

Unit 4 :- Multiple Correlation and Regression

→ Partial correlation coefficient

It is the statistical technique of studying relationship between a dependent variable and an independent variable after controlling the effect of other independent variables.

Types:

- i) Zero Order Partial Correlation Coefficient.
- ii) First Order Partial Correlation Coefficient
- (iii) Higher Order Partial Correlation Coefficient

(1) Zero order Partial correlation coefficient

It studies the relationship b/w two variables but without controlling other variables on both correlating variables.

(2) First order partial correlation coefficient

→ It studies the relationship between two variables after removing the association of third variable with both of those correlating variables.

(3) Higher order partial correlation coefficient

→ It studies the degree of association bet'n two variables after controlling for the effect of two or more variables.

Imp

(2) First order partial correlation coeff.

Consider the study of 3 variables; y , x_1 & x_2 . where y being a dependent variable & x_1 & x_2 being independent variables. We may be interested to know correlation b/w y & x_1 , y & x_2 & x_1 & x_2 . But correlation b/w any two variables may be partly due to the correlation b/w 3rd variables with both correlating variables.

In such situation, we control the effect of 3rd variable on both correlating variable.

y = B.P of person

x_1 = Age of person

x_2 = Weight "

$r_{y \cdot x_1}$ = Partial correlation coefficient b/w y and x_1 after controlling for the effect of x_2

$$= \frac{r_{y1} - r_{y2} \cdot r_{12}}{\sqrt{1 - r_{y2}^2} \sqrt{1 - r_{12}^2}}$$

$r_{y \cdot x_1}$ = y and x_2 $\dots x_1$

$$= \frac{r_{y2} - r_{y1} \cdot r_{12}}{\sqrt{1 - r_{y1}^2} \sqrt{1 - r_{12}^2}}$$

$r_{12 \cdot y}$ = x_1 and x_2 after y

$$= \frac{r_{12} - r_{y1} \cdot r_{y2}}{\sqrt{1 - r_{y1}^2} \sqrt{1 - r_{y2}^2}}$$

$$\begin{array}{l|l} r_y = r_y & -1 \leq r_{y \cdot x_1} \leq +1 \\ r_{y2} = r_{2y} & -1 \leq r_{y2 \cdot 1} \leq +1 \\ r_{12} = r_{21} & -1 \leq r_{12 \cdot y} \leq +1 \end{array}$$

$$r_{y1} = \frac{n \sum yx_1 - \sum y \cdot \sum x_1}{\sqrt{n \sum y^2 - (\sum y)^2} \sqrt{n \sum x_1^2 - (\sum x_1)^2}}$$

$$r_{y2} = \frac{n \sum yx_2 - \sum y \cdot \sum x_2}{\sqrt{n \sum y^2 - (\sum y)^2} \sqrt{n \sum x_2^2 - (\sum x_2)^2}}$$

$$r_{12} = \frac{n \sum x_1 x_2 - \sum x_1 \cdot \sum x_2}{\sqrt{n \sum x_1^2 - (\sum x_1)^2} \sqrt{n \sum x_2^2 - (\sum x_2)^2}}$$

Numerical

The salaries of IT workers are expected to be dependent on years of schooling and yrs. of work experience. The following table gives information on the annual salary (in thousands. of dollar) for 12 persons.

S.N	Ann. Salary (in thds. of dollars)	Years of schooling	Years of experience
1.	52	16	6
2.	44	12	10
3.	48	13	15
4.	77	20	8
5.	68	18	11
6.	48	16	2
7.	59	14	12
8.	83	18	4
9.	28	12	6
10.	61	16	9
11.	27	12	2

12.

69

16

18

- (a) Compute simple correlation coefficient.
 (b) Compute partial correlation coefficient & interpret the results.

Soln:-

Here, y = Annual Salary x_1 = yrs. of schooling x_2 = yrs. of experience

Computational Result

$\sum y = 634$	$\sum y^2 = 36086$	$\sum yx_1 = 10036$
$\sum x_1 = 183$	$\sum x_1^2 = 2869$	$\sum yx_2 = 5865$
$\sum x_2 = 103$	$\sum x_2^2 = 1155$	$\sum x_1 x_2 = 1569$
$n = 12$		

Computational sample correlation coefficient

$$\begin{aligned}
 r_y &= \frac{n \sum yx_1 - \sum y \sum x_1}{\sqrt{n \sum x_1^2 - (\sum x_1)^2} \sqrt{n \sum x_2^2 - (\sum x_2)^2}} \\
 &= \frac{12 \times 10036 - 634 \times 183}{\sqrt{12 \times 2869 - (183)^2} \sqrt{12 \times 1155 - (103)^2}} \\
 &\approx 0.8164
 \end{aligned}$$

$$r_{Y_2} = \frac{n \sum Y X_2 - \sum Y \sum X_2}{\sqrt{n \sum Y^2 - (\sum Y)^2} \sqrt{n \sum X_2^2 - (\sum X_2)^2}}$$

$$= 0.505$$

$$r_{12} = \frac{n \sum X_1 X_2 - \sum X_1 \sum X_2}{\sqrt{n \sum X_1^2 - (\sum X_1)^2} \sqrt{n \sum X_2^2 - (\sum X_2)^2}}$$

$$= -0.0120$$

$$r_{Y_{1.2}} = \frac{r_{Y_1} - r_{Y_2} \cdot r_{12}}{\sqrt{1 - r_{Y_2}^2} \sqrt{1 - r_{12}^2}}$$

$$= \frac{0.816 - 0.505 \times (-0.0120)}{\sqrt{1 - (0.505)^2} \sqrt{1 - (-0.0120)^2}} = 0.952$$

$$r_{Y_{2.1}} = \frac{r_{Y_2} - r_{Y_1} \cdot r_{12}}{\sqrt{1 - r_{Y_1}^2} \sqrt{1 - r_{12}^2}}$$

$$= \frac{0.505 - 0.8164 \times (-0.0120)}{\sqrt{1 - (0.8164)^2} \sqrt{1 - (0.0120)^2}} = 0.891$$

$$r_{12 \cdot Y} = \frac{r_{12} - r_{Y_1} \cdot r_{Y_2}}{\sqrt{1 - r_{Y_1}^2} \sqrt{1 - r_{Y_2}^2}}$$

$$= \frac{-0.0120 - 0.8164 \times 0.505}{\sqrt{1 - (0.8164)^2} \sqrt{1 - (0.505)^2}} = -0.849$$

Simple Correlation

$$r_{y_1} = 0.816$$

$$r_{y_2} = 0.505$$

$$r_{12} = -0.0120$$

Partial correlation

$$r_{y_1 \cdot 2} = 0.952$$

$$r_{y_2 \cdot 1} = 0.891$$

$$r_{12 \cdot y} = -0.849$$

Interpretation :-

1. There is a very strong positive correlation betn annual salary and years of education after controlling for years of experience i.e. among the people with same yrs. of experience.
2. There is a very strong positive correlation betn annual salary and yrs. of experience after controlling for the effect of yrs. of education i.e. among the people with same yrs. of education.
3. There is a strong negative correlation betn yrs. of education and yrs. of schooling after controlling for the effect of Annual salary. It means that people with less education has to work for longer period in order to get the same salary as the people with higher degree or more yrs. of education.

Multiple linear regression

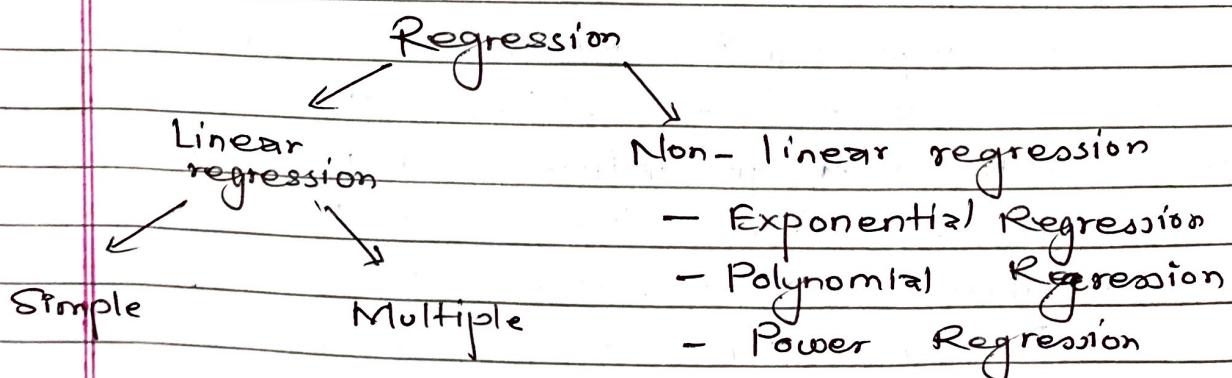
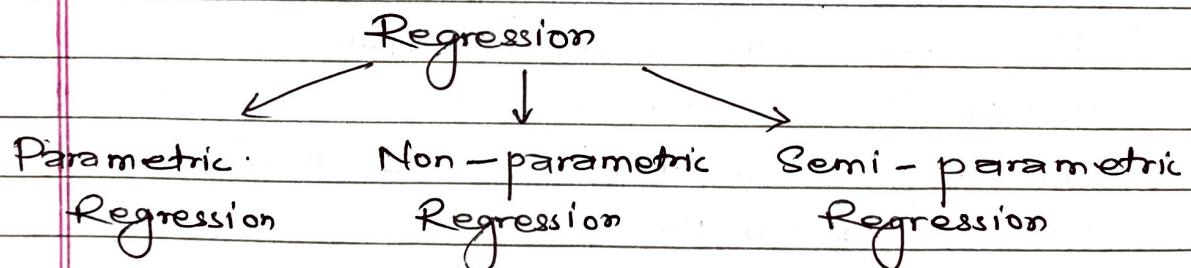
In multiple linear regression, we assume that there is a linear relationship between dependent variable y and set of ' p '

↳ Independent variables x_1, x_2, \dots, x_p

Multiple - More than one IV.

Linear - Linear relationship

The main advantage of multiple linear regression is that it allows us to estimate or predict that values of dependent with more informations from the several independent variables.



Model

The multiple linear regression model with "p" independent variable is given by,

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + e$$

The multiple linear regression model with two independent variable is given by,

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + e$$

where,

Y = Dependent Variable

x_1 = First Independent variable

x_2 = Second Independent variable

β_0 = Pop" y -intercept

β_1 = Pop" slope of y with x_2 holding effect of x_1 constant

β_2 = Pop" slope of y with x_1 holding effect of x_2 constant

e = Error of prediction

Observed y Estimated

$$Y = Y - \hat{Y}$$

Deterministic Model - No error

~~#~~ Statistical Model — with error

Equation

$$\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2, e = y - \hat{Y}$$

$$\hat{Y} = b_0 + b_1 X_1 + b_2 X_2$$

$\downarrow \quad \downarrow \quad \downarrow$
 $b_0 \quad b_1 \quad b_2$

$$E(\hat{\beta}_0) = E(b_0) = \beta_0$$

$$E(\hat{\beta}_1) = E(b_1) = \beta_1$$

$$E(\hat{\beta}_2) = E(b_2) = \beta_2$$

where,

\hat{Y}_i = Estimated value of dependent variable

X_{i1} = first independent variable

x_{12} = Second independent variable

b_0 : Sample y-intercept

b_1 = Sample slope of y on x_1 , holding the effect of x_2 constant

$b_2 = " \quad " \quad " \quad " \quad y \text{ on } x_2 \quad "$
" $x_1 \quad "$

Interpretation of b_0

b_0 measures the average value of γ for $x_1 = 0, x_2 = 0$, it can be either +ve or -ve.

Interpretation of b_1

b) measure the average value of y for unit charge in X_1 , keeping the effect of X_2 constant. It can be either positive or negative.

Interpretation of b_2

The b_2 measures the average values of Y for unit change in X_2 keeping the effect of X_1 constant. It can be either +ve or -ve.

Estimating the regression coefficient

The values of regression coefficient $\hat{B}_0 = b_0$, $\hat{B}_1 = b_1$, $\hat{B}_2 = b_2$ are not known; and they are estimated from sample data.

The two common methods for estimating the regression coefficients are:-

- (1) Ordinary least square method (OLS method)
- (2) Maximum likelihood estimation method (MLE method)

Ordinary least-square method

The principle behind least square method is to find the optimum value of b_0 , b_1 and b_2 so that the sum of residuals or error is minimum.

Residual or Error(e_i) = $Y_i - \hat{Y}_i$, $i = 1, 2, \dots, n$

$$\text{Let } S = \sum_{i=1}^n e_i^2$$

$$= \sum_{i=1}^n \{Y_i - \hat{Y}_i\}^2$$

$$= \sum_{i=1}^n \{Y_i - b_0 - b_1 X_{i1} - b_2 X_{i2}\}^2$$

The least square method minimum as function partially differentiating function wrt b_0 , we get,

$$\frac{\partial S}{\partial b_0} = 0 = \frac{\partial}{\partial b_0} \left\{ \sum_{i=1}^n (y_i - b_0 - b_1 x_{i1} - b_2 x_{i2})^2 \right\} = 0$$

$$\Rightarrow \sum y_i = n b_0 + b_1 \sum x_{i1} + b_2 \sum x_{i2}$$

Partially diff'n s wrt b_1 , we get,

$$\frac{\partial S}{\partial b_1} = 0 \Rightarrow \sum_{i=1}^n y_i x_{i1} = b_0 \sum_{i=1}^n x_{i1} + b_1 \sum_{i=1}^n x_{i1}^2 + b_2 \cdot \sum_{i=1}^n x_{i1} x_{i2}$$

Finally, partially diff'n s wrt b_2 , we get,

$$\frac{\partial S}{\partial b_2} = 0 \Rightarrow \sum_{i=1}^n y_i x_{i2} = b_0 \sum_{i=1}^n x_{i2} + b_1 \sum_{i=1}^n x_{i1} x_{i2} + b_2 \sum_{i=1}^n x_{i2}^2$$

\therefore The normal eqns of estimating b_0 , b_1 & b_2 are,

$$\sum y_i = n b_0 + b_1 \sum x_{i1} + b_2 \sum x_{i2}$$

$$\sum y_i x_{i1} = b_0 \sum x_{i1} + b_1 \sum x_{i1}^2 + b_2 \sum x_{i1} x_{i2}$$

$$\sum y_i x_{i2} = b_0 \sum x_{i2} + b_1 \sum x_{i1} x_{i2} + b_2 \sum x_{i2}^2$$

Matrix form

$$\begin{bmatrix} \sum y_i \\ \sum y_i x_{i1} \\ \sum y_i x_{i2} \end{bmatrix} = \begin{bmatrix} n & \sum x_{i1} & \sum x_{i2} \\ \sum x_{i1} & \sum x_{i1}^2 & \sum x_{i1} x_{i2} \\ \sum x_{i2} & \sum x_{i1} x_{i2} & \sum x_{i2}^2 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix}$$

Regression modelling of annual salary data

$$\rightarrow y = \text{Annual salary } (\times 1000 \$)$$

x_1 = years of schooling

x_2 = years of education

The regression eqn of salary on yrs of schooling and yrs. of experience is,

$$\hat{Y}_i = b_0 + b_1 X_{i1} + b_2 X_{i2} \rightarrow \textcircled{1}$$

Computational Result

$$\begin{array}{lll} \sum Y_i = 634 & \sum Y^2 = 36086 & \sum Y X_1 = 10036 \\ \sum X_{i1} = 183 & \sum X_{i1}^2 = 2869 & \sum X X_2 = 5865 \\ \sum X_{i2} = 102 & \sum X_{i2}^2 = 1155 & \sum X_1 X_2 = 1569 \end{array}$$

$$n = 12$$

∴ The normal eqns for the estimating b_0 , b_1 , and b_2 are,

$$\sum Y_i = n b_0 + b_1 \sum X_{i1} + b_2 \sum X_{i2} \quad \textcircled{II}$$

$$\sum Y_i X_{i1} = b_0 \sum X_{i1} + b_1 \sum X_{i1}^2 + b_2 \sum X_{i1} X_{i2} \quad \textcircled{III}$$

$$\sum Y_i X_{i2} = b_0 \sum X_{i2} + b_1 \sum X_{i1} X_{i2} + b_2 \sum X_{i2}^2 \quad \textcircled{IV}$$

Substituting the value in the above eqn we get,

$$634 = 12 b_0 + 183 b_1 + 102 b_2 \rightarrow \textcircled{V}$$

$$10036 = 183 b_0 + 2869 b_1 + 1569 b_2 \rightarrow \textcircled{VI}$$

$$5865 = 102 b_0 + 1155 b_1 + 1569 b_2 \rightarrow \textcircled{VII}$$

Solving \textcircled{V} , \textcircled{VI} & \textcircled{VII} , we get,

$$b_0 = -33.005$$

$$b_1 = 4.7321$$

$$b_2 = 1.5925$$

The fitted eqn is,

$$\hat{Y} = -33.005 + 4.7321 X_1 + 1.5923 X_2$$

Estimated annual salary = $-33.0005 + 4.7321 \times \text{year of schooling} + 1.5925 \times \text{year of experience}$

Interpretation $b_0 = -33.005$

When $X_1 = 0, X_2 = 0, \hat{Y} = -33.0005$ (not interpretable)

Interpretations of $b_1 = 4.7321$

For every increment of yrs. of schooling by an amount of 1 yrs., the annual salary increases in average by an amount of $4.7321 (\times 1000 \$) = \4732.1 , provided that the effect of yrs. of experience is kept constant.

Interpretations of $b_2 = 1.5925$

For every increment of yrs. of schooling by an amount of 1 yrs., the annual salary increases in average by an amount of $1.5925 (\times 1000 \$) = \1592.5 , provided that the effect of yrs. of experience is kept constant.

Prediction

What is the likely annual salary for the person who has completed 14 yrs. of education and has 6 yrs. of experience?

Soln:-

$$\hat{Y} = -33.0005 + 4.7321X_1 + 1.5925X_2$$

Put,

$$X_1 = 14, X_2 = 6$$

$$\begin{aligned}\hat{Y} &= -33.0005 + 4.7321 \times 14 + 1.5925 \times 6 \\ &= 42.8029 (\times 1000 \$) \\ &= \$42802.9\end{aligned}$$

Alternate methods for estimating regression coefficient

1. Let,

$$W = \begin{vmatrix} 1 & r_{y_1} & r_{y_2} \\ r_{y_1} & 1 & r_{12} \\ r_{y_2} & r_{12} & 1 \end{vmatrix}$$

$$b_1 = \frac{s_y}{s_1} \begin{vmatrix} r_{y_1} & r_{12} \\ r_{12} & 1 \end{vmatrix} \quad b_2 = \frac{s_y}{s_2} \begin{vmatrix} 1 & r_{y_2} \\ r_{12} & r_{y_2} \end{vmatrix} \quad s_1 = S.D \text{ of } Y$$

$$\begin{vmatrix} 1 & r_{y_2} \\ r_{12} & r_{y_2} \end{vmatrix} \quad s_1 = S.D \text{ of } X_1$$

$$\begin{vmatrix} 1 & r_{12} \\ r_{12} & 1 \end{vmatrix} \quad s_2 = S.D \text{ of } X_2$$

2. Matrix Approach :-

$$\text{Let } Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{bmatrix} \quad n \times 1 \quad \beta = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} \quad 3 \times 1$$

$$\beta = (X'X)^{-1}(X'Y)$$

$$X = \begin{bmatrix} 1 & X_{11} & X_{12} \\ 1 & X_{21} & X_{22} \\ \dots & \dots & \dots \\ 1 & X_{n1} & X_{n2} \end{bmatrix} \quad n \times 3$$

$$\beta_3 = \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} \quad 3 \times 1$$

Evaluation of filled equation:-

The filled equation of multiple linear regression must be estimated before it is used for prediction purpose.

The following are the methods of calculating filled equations:-

- (1) Multiple correlation coefficient
- (2) Coefficient of multiple determination
- (3) Standard Error of estimate
- (4) F - test for whole world
- (5) t - test for individual slope
- (6) Residual Analysis.



Multiple Correlation Coefficient

The correlation between observed Y values with its estimated values \hat{Y} as given by multiple linear regression eqn $\hat{Y} = b_0 + b_1 X_1 + b_2 X_2$ is called multiple correlation coefficient. It is denoted by $R_{Y,12}$

$$R_{y_12} = \frac{\text{Cov}(Y, \hat{Y})}{\sqrt{V(X)} \sqrt{V(\hat{Y})}} = \frac{n \sum Y \hat{Y} - \sum Y \cdot \sum \hat{Y}}{\sqrt{n \sum Y^2 - (\sum Y)^2} \sqrt{n \sum \hat{Y}^2 - (\sum \hat{Y})^2}}$$

Alternatively,

$$R_{y_12} = \frac{r_{Y_1}^2 + r_{Y_2}^2 - 2r_{Y_1} \cdot r_{Y_2} \cdot r_{12}}{1 - r_{12}^2}$$

$0 \leq R_{y_12} \leq 1$ | High value of R_{y_12} indicates good fit

- | | | | |
|-------------------|----------|-----------------------------|----------|
| 1. $R_{y_12} = 1$ | 100% fit | 5. $R_{y_12} \rightarrow 1$ | Good fit |
| 2. $R_{y_12} = 0$ | No fit | 4. $R_{y_12} \rightarrow 0$ | Poor fit |

2. Coefficient of multiple determination

then,

$$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \text{Total sum of squares}$$

The SST can be split into two components :-

1. SSR

$$SST = SSR + SSE$$

2. SSE

$$SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2 = \text{Sum of squares due to regression}$$

$$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \text{Sum of squares due to errors}$$

If regression fits well, large portion SST is due to SSR. Then, ratio of SSR to the SST is a good measure of aptness of the fitted eqn.

Coefficient of determination is given by,

$$R^2 y_{12} = \frac{SSR}{SST} \quad | \quad 0 < R^2 y_{12} \leq 1$$

Again,

$$SST = SSR + SSE$$

SST } Given in
SSR } exam
SSE }

$$\therefore R^2 y_{12} = \frac{SSR}{SST} = \frac{SST - SSE}{SST} = 1 - \frac{SSE}{SST}$$

$$R^2 y_{12} = 1 - \frac{SSE}{SST}$$

Alternatively,

$$R^2 y_{12} = \frac{r_{y_1}^2 + r_{y_2}^2 - r_{y_1} \cdot r_{y_2} \cdot r_{12}}{1 - r_{12}^2}$$

y = Annual salary

x_1 = Yrs. of schooling

x_2 = Yrs. of experience

$$\hat{y} = -33.0005 + 4.7321x_1 + 1.5925x_2$$

$$r_{y_1} = 0.816$$

$$r_{y_2} = 0.505$$

$$r_{12} = -0.012$$

$$R^2 y_{12} = 0.926 = 92.6\%$$

About 93% variation in annual salary explained by variation in yrs. of schooling and yrs of experiences. i.e:

the linear rel'n bet'n annual salary . yrs. of schooling and yrs. of experience. Remaining 7% variation in annual salary given by the factors/variables other than the yrs. of schooling and yrs. of experience.

3. Standard Error of Estimate

The standard error of estimate is the deviation of the residuals and it is denoted by,

$$S_e = \sqrt{\frac{SSE}{n-p-1}}$$

$$= \sqrt{\frac{SSB}{n-2-1}}$$

$$= \sqrt{\frac{SSE}{n-3}}$$

$p=2$ In our study
(X_1, X_2)

SSE = Error sum of squares
 $= \sum_{i=1}^n (X_i - \hat{X}_i)^2$

Low value of S_e indicates good fit and high value of S_e indicates poor fit.

$$0 \leq S_e < \infty$$

$S_e = 0$ perfect fit 100% fit

S_e measures the variability of points around the regression plane.

4. Overall or global test of model accuracy:-

The test evaluate the fitted eqn. as whole i.e. it determines whether there is linear relationship betⁿ dependent variable Y and independent X_1 and X_2 .

$H_0: \beta_1 = \beta_2 = 0$ (There is no relⁿ betⁿ Y , X_1 and X_2)

$H_1: \beta_i \neq 0$ (for at least one $i = 1, 2$)

Test Statistics

The test statistic is given by,

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

MSR = Variable due to regression

$$= \frac{SSR}{p} \quad | \quad p = \text{No. of independent variable in the model}$$

$$= \frac{SSR}{2} \quad | \quad SSR = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

MSE = Variance due to error

$$= \frac{SSE}{n-p-1} \quad | \quad SSE = \sum_{i=1}^n (y_i - \hat{Y}_i)^2$$

$$= \frac{SSE}{n-3}$$

ANOVA Table

Source of variation	Degrees of freedom	Sum of squares	Mean sum of squares	F - ratio (cal F)
Regression	$p=2$	$SSR = \sum_{i=1}^n (Y_i - \bar{Y})^2$	$MSR = \frac{SSR}{2}$	$F = \frac{MSR}{MSE}$
Error	$n-p-1$ $= n-2-1$ $= n-3$	$SSE = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$	$MSE = \frac{SSE}{n-3}$	
Total	$n-1$	$SST = \sum_{i=1}^n (Y_i - \bar{Y})^2$	\bar{MST}	

Decision Rule:-

Reject H_0 if $cal{F} > F_{\alpha}(2, n-3)$

↙ ↘
 numerator Denominator
 d.f d.f

e.	y	x_1	x_2
52	16	6	52. 2686
44	12	10	39. 7102
48	13	15	52. 4048
77	20	8	74. 3820
68	18	11	69. 6953
48	16	2	45. 8986
59	14	12	52. 3594
53	18	4	58. 5478
28	12	6	33. 3402
61	16	9	57. 0461
27	12	2	26. 1702
69	16	18	71. 3786

$$\hat{y} = -33 + 4.7821x_1 + 1.5925x_2$$

$$(y - \hat{y})^2$$

0.0714596
18. 40238404
19. 40226304
6. 853924
2. 87404209
4. 41588196
44. 09756836
80. 77808484
28. 51773604
15. 63332522
0. 00088804
5. 65773796

$$(\hat{y} - y)^2$$

0. 3189
172. 2158
0. 1836
464. 3465
284. 3270
48. 0901
0. 2246
32. 6555
379. 9809
17. 7477
668. 8999
343. 9282

$$SSR = 2412.9187$$

$$SSE = 176.7059$$

$$SST = 2589.6246$$

$$SSE = 176.7059$$

$$SSR = 2412.9187$$

$$\bar{Y} = \frac{\sum Y}{n} = 52.833$$

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

$$= \frac{2412.9187}{2589.6246}$$

$$= 0.9317$$

$$= 93.17\%$$

F-test for whole model

(Overall test of model accuracy)

Step 1:- Null and Alternative Hypothesis

$H_0: \beta_1 = \beta_2 = 0$ (There is no linear rel' bet' γ , x_1 & x_2)

(x_1 and x_2 are not significant prediction for γ)

$H_1: \beta_j \neq 0$ for at least one $j=1, 2$

(At least one indep. variable is significant prediction for γ)

Step 2: Choice of α for the test

$$\alpha = 5\%$$

Step 3: Test statistic

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

, where $MSR = SSR/p$ $p=2$

$$MSF = SSE/n-p-1$$

The test statistic, F follows F distn with $p=2$
d.f in numerator and $n-p-1=n-3$ d.f in denominator.

$$F \sim F_{\alpha} \{2, n-3\}$$

Step 4: Calculated F

Source of variation	Degrees of freedom	Sum of squares	Mean sum of squares	Cal F
Regression	$p=2$	$SSR = 2412.91352$	$MSR = \frac{2412.91352}{2} = 1206.459$	$F = \frac{MSR}{MSE}$
Error	$n-p-1$ $= 12-2-1$ $= 9$	$SSE = 176.7059$	$MSE = \frac{176.7059}{9} = 19.6339$	
Total	$n-1 = 12-1$ $= 11$	$SST = 2589.6246$		

$$\therefore \text{Cal } F = 61.4478$$

Step 5: Tab F

Numerator d.f = 2

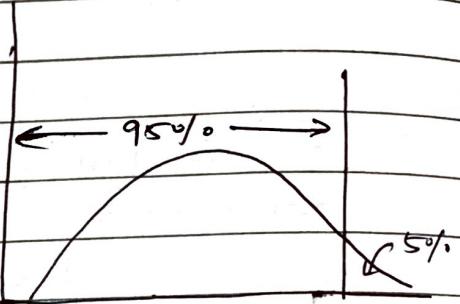
Denominator d.f = 9

$$\alpha = 5\%$$

$$F_{0.05} (2, 9) = 4.26$$

AR:- $F < 4.26$

RR:- $F \geq 4.26$



Step 6:- Statistical Decision

Since, $\text{Cal } F = 61.4478 \gg \text{Tab } F = 4.26$, we strongly rejected H_0 at 5% level of significance.

Step 7:- Conclusion

At least one slope is significant i.e. at least one IV is significant prediction of Y.

t - test for individual slope

If f-test test of overall goodness of model reject H_0 then at least one IV is significant for Y . To know which I.V is significant for Y , we have to perform t - test for individual slope.

Hypothesis

$$H_0: \beta_j = 0 \quad (x_j \text{ is not significant for } Y)$$

$$H_1: \beta_j \neq 0 \quad (x_j \text{ is significant for } Y)$$

$$j = 1, 2$$

$$j = 1 \quad x_1$$

$$j = 2 \quad x_2$$

Test Statistic

$$t = \frac{b_j}{SE(b_j)} \quad \left| \quad \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)}$$

where,

b_j = Sample slope of Y on x_j holding others

$H_1: \beta_1 \neq 0$ (There is linear relationship, i.e. x_1 is significant of y)

Step 2:- Choice of α

$$\alpha = 5\%$$

Step 3:- Test Statistic

$$t = \frac{b_1}{SE(b_1)} \quad | \quad SE(b_1) = \sqrt{\frac{MSE}{SSx_1}}$$

$SSx_1 = \sum(x_1 - \bar{x}_1)^2$

Given

The test statistic has t dist' with $n-p-1 = n-3$ df

Step 4:- Calculated F

$$MSE = 19.6339$$

$$SSx_1 = \sum(x_1 - \bar{x}_1)^2 \quad | \quad (x - \bar{x})^2$$

$\sum x_1 = 18.3$

$\bar{x}_1 = 15.25$

$E(x_1 - \bar{x}_1)^2 = 78.25$

$$SE(b_1) = \sqrt{\frac{MSE}{SSx_1}}$$

$$= \sqrt{\frac{19.6339}{78.25}}$$

$$= 0.500911629$$

Now,

$$t = \frac{b_1}{\text{SE}(b_1)}$$

$$= \frac{4.7321}{0.5009}$$

$$= 9.447$$

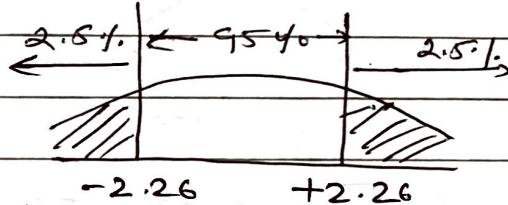
Step 5 :- Tabulated t

The test is two sided, and $\alpha = 5\%$.

$$\text{D.f} = n - 3$$

$$= 12 - 3$$

$$= 9$$



$$\text{AR: } -2.26 < t < +2.26$$

$$\text{RR: } t \geq 2.26 \text{ or } t \leq -2.26$$

Step 6 :- Statistical Decision

Since, cal t = 9.447 > upper critical t = 2.26
we reject H₀.

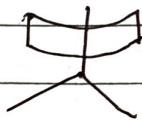
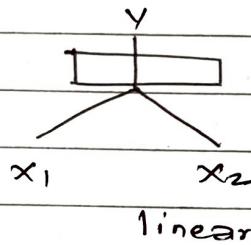
Step 7 :- Conclusion

The variable x₁ is significant predictor for y.

Assumption of multiple linear regression

1. Variables have non-zero variance.
2. Variables are measured without error.
3. There must be a linear relationship between dependent variable and independent variable.

y, x_1, x_2



Non-linear

linear

4. The errors are normally and independently distributed.

5. Equal variance Assumption (Homoscedasticity)

This assumption states that the variance of error terms are similar across the values of the independent variables.

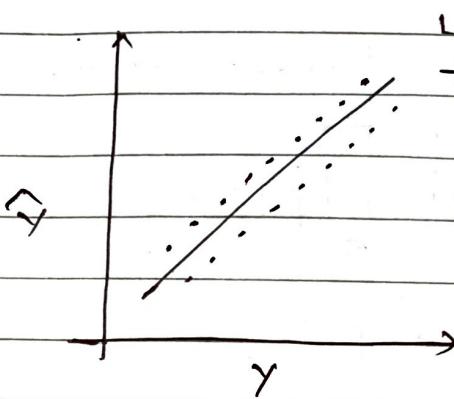
6. Non multicollinearity states the independent variables are not highly correlated with each other.

Residual Analysis:-

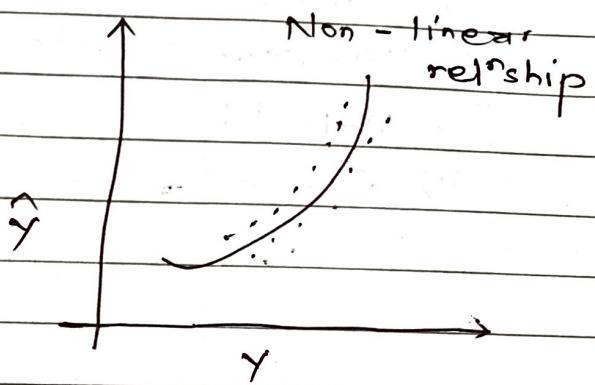
Linearity Assumption:-

This assumption is test by using following two graphs:-

- (a) A graph of γ vs $\hat{\gamma}$
 (b) A graph of e vs \hat{y}



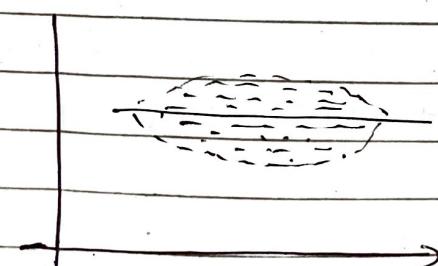
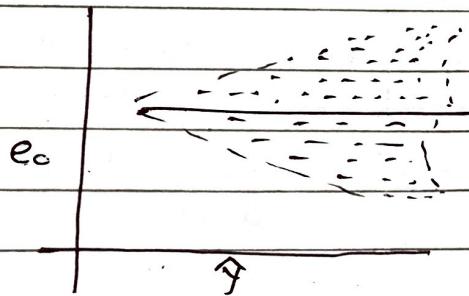
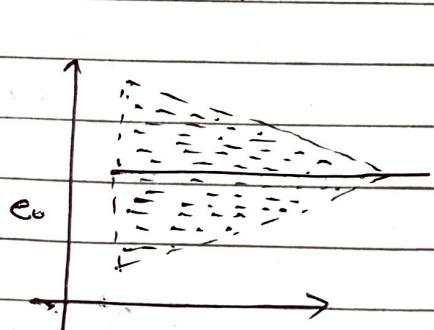
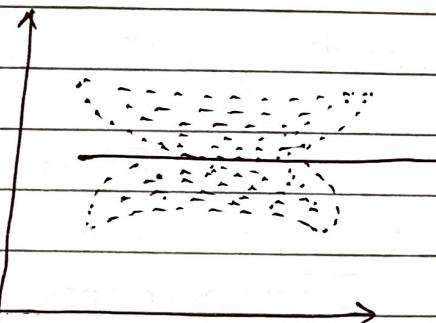
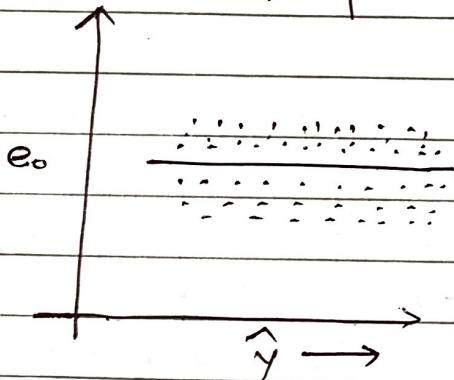
Linear
Relationship



Non - linear
relationship

Homoscedasticity
Linear Relationship

Non linear rel^{se}ship
Heteroscedasticity



Normality of error:-

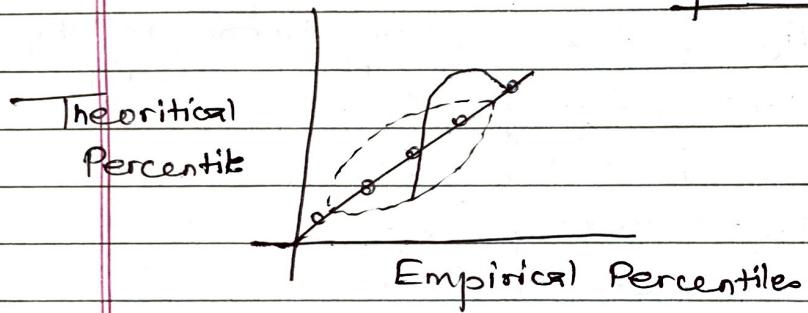
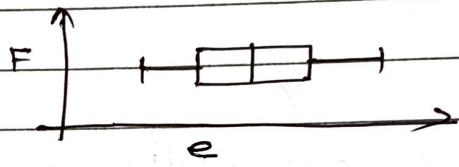
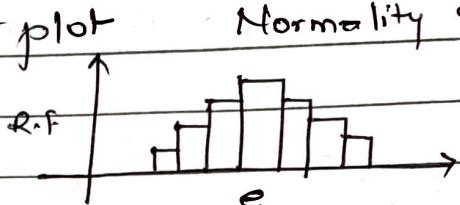
(1) Histogram

(2) Box and whisker plot

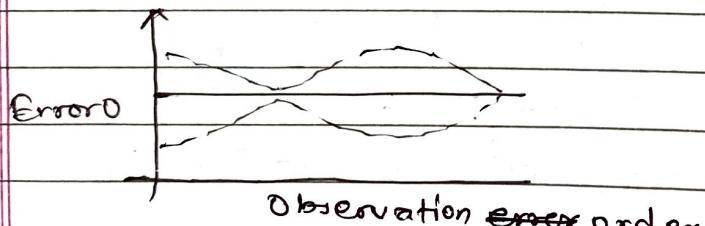
{(3) P-P Plot

(4) Q-Q Plot

Normality of error



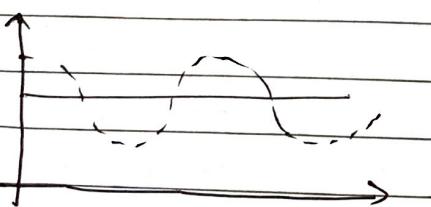
Independence of Errors:-



y-axis → Error

x-axis → Observation

order



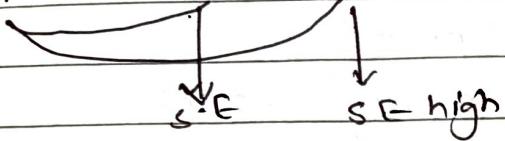
Autocorrelation

or
Serial correlation

Dependence of
error

Non-^{multi}collinearity

$$Y = b_0 + b_1 X_1 + b_2 X_2$$



Multi-collinearity

It is desirable that there is a strong correlation b/w dependent variable and independent variable. Multicollinearity is the state in which there is a moderate to high inter correlations or interassociations among the independent variables (predictors). If multicollinearity exists two or more predictors the reliability of the regression coefficients is reduced and therefore prediction and statistical inferences are not reliable.

Effects of multi-collinearity

- 1. The standard error of regression coefficient are likely to be high. There will be change in signs as well magnitude of regression coefficients from sample to sample.
- 2. The estimated regression coefficients of any one variable depends on which other independent variables are included in the model.