$$F = \frac{SSR/k-1}{SSE/(n-k)} = 196.36$$

\* At 5% significance level we reject the null hypothesis and conclude that the regression model is statistically significant.

### Example 8

Using the data in example 2 test whether the estimated regression model is statistically significant at 1% significance level.

$$H_0: r^2 = 0$$

$$H_1: r^2 \neq 0$$

F-statistic

$$F = \frac{\frac{SSR}{k-1}}{\frac{SS\hat{e}}{n-k}}$$

$$= \frac{1594.12}{\frac{3-1}{234.01}}$$

$$= \frac{1594.12}{103}$$

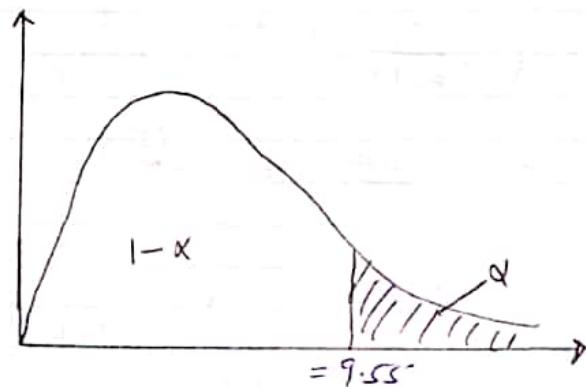
$$= \frac{797.06}{33.43}$$

$$= 23.84$$

F-distribution

$$F_{k-1, n-k} = F_{2, 7}$$

$$= 9.55$$



At 1% significance level we reject the null hypothesis and conclude that the regression model is statistically significant.

$r^2$  = Power of the model.

Evaluation of the regression Model using the slope / B-coefficient.

In this case we evaluate how well a given independent variable predicts the dependent variable holding other factors constant.

The following measures are used:-

- i) Correlation coefficient ( $r, f$ )
- ii) The Std error of the slope ( $\sigma_b$ )
- iii) t or Z statistic

Correlation coefficient ( $r, f$ )

Correlation is the measure of the degree and nature of relationship shown between two variables.

Correlation is measured using correlation coefficient which ranges from -1 and 1.

A negative correlation implies that if one of the variables is increasing the other variable is decreasing and vice versa.

A positive correlation implies that the two variables vary in the same direction i.e., if one of the variables is increasing the other variable is also increasing and vice versa.

The nature of correlation can be obtained by checking the nature of the coefficient of the independent variable in the regression model.

e.g

$$y = 4.2 - 3.6x_1 + 8.2x_2$$

$x_1$	$x_2$	$y$
2	3	21.6 > 3.6
3	3	18.

Interpretation of coefficients

-3.6 Holding other factors constant if  $x_1$  was increased by 1 unit the estimated value of  $y$  would decrease by 3.6 units

8.2 This implies if  $x_2$  was increased by 1 unit holding other factors constant the estimated value of  $y$  will increase by 8.2 units.

Std. error of the slope ( $s_b$ )

This is the average error term for the  $b$  coefficient.

$$s_b = \frac{s_e}{\sqrt{\sum (x - \bar{x})^2}}$$

The smaller the standard error the better the predictor variable in predicting the dependent variable.

The std error of the slope is used in construction Confidence Interval for the  $b$ -coefficient.

$$CI = \hat{b} \pm Z_{\alpha/2} s_b \quad - \text{for large sample size } n \geq 30$$

$$CI = \hat{b} \pm t_{\alpha/2} s_b \quad - \text{for small sample size}$$

$$\hat{b} - Z_{\alpha/2} s_b \leq b \leq \hat{b} + Z_{\alpha/2} s_b.$$

Example 9.

Construct a 95% Confidence interval for the  $b$ -coefficient using the data in example 1.

$x$	$y$	$(x - \bar{x})^2$
14	376	265.69
18	398	151.29
22	409	68.89
27	458	10.89
28	437	5.29
32	459	2.89
36	471	32.49
39	482	75.69
42	491	136.89
45	519	216.09
303		966.1
		$\bar{x} = 30.3$

$$s_b = \frac{9.40}{\sqrt{966.1}} = \frac{9.40}{31.08} = 0.3024$$

$$CI = 4.24 \pm \frac{2.306}{1.96 \times 0.3024}$$

$$= 4.24 \pm 0.39$$

$$= 3.85 \leq b \leq 4.83$$

$$3.54 \leq b \leq 4.94$$

11

Example 10.  
Construct 99% C.I for the b coefficient of the estimated regression model in example 2.

$x_1$	$x_2$	$y$	$(x_1 - \bar{x}_1)^2$	$(x_2 - \bar{x}_2)^2$
2	3	18	102.01	81
4	8	12	65.61	49
6	12	15	37.21	25
8	4	24	16.81	9
10	1	32	4.41	1
12	5	25	0.01	1
14	6	29	3.61	9
16	7	41	15.21	25
18	9	48	34.81	49
20			62.41	81
$\bar{x}_1 = 11.0$		<u>55</u>	<u>342.10</u>	<u>330.</u>
$\bar{x}_2 = 6.6$				<u>110.40</u>

$$S_e = 5.782$$

$$S_{b_1} = \frac{S_e}{\sqrt{\sum(x_1 - \bar{x}_1)^2}}$$

$$\begin{aligned} S_{b_1} &= \frac{S_e}{\sqrt{\sum(x_1 - \bar{x}_1)^2}} \\ &= \frac{5.782}{\sqrt{342.1}} \quad \frac{5.782}{\sqrt{330}} \\ &= 0.313 \quad = 0.318 \end{aligned}$$

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

$$= 2.26 \pm 2.576 \times 0.313$$

$$= 2.26 \pm 0.81$$

$$1.45 \leq b_1 \leq 3.07$$

$$= 2.26 \pm 3.499 \times 0.318$$

$$= 2.26 \pm 1.098 \times 1.11$$

$$= 1.165 \leq b_1 \leq 3.355$$

$$1.15 \leq b_1 \leq 3.37$$

$$S_{b_2} = \frac{S_e}{\sqrt{\sum(x_2 - \bar{x}_2)^2}}$$

$$= \frac{5.782}{\sqrt{110.40}}$$

$$= 0.550$$

$$CI = b_2 \pm t_{\alpha/2} S_{b_2}$$

$$= -0.33 \pm 2.576 \times 0.55$$

$$= -0.33 \pm 1.42$$

$$= -1.75 \leq b_2 \leq 1.09$$

$$= -0.33 \pm 3.499 \times 0.55$$

$$= -0.33 \pm 1.92$$

$$= -2.25 \leq b_2 \leq 1.59$$

### 3. t or Z Statistics

This statistic is used to test whether a given predictor variable is statistically significant in predicting the dependent variable.

$$t \text{ or } Z = \frac{b}{S_b}$$

The hypothesis is stated as follows:-

$H_0: b = 0$  (predictor variable is not statistically significant)

$H_1: b \neq 0$  (predictor variable is statistically significant)

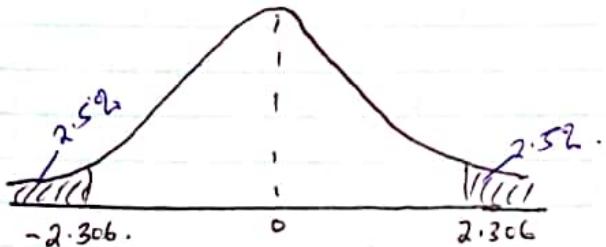
#### Example 11

At 5% significance level test whether prdn units is a statistically significant predictor variable for production cost using the data in example 1.

$H_0: \beta = 0$  (prdn units is not a statistically significant predictor variable for prdn cost)

$H_1: \beta \neq 0$  (prdn units is a statistically significant predictor variable for prdn cost).

$$t = \frac{b}{S_b} = \frac{4.24}{0.3024} = 14.02$$



$$t_{\alpha/2} = t_{0.975, 8} = 2.306$$

Conclusion: At 5% significance level we reject the null hypothesis and conclude that prdn units is a statistically significant predictor variable for prdn cost.

#### Example 12

At 1% significance level test whether  $X_1$  and  $X_2$  are significant predictor variables for  $y$  using the data in example 2.

$H_0: \beta_1 = 0$  ( $X_1$  is not statistically significant predictor of  $y$ )

$H_1: \beta_1 \neq 0$  ( $X_1$  is statistically significant predictor of  $y$ )

$$t = \frac{b_1}{S_{b_1}} = \frac{2.26}{0.318} = 7.11$$

$$t_{\alpha/2} = t_{0.995, 7} = \underline{\underline{3.499}}$$

At 1% significance level we reject the null hypothesis and conclude that  $x_2$  is statistically significant predictor of  $y$ .

$H_0: \beta_2 = 0$  ( $x_2$  is not a statistically significant predictor of  $y$ )  
 $H_1: \beta_2 \neq 0$  ( $x_2$  is statistically significant predictor of  $y$ )

$$t_2 = \frac{\beta_2}{\text{S}_{\beta_2}} = \frac{-0.33}{0.55} = -0.6$$

$$t_{df_2} = 2.995 \quad 3.499.$$

At 1% significance level we fail to reject the null hypothesis and conclude that  $x_2$  is not statistically significant predictor of  $y$ .

### Example Question

A study was conducted to examine the predictive ability of economists and finance specialists on their future performance of their country's economy. Equal random samples of economist and finance specialists were asked to make predictions of the percentage change in the gross domestic product (GDP) of their countries. A year later the actual percentage change in GDP ( $y$ ) was regressed against the economists prediction ( $x_1$ ) and the finance specialists prediction  $x_2$  using a statistical software. Below is part of the output:

Parameters estimates

Variable	Co-efficient	Standard Error	t-ratio
Constant	0.983	4	0.34
$x_1$	B	0.257	2.42
$x_2$	0.587	0.225	2.61

### Analysis of Variance (ANOVA)

Source	Degrees of Freedom	Sum of Squares	Mean of Sum of Squares	F-value
Regression	2	2904.8	D	E
Error	11	100.1		
Total	13	C	9.1	

Required

- Compute the values A, B, C, D, E
- The estimated regression model

5mks  
2mks

- iii) The explanatory power of the model  
 iv) Test the significance of each of the predictor variables at 5% significance level (6mks)  
 v) Comment on the statistical significance of the model at 5% significance level.

$$\begin{array}{l}
 \text{A} \quad b = t \times s_b \\
 t = \frac{b}{s_b} = \frac{2.42}{0.257} = 9.42 \\
 s_b = \frac{b}{t} = \frac{0.622}{9.42} = 0.066 \\
 \text{B} \quad b = t \times s_b = 2.42 \times 0.257 = 0.622 \\
 \text{C} \quad C = SST = SSR + SSE = 2904.8 + 100.1 = 3004.9 \\
 \text{D} \quad MSR = \frac{SSR}{k-1} = \frac{2904.8}{2} = 1452.4
 \end{array}$$

$$E = F = \frac{MSR}{MSE} = \frac{1452.4}{9.1} = 159.60.$$

$$\hat{Y} = 0.983 + 0.622X_1 + 0.587X_2.$$

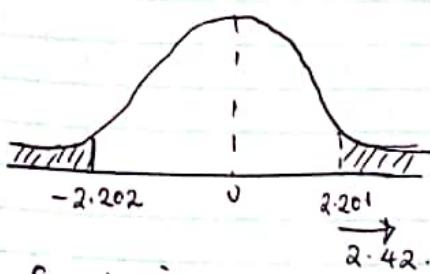
$$r^2 = \frac{SSR}{SST} \times 100 = \frac{2904.8}{3004.9} \times 100 = 96.67\%$$

96.67% of the change in GDP is explained by the economists and finance specialists prediction while the remaining proportion is explained by other factors including the error term.

v) For  $X_1$

$$t_{\text{table}} = 2.201 \quad H_0: X_1 = 0$$

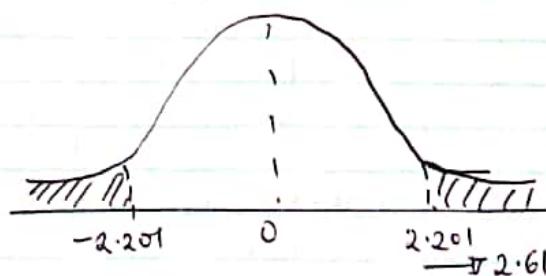
$$t_C = 2.42. \quad H_P: X_1 \neq 0$$



for  $X_2$

$$t_{\text{table}} = 2.201 \quad H_0: X_2 = 0$$

$$t_C = 2.61 \quad H_P: X_2 \neq 0$$



Conclusion

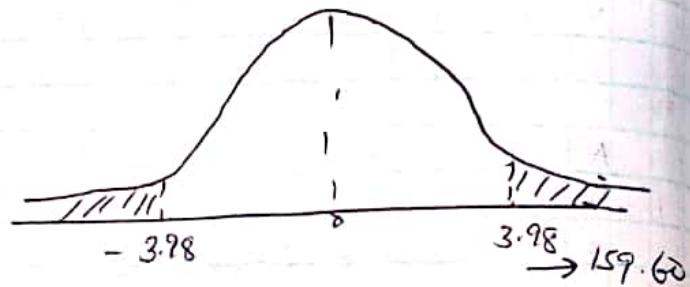
At 5% significance level, we reject the null hypothesis and conclude that the economists prediction is a statistically significant predictor of GDP variable of GDP.

At 5% Significance level, we reject the null hypothesis and conclude that the finance specialists prediction is a statistically significant predictor variable of GDP.

$H_0: r=0$  (regression model is not statistically significant)  
 $H_1: r \neq 0$  (regression model is statistically significant)

$$f_{k-1, n-k} = f_{2, 9} = 3.98$$

$$F_C = 159.60$$



At 5% significance level, we reject the null hypothesis and conclude that the regression model is statistically significant.

### Question 2.

KCA University research department plans to do an investigation on the impact of student's final score and their linguistic aptitude at entry point. For the purposes of this survey, the following sample data has been assembled for 12 subject areas under investigation. The University applies a 5% significance level in this kind of statistical analysis.

Subject	Linguistic Aptitude	Final Score
1	5	5
2	10	20
3	6	4
4	8	15
5	4	11
6	4	9
7	3	12
8	10	18
9	2	7
10	6	2
11	7	14
12	9	17

Required:

- Obtain the estimated regression model
- Comment on the meaning of the coefficient values
- Comment on the significance of linguistic aptitude in predicting the final score
- Is the regression model statistically significant.

$$\hat{y} = a + bx$$

$$a = \frac{\sum y - b \sum x}{n}$$

$$b = \frac{n \sum xy - \sum x \sum y}{n \sum x^2 - (\sum x)^2}$$

$x$	$y$	$(x-\bar{x})^2$	$xy$	$x^2$	$\hat{y}$	$(y-\hat{y})^2$	$(\bar{y}-\hat{y})^2$	$(y-\bar{y})^2$
5	5	5.76	25	25	9.46	19.89	2.72	38.07
10	20	6.76	200	100	16.71	48.82	30.69	77.97
6	4	1.96	24	36	10.91	44.75	0.07	51.41
8	15	0.36	120	64	13.81	33.76	6.97	14.67
4	11	11.56	44	16	8.01	16.84	9.99	0.03
4	9	11.56	36	16	8.01	16.84	9.99	4.71
3	12	19.36	36	9	6.56	29.59	20.34	0.69
10	18	6.76	180	100	16.71	1.66	30.69	46.65
2	7	29.16	14	4	5.11	3.57	36.72	17.37
6	2	1.96	12	36	10.91	79.39	0.07	84.09
7	14	0.16	98	49	12.36	2.69	1.42	8.01
9	17	2.56	153	81	15.26	3.03	16.73	33.99
74	134	86.92540	942	536		209.77	166.6	377.68

$$b = \frac{12 \times 942 - 74 \times 134}{12 \times 536 - 74^2}$$

$$= 1.45$$

$$a = \frac{134 - 1.45(74)}{12}$$

The minimum final score is 2.21

ii) Coefficient of  $x = 1.45 \rightarrow$  When the value of linguistic aptitude is increased by one the final score is increased by 1.45.

iii)  $H_0: b = 0$  (linguistic aptitude is not a statistically significant predictor variable of the final score)  
 $H_1: b \neq 0$  (Linguistic aptitude is statistically significant predictor variable of the final score)

$$S_e = \sqrt{\frac{\sum (y - \hat{y})^2}{n-k}}$$

$$= \sqrt{\frac{209.77}{12-2}}$$

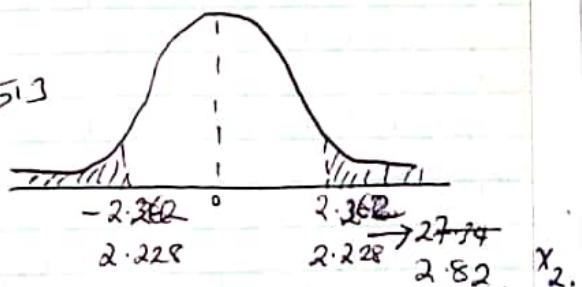
$$= 4.58$$

$$t = \frac{b}{S_b}$$

$$= \frac{1.45}{0.513}$$

$$= \frac{27.34}{2.82}$$

$$t_{0.975} = 2.262$$



$$S_b = \frac{S_e}{\sqrt{\sum (x-\bar{x})^2}}$$

$$= \frac{4.58}{\sqrt{86.368.93}}$$

$$= 0.053$$

$$0.513$$

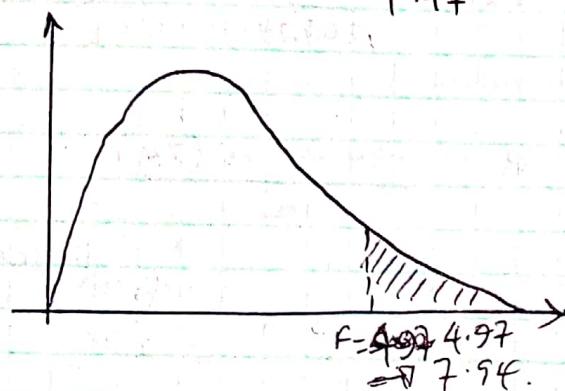
At 5% significance level, we reject the null hypothesis and conclude that linguistic aptitude is statistically significant predictor variable of the final score.

- iv)  $H_0: r^2 = 0$  (regression model is not statistically significant)  
 $H_1: r^2 \neq 0$  (regression model is statistically significant)

$$r^2 = \frac{SSR}{SST} = \frac{166.6}{377.68} \times 100 = 44.11\%$$

$$F \text{ value} = \frac{\frac{SSR}{k-1}}{\frac{SSE}{n-k}} = \frac{\frac{166.6}{1}}{\frac{209.77}{12-2}} = 7.94$$

$$F_{k-1, n-k} = F_{1, 10} = 4.97$$



At 5% significance level,  
we reject the null hypothesis  
and conclude that the  
regression model is  
statistically significant.

### Example 13.

The job ~~satisfaction~~ performance of an employee is said to be influenced by

(i) remuneration / benefits ( $x_1$ )

(ii) experience ( $x_2$ )

(iii) gender ( $x_3$ )

(iv) academic qualifications ( $x_4$ )

The following data has been obtained from the HR records relating to 20 employees.

The performance is believed to assume a linear regression model.

	$y$	$x_1$	$x_2$	$x_3$	$x_4$	$\Sigma$
$y$ -1-10	8.4	7	5	M 0	4	
$x_1$ -3-10	6.2	6	3	M 0	4	
	3.9	5	1	F 1	3	
$x_2$ -1-6	4.5	3	2	M 0	2	
$x_3$	F, M	6.2				
	5.2	4	4	F 1	1	
$x_4$ -1-5	3.5	7	6	F 1	2	
	8.0	9	2	M 0	5	
	4.7	6	4	F 1	4	
	3.3	4	2	F 1	5	
	7.7	3	6	M 0	4	
	8.1	7	6	F 1	3	
	7.3	6	5	F 1	3	

$y$	$x_1$	$x_2$	$x_3$	$x_4$
6.0	8	4	M 0	2
7.5	7	5	F 1	3
8.8	6	6	M 0	4
3.4	5	1	F 1	5
7.2	5	6	M 0	3
7.0	6	5	M 0	4
8.6	4	6	M 0	3
9.9	8	6	F 1	5

$$\hat{y} = 1 + 0.25x_1 + 0.76x_2 - 1.384x_3 + 0.459x_4$$

Required :

- (i) Interpret the coefficients for each of the variables including the constant.

Intercept - 1. The minimum job performance level is 1.

$x_1 = 0.248$ . If the remuneration scale increases by 1 the job performance of an employee will increase by 0.248

$x_2 = 0.76$  - If the experience scale increases by 1 the employee performance would increase by 0.76

$x_3$  - Holding other factors constant the performance of a male employee is higher than the performance of a female employee by 1.384

$x_4$  - If the academic qualification scale of the employee is increased by 1 the job performance would increase by 0.459.

Test the significance of each of the predictor variables at 5% significance level

Construct a 90% CI for  $x_1$  and  $x_2$  and comment on your results

Test the power of model

05/03/2020

### Dummy Variable

A dummy variable is used in regression analysis to indicate absence or presence of some categorical outcome.

A dummy variable will take value 0 or 1. A dummy independent variable which for some observations has a value of zero will make the variables coefficient to have no role or influence in the dependent variable whereas when the dummy variable takes a value of 1 its coefficients will have an influence on the dependent variable.

Example of dummy variables

Gender can either be male or female  
Political affiliation e.g. democratic or republican  
Academic qualification graduate or non-graduate

Example

The following data relates to the final score of the students as influenced by gender and class attendance

Final Score out of 100	Gender	Class attendance 0-13
48	F	10
64	M	11
72	F	11
55	F	8
69	M	9
57	F	7
46	M	12
66	M	8
50	F	9
61	F	10

Required

- Regression model for the above data
- Interpret the meaning of the coefficients for each of the independent variables
- Test whether each of the predictor variable is statistically significant at 1% significance level.
- Comment on the significance of the regression model at 5% significance level.

$y$	$x_1$	$x_2$	$x_1x_2$	$x_1^2$	$x_2^2$	$x_1y$	$x_2y$
48	1	10	10	1	100	48	480
64	0	11	0	0	121	0	704
72	1	11	11	1	121	72	792
55	1	8	8	1	64	55	440
69	0	9	0	0	81	0	621
57	1	7	7	1	49	57	399
46	0	12	0	0	144	46	552
66	0	8	0	0	64	0	528
50	1	9	9	1	81	50	450
<u>61</u>	<u>1</u>	<u>10</u>	<u>10</u>	<u>1</u>	<u>100</u>	<u>61</u>	<u>610</u>
<u>588</u>	<u>6</u>	<u>95</u>	<u>55</u>	<u>6</u>	<u>925</u>	<u>343</u>	<u>5576</u>

$$\sum y = n \beta_0 + \beta_1 \sum x_1 + \beta_2 \sum x_2$$

$$\sum x_1 y = \beta_0 \sum x_1 + \beta_1 \sum x_1^2 + \beta_2 \sum x_1 x_2$$

$$\sum x_2 y = \beta_0 \sum x_2 + \beta_1 \sum x_1 x_2 + \beta_2 \sum x_2^2$$

$$\begin{bmatrix} \sum y \\ \sum x_1 y \\ \sum x_2 y \end{bmatrix} = \begin{bmatrix} n & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 & \sum x_1 x_2 & \sum x_2^2 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{bmatrix}$$

$$\beta_0 = \frac{\begin{vmatrix} \sum y & \sum x_1 & \sum x_2 \\ \sum x_1 y & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 y & \sum x_2 x_1 & \sum x_2^2 \end{vmatrix}}{\begin{vmatrix} n & \sum x_1 & \sum x_2 \\ \sum x_1 & \sum x_1^2 & \sum x_1 x_2 \\ \sum x_2 & \sum x_2 x_1 & \sum x_2^2 \end{vmatrix}}$$

$$= \frac{\begin{vmatrix} 588 & 6 & 95 \\ 343 & 6 & 55 \\ 5576 & 55 & 925 \end{vmatrix}}{\begin{vmatrix} 10 & 6 & 95 \\ 6 & 6 & 55 \\ 95 & 55 & 925 \end{vmatrix}}$$

$$= \frac{(10)(5576 - 55 \cdot 925) - 6(343 \cdot 925 - 55 \cdot 55) + 95(6 \cdot 55 - 6 \cdot 6)}{(10)(5576 - 55 \cdot 925) - 6(343 \cdot 925 - 55 \cdot 55) + 95(6 \cdot 55 - 6 \cdot 6)}$$

$$= \frac{588(5576 - 5025) - 6(317275 - 306680) + 95(18865 - 33456)}{10(5576 - 5025) - 6(317275 - 306680) + 95(18865 - 33456)}$$

$$= \frac{34985}{500}$$

$$= 69.97$$

$$\beta_1 = \frac{\begin{vmatrix} 10 & 588 & 95 \\ 6 & 343 & 55 \\ 95 & 5576 & 925 \end{vmatrix}}{\begin{vmatrix} 10 & 6 & 95 \\ 6 & 6 & 55 \\ 95 & 55 & 925 \end{vmatrix}}$$

$$= \frac{10(317275 - 306680) - 588(5576 - 5025) + 95(33456 - 32585)}{500}$$

$$= \frac{-2405}{500}$$

$$= -4.81$$

$$\beta_2 = \frac{\begin{vmatrix} 10 & 6 & 588 \\ 6 & 6 & 343 \\ 95 & 55 & 5576 \end{vmatrix}}{\begin{vmatrix} 10 & 6 & 95 \\ 6 & 6 & 55 \\ 95 & 55 & 925 \end{vmatrix}}$$

$$\beta_2 = \frac{10(33456 - 18465) + 6(33456 - 32588) + 588(330 - 570)}{500}$$

$$= \frac{-436}{500}$$

$$= -0.872.$$

$$y = 69.97 - 4.81x_1 - 0.872x_2,$$

(i)  $F=1, M=0$   
Score male student 5 classes

$$y = 69.97 - 4.81(0) + 0.872(5)$$

$$= 65.61$$

$F=0, M=1$

$$y = 65.16 + 4.81(1) - 0.872(5)$$

$$= 65.61.$$

ii)  $\beta_0 = 69.97$

$\beta_1 = -4.81 \rightarrow$  Holding other factors constant, the male student final score will be higher by 4.81 than a female student who attended the same number of classes

$\beta_2 = -0.872$  - Holding other factors constant when the class attended is increased by one the final score would decrease by 0.872.

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

$$t = \frac{\beta_1}{S_{\beta_1}}$$

$$S_{\beta_1} = \sqrt{\frac{2(y-\bar{y})^2}{n-k}} / \sqrt{\frac{S(x-\bar{x})^2}{n-k}}$$

$$S_{\beta_1} = \sqrt{\frac{681.77}{7}} = \frac{6.37}{4.05581}$$

$$S_{\beta_2} = \sqrt{\frac{681.77}{22.5}} = 2.08$$