

## Unit 5: Correlation Analysis

If there exists some relationship b/w two variables then the statistical analysis of such data is called bivariate analysis.

Correlation refers to the relationship of two or more variables.

Correlation is a statistical analysis which measures and analyzes the degree or extent to which two variables fluctuate with reference to each other.

The correlation expresses the relationship of two sets of variables upon each other.

Types of correlation Correlation is classified into many types,

(1) +ve & -ve correlation +ve & -ve correlation depends upon the direction of change of variables. If two variables tends to move together in the same direction i.e. an increase in the value of one variable affected by the increase in the value of other variable or a decrease in the value of one variable is affected by a decrease in the value of another variable then the correlation is called +ve or direct correlation. If two variables tends to move together in opp directions so that an increase or decrease in the value of one variable is affected by a decrease or increase in the value of other variable then the correlation is called -ve or inverse correlation.

(2) Simple & multiple correlations When we study only 2 variables the relationship is described as correlation. But in multiple correlation we study more than two variables simultaneously.

- ③ Partial & Total Correlation: The study of 2 variables excluding some other variables is called partial correlation.
- In total correlation all the facts are considered in problem.
- ④ Linear & Non Linear Correlation: If the ratio of change in two variables is uniform then there will be a linear correlation b/w them.
- In a non-linear or curvilinear correlation the amount of change in one variables do not bear a constant ratio of the amount of change in other variable.

### Methods of studying Correlation

- ⑤ Carl Pearson's Coeff of Correlation: This is mathematical method for measuring the magnitude of linear relationship b/w two variables. It is denoted by
- $$r(x,y) = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sqrt{\sum (x-\bar{x})^2 \cdot \sum (y-\bar{y})^2}}$$

### Properties of Coeff of Correlation

- ① Coeff of correlation lies b/w -1 & +1 i.e. the limits for correlation are ~~-1~~  $\pm 1$ .
- ② If  $r=1$ , then the two variables are perfectly +ve if  $r=-1$ , then the correlation is perfectly -ve.
- If  $r>0$ , then there is +ve correlation b/w two variables.
- If  $r<0$ , then there is a -ve correlation b/w two variables.
- If  $r=0$ , then two variables are uncorrelated.

Prob Psychological  
 1. A ~~psychological~~ test & intelligence and engineering ability were applied to 10 students. Here is a record of data showing intelligence ratio & engineering ratio. Cal the coeff of Correlation.

intelligence ratio & eng ratio  
 intelligence ratio = 102 101 100 99 98 95 96 104 92 97 94  
 eng ratio = 101 103 100 98 95 96 96 104 92 97 94

Ques. Is there a linear relationship between intelligence ratio & eng ratio?

Let  $x$  = intelligence ratio

$y$  = eng ratio

$$\Sigma x = 99$$

$$(x - \bar{x}) = 6 \quad 5 \quad 3 \quad 2 \quad 1 \quad 0 \quad -1 \quad -3 \quad -6 \quad -7 \quad -8$$

$$(y - \bar{y}) = 3 \quad 5 \quad 2 \quad 0 \quad -3 \quad -2 \quad 6 \quad -6 \quad -1 \quad -4 = 0$$

$$(x - \bar{x})^2 = 36 \quad 25 \quad 9 \quad 4 \quad 1 \quad 0 \quad 1 \quad 9 \quad 36 \quad 49 = 170$$

$$(y - \bar{y})^2 = 9 \quad 25 \quad 4 \quad 0 \quad 9 \quad 4 \quad 36 \quad 1 \quad 16 = 140$$

$$\cancel{(x - \bar{x})(y - \bar{y})} =$$

$$(x - \bar{x})(y - \bar{y}) = 18 \quad 25 \quad 6 \quad 0 \quad -3 \quad 0 \quad -6 \quad 18 \quad 6 \quad 28 = 92$$

$$r(x, y) = \frac{-92}{\sqrt{170 \cdot 140}} = \frac{-92}{\sqrt{23800}} = \frac{-92}{153.88} = -0.5963$$

Q. Cal. the coeff of correlation to the following data

x	38	45	46	38	35	38	46	32	36	38
y	28	34	38	36	36	26	28	29	25	36

(S)

$$\bar{x} = 39.2$$

$$\bar{y} = 31.4$$

$$(x-\bar{x}) = -1.2 \quad 5.8 \quad 6.8 \quad -1.2 \quad -4.2 \quad -1.2 \quad 6.8 \quad -7.2 \quad -3.2 \quad -1.2$$

$$(y-\bar{y}) = -3.4 \quad 2.6 \quad 6.6 \quad -2.6 \quad 4.6 \quad -5.4 \quad 3.4 \quad -2.4 \quad -6.4 \quad 4.6$$

$$(x-\bar{x})^2 = 1.44 \quad 33.64 \quad 46.24 \quad 1.44 \quad 17.64 \quad 1.44 \quad 46.24 \quad 51.84$$

$$10.24 \quad 1.44 = 211.6$$

$$(y-\bar{y})^2 = 11.56 \quad 6.76 \quad 43.56 \quad 6.76 \quad 21.16 \quad 29.16 \quad 11.56 \quad 5.76 \quad 40.96$$

$$21.16 = 198.4$$

$$(x-\bar{x})(y-\bar{y}) = 4.08 \quad 15.08 \quad 44.88 \quad -19.32 \quad 6.48 \quad -23.12 \quad 17.28 \quad 20.48 \quad -5.52$$

$$r(x,y) = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sqrt{\sum (x-\bar{x})^2 \cdot \sum (y-\bar{y})^2}} = \frac{57.2}{\sqrt{(211.6)(198.4)}} = 0.2792$$

Spearman's Rank correlation. This method is based on ranks and is useful in dealing with qualitative characteristics such as intelligence, honest, kindness etc. It cannot be measured quantitatively as increase of Karl Pearson's coeff of correlation. It is based on the ranks given to the observations. The formula for Spearman's Rank correlation coeff is given by

$$P = 1 - \frac{6 \sum d^2}{n(n^2-1)}$$

In case common ranks are given to the repeated values then we can use the following formula after cal. the correction factor & each repeated values

$$\therefore P = 1 - \frac{c[\sum d^2 + CF]}{n(n^2-1)}$$

### Properties of Rank correlation

- (1) The value of Spearman's rank correlation is always lies b/w -1 & +1
- (2) If  $r_s = 1$  then there is a complete agreement in the order of the ranks & the directions of the ranks are same.
- (3) If  $r_s = -1$  then there is a complete disagreement in the order of the ranks & the direction of ranks are opposite.

Prob

Find rank correlation to the following data

$$x_r = 75, 95, 85, 80, 66, 52, 99, 72$$

$$y_r = 86, 90, 95, 84, 72, 99, 55, 64$$

Sol

$$n = 8$$

$$R_1 = 5, 2, 3, 4, 7, 8, 1, 6$$

$$R_2 = 1, 3, 2, 5, 6, 1, 8, 7$$

$$d = R_1 - R_2 = 1, -1, 1, 1, 7, -7, 1$$

$$d^2 = 1, 1, 1, 1, 49, 49, 1 = 104$$

$$P = 1 - \left[ \frac{6 \sum d^2}{n(n^2-1)} \right] = 1 - \left( \frac{6 \cdot 24}{504} \right) = 1 - 0.2881$$

$R_1 \rightarrow$  highest ranks /  
 & low ranks of x  
 $R_2 \rightarrow$  ranks / highest rank  
 of y

Spearman

2. Cal rank correlation coeff of following data

xr 100 105 95 100 100 95 107 132 150 132

yr 95 100 100 96 96 95 105 96 145 136

xr	yr	$R_1$	$R_2$	$d = R_1 - R_2$	$d^2$
100	95	7	9.5	-2.5	6.25
105	100	5	2.5	0.5	0.25
95	100	9.5	2.5	7	49
100	96	7	7	0	0
100	96	7	7	0	0
95	95	9.5	9.5	0	0
107	105	4	3	-1	1
132	96	2.5	7	-4.5	20.25
150	145	1	1	0	0
132	136	2.5	2	-0.5	0.25

② In xr series 132 is repeated 2 times

$$\text{Correction factor} = \frac{1}{12}(n^3 - n)$$

In xr series 100 is repeated 3 times

$$M = 3, C.F. = \frac{1}{12}(n^3 - n) = 24$$

In x series 95 is repeated 2 times

~~m=2, CF=0.5~~

In y series 100 is repeated 2 times

~~m=2, CF=0.5~~

In y series 96 is repeated 3 times

~~m=3, CF=2/3 = 0.666~~

In y series 95 is repeated 2 times

~~m=2, CF=0.5~~

NOW Spearman rank correlation coeff

$$\rho = 1 - \frac{6(2d^2 + CF)}{n(n^2 - 1)}$$

$$= 1 - \left[ \frac{6(53 + 0.666)}{10(100 - 1)} \right] = 0.6424$$

3. Cal S.R.C coeff to the following data

x	y	R <sub>1</sub>	R <sub>2</sub>	d = R <sub>1</sub> - R <sub>2</sub>	d <sup>2</sup>	n=12
95	68	9	5.5	-3.5	12.25	
63	66	11	9	1.5	2.25	
67	68	6.5	5.5	1	1	
64	65	10	11.5	-1.5	2.25	
68	69	6.5	3	1.5	2.25	
62	66	12	9.5	2.5	6.25	

70	68	2	5.5	-3.5	12.25
66	65	8	10.5	-3.5	12.25
68	71	4.5	4	3.5	12.25
67	67	6.5	8	-1.5	2.25
69	68	3	5.5	-2.5	6.25
71	70	1	2	-1	1
					72.5

In x series 68 is repeated 2 times

$$m=2$$

$$CF = \frac{1}{12^2} (m^3 - m) = 0.5$$

In x series 67 is repeated 2 times

$$m=2, CF=0.5$$

In y series 68 is repeated 4 times

$$m=4, CF=5$$

In y series 66 is repeated 2 times

$$m=2, CF=0.5$$

In y series 65 is repeated 2 times

$$m=2, CF=0.5 = 7$$

$$\ell = 1 - \frac{6(\Sigma d^2 + CF)}{n(n^2-1)} = 1 - \left( \frac{6(72.5+7)}{12(12^2-1)} \right) \\ = 0.722$$

Prob

1. 10 competitors in a musical test were ranked by 3 judges A B C in the following order using rank correlation method discuss which pair of judges has the nearest approach to common likings in music.

(RA) ranks given by judge A + 1 6 5 10 3 2 4 9 7 8

(RB) judge B + 3 5 8 4 7 10 2 1 6 9

(RC) judge C + 6 4 9 8 11 2 10 5 7 3

$$d_1 = R_A - R_B = d_1^2 = 4 + 1 + 9 + 36 + 16 + 64 + 4 + 64 + 1 + 1 = 200$$

$$d_2 = R_B - R_C = d_2^2 = 9 + 1 + 16 + 36 + 64 + 64 + 16 + 1 + 36 = 244$$

$$d_3 = R_C - R_A = d_3^2 = 25 + 4 + 16 + 4 + 0 + 36 + 16 + 0 + 25 = 130$$

$$\rho(A, B) = 1 - \frac{6 \sum d_1^2}{n(n^2-1)} = 1 - \frac{1200}{10(10^2-1)} = 1 - 0.2121$$

$$\rho(B, C) = 1 - \frac{6 \sum d_2^2}{n(n^2-1)} = 1 - \frac{1464}{990} = 1 - 0.4788$$

$$\rho(C, A) = 1 - \frac{6 \sum d_3^2}{n(n^2-1)} = 1 - \frac{780}{990} = 0.2121$$

The rank correlation of C & A has the max value so we conclude that the judges A & C has the nearest approach to common likings in music.

Regression analysis The study of correlation measures the direction of the relationship b/w two variables i.e. correlation we cannot estimate the value of 1 variable when the value of other variable is given in regression we can estimate the value of 1 variable with the value of other variable which is known as the statistical method which helps us to estimate the unknown value of 1 variable from the known value of related variable is called regression.

Regression is an avg relationship b/w two variables if the line described in avg relationship b/w 2 variables is known as a line of regression.

Regression eq's, Regression eq is an algebraic expression of the regression lines. Reg may be classified into 2 types

① Regression eq of  $y$  on  $x$  can be written as  $y - \bar{y} = b_{yx}(x - \bar{x})$

where  $b_{yx}$  is the P. coeff of  $y$  on  $x$  and is  $(b_{yx} = r \cdot \frac{\sigma_y}{\sigma_x})$

defined as  $\rightarrow$   $b_{yx} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$

$(x - \bar{x})(y - \bar{y})$  is known as the product moment of  $x$  &  $y$

② R.eq of  $x$  on  $y$  can be written as  $(x - \bar{x}) = b_{xy}(y - \bar{y})$  where  $b_{xy}$  is the P. coeff of  $x$  on  $y$  and is defined

$$as \quad b_{xy} = r \cdot \frac{\sigma_x}{\sigma_y}$$

$$\gamma = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \times \sum (y - \bar{y})^2}}$$

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

$$\bar{y} = \sqrt{\frac{\sum (y - \bar{y})^2}{n}}$$

Properties of r-coeff & The angle b/w two regression lines is given by  $\tan \theta = \frac{r}{\sqrt{1-r^2}} \left( \frac{1-r^2}{r} \right)$

If  $r=0$  then  $\tan \theta = \infty$

$$\theta = \frac{\pi}{2}$$

If  $r=\pm 1$  then  $\tan \theta = 0$

$$\theta = 0$$

→ The G.M of two regression coeff is equal to coeff of correlation i.e.  $r_{xy} = \sqrt{r_x r_y}$

→ If 1-r.coeff is +ve then the other R.coeff is also +ve

Prob

1. A Panel of 2 judges P & Q graded 7 performances by independently awarding marks as follows the 8th performance which judge Q could not attend was awarded 37 marks by judge P if judge Q has also been present how many marks would be expected to have been awarded by him to the 8th performance.

(x) marks by judge P = 46 42 44 40 43 41 45

(y) Judge Q = 40 38 36 35 39 37 41

$$\bar{x} = 43$$

$$\bar{y} = 38$$

$$(x-\bar{x}) = 3+1 \quad 1-3 \quad 0-2 \quad -2$$

$$(y-\bar{y}) = 2+0+2+3+1+1+3$$

$$(x-\bar{x})^2 = 9+1+1+9+0+4+4 = 28$$

$$(y-\bar{y})^2 = 4+0+4+9+1+1+9 = 28$$

$$(x-\bar{x})(y-\bar{y}) = 6+0+2+9+0+2+6 = 21$$

$$\textcircled{1} r = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sqrt{\sum (x-\bar{x})^2 \times \sum (y-\bar{y})^2}} = \frac{21}{\sqrt{28 \times 28}} = 0.75$$

$$\textcircled{2} \bar{x} = \sqrt{\frac{28}{7}} = 2$$

$$\textcircled{3} \bar{y} = \sqrt{\frac{\sum (y-\bar{y})^2}{n}} = \sqrt{\frac{28}{7}} = 2$$

$$\textcircled{4} b_{yx} = r \cdot \frac{\bar{y}}{\bar{x}} = 0.75 \cdot \frac{2}{2} = 0.75$$

Regression eq of  $y$  on  $x = y - \bar{y} + b_{yx}(x - \bar{x})$

$$= y - 2 + 0.75x - 32.25$$

$$= y + 0.75x - 32.25 + 38$$

$$= y + 0.75x + 5.75$$

$x$  is given  
in question

$$x = 37$$

$$y = (0.75)37 + 5.75 = 33.75$$

$$= 33.5$$

2. Price index of cotton and price index of wool are given below for 12 months in a year. Fit a linear regression b/w two indexes and also estimate the price index of wool when the price index of cotton is 100 and also estimate the price index of cotton when the price index of wool is 80.

$$(x) \text{P.i of cotton} = 78 \quad 77 \quad 85 \quad 88 \quad 87 \quad 82 \quad 81 \quad 77 \quad 76 \quad 73 \quad 83 \quad 97 \quad 93$$

$$(y) \text{P.i of wool} = 84 \quad 82 \quad 82 \quad 85 \quad 89 \quad 90 \quad 88 \quad 92 \quad 83 \quad 87 \quad 92 \quad 99$$

$$\Sigma = 83.6667, \bar{x} = 88.4167$$

$$(x - \bar{x}) = -5.6667 \quad -6.6667 \quad 1.3333 \quad 4.3333 \quad 3.3333 \quad -1.6667 \quad -2.6667 \\ -6.6667 \quad -7.6667 \quad -0.6667 \quad 13.3333 \quad 9.3333$$

$$(y - \bar{y}) = -4.4167 \quad -6.4167 \quad -6.4167 \quad -3.4167 \quad 0.5833 \quad 1.5833 \\ -0.4167 \quad 3.5833 \quad -5.4167 \quad 0.5833 \quad 9.5833 \quad 10.5833$$

$$(x - \bar{x})^2 = 32.1115 \quad 46.4449 \quad 1.7777 \quad 18.7775 \quad 11.1109 \quad 2.7779 \quad 7.1113$$

$$46.4449 \quad 58.7783 \quad 0.4445 \quad 177.7769 \quad 87.1105 = 486.6679$$

$$(y - \bar{y})^2 = 19.5072 \quad 41.1740 \quad 41.1740 \quad 11.6738 \quad 0.3402 \quad 2.5068$$

$$0.1736 \quad 12.8400 \quad 29.3406 \quad 0.3402 \quad 91.8396 \quad 112.0062 = 302.9162$$

$$(x - \bar{x})(y - \bar{y}) = 25.0281 \quad 42.7782 \quad -8.5554 \quad -14.8656 \quad 1.9443 \quad -2.6389$$

$$1.1112 \quad -23.8888 \quad 41.5282 \quad -0.3889 \quad 127.7740 \quad 98.7771$$

$$= 288.5895$$

$$\textcircled{1} \quad r = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sqrt{\sum (x-\bar{x})^2} \times \sqrt{\sum (y-\bar{y})^2}} = \frac{6665}{\sqrt{282.6675} \times \sqrt{176619.6649}} = 0.6869$$

$$\textcircled{2} \quad \sigma_x = \sqrt{\frac{486.6679}{12}} = 6.3683$$

$$\textcircled{3} \quad \sigma_y = \sqrt{\frac{362.9102}{12}} = 5.4994$$

$$\textcircled{4} \quad b_{yx} = \frac{8.15}{\frac{\sigma_y}{\sigma_x}} = 0.6869 \times \frac{5.4994}{6.3683} = 0.5932$$

$$\textcircled{5} \quad b_{xy} = \frac{8.15}{\frac{\sigma_x}{\sigma_y}} = 0.6869 \times \frac{6.3683}{5.4994} = 0.7954$$

regression eq of y on x =  $y - \bar{y} = b_{yx}(x - \bar{x})$

$$= y - 88.4167 + 0.5932(x - 83.6667)$$

$$= y - 88.4167 + 0.5932x + 49.6311$$

$$= y = 0.5932x + 38.7856$$

$$= x = 100 \\ y = 0.5932(100) + 38.7856$$

x & y are in  
the question

regression eq of x on y =  $x - \bar{x} = b_{xy}(y - \bar{y})$

$$= x - 83.6667 = 0.7954(y - 88.4167)$$

$$= x = 0.7954y - 70.3266 + 83.6667$$

$$= x = 0.7954y + 13.3401$$

$$= y = 80$$

$$x = 76.9721$$

3. The eq of 2 regression lines are  $7x - 16y + 9 = 0$ ,  $5y - 4x - 3 = 0$ .  
find the means of  $x$  &  $y$  and also find corr of correlation

Sol

w.e.t the 2 regression lines passes through the avg of  
 $x$  &  $y$  i.e.  $7\bar{x} - 16\bar{y} + 9 = 0$ ,  $5\bar{y} - 4\bar{x} - 3 = 0$

$$7\bar{x} - 16\bar{y} + 9 = 0$$

$$28\bar{x} - 64\bar{y} + 36 = 0$$

$$\frac{28\bar{x} - 35\bar{y} + 21 = 0}{-29\bar{x} + 15 = 0}$$

$$\bar{y} = 0.5172$$

$$\bar{y} = \frac{15}{29}$$

$$\bar{x} = 7\bar{x} - 16(0.5172) + 9 = 0$$

$$\bar{x} = 0.103$$

Let  $7x - 16y + 9 = 0$  be the eq of  $x$  on  $y$

$$7x = 16y - 9$$

$$x = \frac{16}{7}y - \frac{9}{7}$$

which is in the form of  $x = by + c$

$$by = \frac{16}{7}y$$

Let  $5y - 4x - 3 = 0$  be the eq of  $y$  on  $x$

$$5y = 4x + 3$$

$$y = \frac{4x}{5} + \frac{3}{5}$$

which is in the form of  $y = bx + c$

$$bx = \frac{4}{5}$$

$$r = \sqrt{b_{yx} b_{xy}}$$

$$\Rightarrow \sqrt{\frac{16}{7} \times \frac{4}{5}} = 1.3522$$

$r = 1.3522$  which is incorrect because  $r$  should be  $\pm 1$  or  $-1$  &  $1$  so consider  $7x - 16y + 9 = 0$  as an eq of  $y$  on  $x$

$$5y - 7x - 3 = 0 \text{ as eq of } x \text{ on } y$$

$$-7x - 16y + 9 = 0$$

$$-16y = -7x - 9$$

$$y = \frac{7}{16}x + \frac{9}{16}$$

$$b_{yx} = \frac{7}{16}$$

$$5y - 7x - 3 = 0$$

$$-4x = -5y + 3$$

$$x = \frac{5}{4}y - \frac{3}{4}$$

$$b_{xy} = \frac{5}{4}$$

$$r = \sqrt{b_{yx} b_{xy}} = 0.7395$$

Ques 4. The following data based on 450 students are given for marks in statistics & economics. At a certain examination the mean marks in statistics is 40. The mean marks in economics is 80. The SD of marks in statistics is 12. Variance of marks in economics is 250 & sum of products deviations of marks from their respective means is 62075 estimate the avg marks in economics who obtain 50 marks in statistics.

Sol + Let  $x$  = marks in statistics  
 $y$  = marks in economics

$$N = 450$$

$$\bar{x} = 40, \bar{y} = 80, \sigma_x = 12, \sigma_y = 16$$

$$\sigma^2 = 256 = 16$$

$$\sum (x - \bar{x})(y - \bar{y}) = 42075$$

$$r_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}} = \frac{42075}{\sqrt{12} \times \sqrt{16}} = 0.4870$$

$$b_{yx} = r \cdot \frac{\sigma_y}{\sigma_x} = 0.4870 \times \frac{16}{12} = 0.6493$$

The regression equation of  $y$  on  $x$  is

$$y - \bar{y} = b_{yx}(x - \bar{x})$$

$$y - 80 = 0.6493(x - 60)$$

$$y - 80 = 0.6493x - 25.9720$$

$$y = 0.6493x + 54.0280$$

5. The height of mothers & daughters are given in the following table. Estimate the average height of daughter when the height of mother is 64.5 inches.

(x) height of mother (in inches) - 62 63 64 64 65 66 68 70

(y) daughter + 64 65 67 69 68 71 69 65

$$\bar{x} = 65.25$$

$$\bar{y} = 66.25$$

$$(x-\bar{x}) = -3.25 \quad -2.25 \quad -1.25 \quad -0.25 \quad 0.75 \quad 2.75 \quad 4.75$$

$$(y-\bar{y}) = -2.25 \quad -1.25 \quad -0.25 \quad 0.75 \quad 1.75 \quad 2.75 \quad 4.75$$

$$(x-\bar{x})^2 = 10.5625 \quad 8.0625 \quad 1.5625 \quad 1.5625 \quad 0.0625 \quad 0.5625 \quad 7.5625$$

$$\sum (x-\bar{x})^2 = 22.5625 \quad 69.5$$

$$(y-\bar{y})^2 = 8.0625 \quad 1.5625 \quad 27.5625 \quad 0.5625 \quad 3.0625 \quad 82.5625 \quad 1.5625$$

$$\sum (y-\bar{y})^2 = 13.0625 \quad 7.5625 \quad 69.5$$

$$(x-\bar{x})(y-\bar{y}) = 7.3125 \quad 2.8125 \quad 6.5625 \quad -0.9375 \quad -0.4375 \quad 3.5625 \quad -3.4375$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum (x-\bar{x})(y-\bar{y}) = \frac{13.0625}{8} = 1.6325$$

$$\text{① } r = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \text{Var}(y)}} = \frac{1.6325}{\sqrt{28.5}} = \frac{1.6325}{5.3452} = 0.3059$$

$$\text{② } \sigma_x = \sqrt{\frac{69.5}{8}} = \sqrt{8.6875} = 2.9475$$

$$\text{③ } \sigma_y = \sqrt{\frac{69.5}{8}} = \sqrt{8.6875} = 2.9475$$

$$\text{④ } b_{yx} = r \cdot \frac{\sigma_y}{\sigma_x} = 0.3059 \cdot \frac{2.9475}{2.9475} = 0.5758$$

$$\text{Regression equation of } y \text{ on } x \Rightarrow y = b_{yx}(x-\bar{x})$$

$$= y - 66.25 = 0.5758(x - 65.25)$$

$$= y - 66.25 = 0.5758x - 37.5680$$

$$x = 64.5$$

$$= y - 66.25 = 37.1391 - 37.5680$$

$$y = 65.8211$$

## Analysis of Variance (ANOVA)

ANOVA is a statistical tool to test the homogeneity of

different groups based on their differences. ANOVA is

used to determine the diff b/w the means of samples

by analysing the variation within each of the samples.

ANOVA may be classified into two ways.

① one way classification (O.W.C)

② Two way classification (T.W.C)

① O.W.C is a shortcut method where a single factor is considered and its effect on samples is observed. This method is performed when the mean of the samples are equal or not, for eg. To study the effect of diff levels of fertilisers on a plant growth.

② To study the effect of diff levels of a medicine on a disease.

Consider a single factor and its effect on samples

or

$$x_1 \quad x_2 \quad x_3 \quad \dots \quad x_k$$

$$x_{11} \quad x_{21} \quad x_{31} \quad \dots \quad x_{k1}$$

$$x_{12} \quad x_{22} \quad x_{32} \quad \dots \quad x_{k2}$$

$$\vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots$$

$$x_{1n} \quad x_{2n} \quad x_{3n} \quad \dots \quad x_{kn}$$

To test there is only significant diff b/w samples or treatments means let us consider the null hypothesis as

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

$$\mu_1 \neq \mu_2 \neq \mu_3 \neq \dots \neq \mu_k$$

To test this hypothesis we can use F-ratio

test as follows

①  $\Sigma x_i^2$

② find grand total G

\* Correction factor for  $G^2 = \frac{G^2}{n}$

\* Cal total sum of squares TSS =  $\Sigma x_i^2 - CF$

\* Cal sum of squares for treatments as  $SST = \frac{\sum x_i^2}{n_i} - CF$

\* find the sum of squares of error or (SSE)

$$SSE = TSS - SST$$

\* Now cal the mean sum of squares for treatment as MST

$$MST = \frac{SST}{k-1}$$

\* Mean sum of square error is

$$MSSE = \frac{SSE}{n-k}$$

Here n = total no. of obs, k = no. of samples (or) treatments

$k-1$  is treatment degrees of freedom that is diff betw. no. of treatments  
 $n-k$  is error degrees of freedom  $\Rightarrow n - k = 15 - 3 = 12$

Now cal F-ratio or  $F = \frac{MSST}{MSSE}$

Now compare the cal value of  $F$  at 5% or 1% level of significance with  $(k-1, n-k)$  degrees of freedom if the table value of  $F$  is greater than the cal value of  $F$  then we accept our null hypothesis otherwise we reject the hypothesis.

Prob

1. Following are the data results of samples tested for medication, exercises and diet from 3 populations. Test the hypothesis that there is any diff b/w these 3 areas at 1% level of confidence

Medication	10	12	9	15	13	= 59
Exercise	16	8	3	6	2	= 31
diet	5	9	12	8	4	= 38

Sol  $k=3$ ;  $n=15$

$H_0$  = There is no significant diff among medication, exercise & diet

$H_1$  = There is a significant diff among medication, exercise & diet

$$\sum x_i^2 = 10^2 + 12^2 + 9^2 + 15^2 + 13^2 + 6^2 + 8^2 + 7^2 + 5^2 + 9^2 + 12^2 + 8^2 + 4^2 = 1162$$

$$G_t = 116$$

$$\text{Total CF} = \frac{G_t^2}{n} = \frac{116^2}{15} = 897.0667$$

$$\begin{aligned} TSS &= \sum x_{ij}^2 - \text{CF} \\ &= 1162 - 897.0667 \\ &= 264.9333 \end{aligned}$$

$$\begin{aligned} SST &= \frac{\sum x_i^2}{n} - \text{CF} \\ &= \frac{(59)^2}{5} + \frac{(19)^2}{5} + \frac{(38)^2}{5} - 897.0667 \\ &= 160.1333 \end{aligned}$$

$$SSE = TSS - SST$$

$$= 264.9333 - 160.1333$$

$$= 104.8$$

$$MSSE = \frac{SST}{k-1} = \frac{160.1333}{2} = 80.0667$$

$$MSSE = \frac{SSE}{n-k} = \frac{104.8}{12} = 8.7333$$

$$F = \frac{MSST}{MSSE} = \frac{80.0667}{8.7333} = 9.1680$$

The table value of F with  $(k-1, n-k) = (2, 12)$  at 1%.

level of significance is 6.93

Hence the table value of F is less than the cal value of F so we reject our null hypothesis.

② Two way classification: T.W.C technique is used when the data are classified on the basis of two factors.

for eg. The agricultural O.P may be classified on the basis of diff varieties of seeds and also on the basis of diff varieties of fertilisers. Let us consider the first factor has 'k' samples and second factor has 'n' samples are listed in a table below:

	$c_1$	$c_2$	$\dots$	$c_k$
$r_1$	$x_{11}$	$x_{12}$	$\dots$	$x_{1k}$
$r_2$	$x_{21}$	$x_{22}$	$\dots$	$x_{2k}$
$r_m$	$x_{m1}$	$x_{m2}$	$\dots$	$x_{mk}$

Let  $R$  = no. of rows

$C$  = no. of columns

$n$  = total no. of obs

To test the significance diff. of effect of two factors

Let us consider  $H_0$  as there is no significant diff.

of effect of first factors & there is no significant diff

of effect of second factor on the given samples

First find

(1)  $\sum x_{ij}^2$ , and then find grand total  $G$

(2) find the correction factor  $CF = \frac{G^2}{nR}$

- ③ find total sum of squares  $TSS = \sum x_i^2 - \bar{x}^2$   
 ④ find the row sum of squares  $SSR = \frac{\sum R_i^2}{n_r} = CR$   
 ⑤ find column sum of squares  $SSC = \frac{\sum C_j^2}{n_c} = CF$   
 ..  
 ⑥ find error sum of squares  $SSE = TSS - SSR - SSC$   
 ⑦ find mean sum of squares of the rows as  $MSSR = \frac{SSR}{R-1}$   
 ⑧ find " " " " " " columns as  $MSC = \frac{SSC}{C-1}$   
 ⑨ find " " " " " " error as  $MSE = \frac{SSE}{n-R-C+1}$   
 ⑩ Now cal F-Ratio for the first factor (row wise) is  $\frac{MSSR}{MSE}$   
 ⑪ cal F-ratio for the 2nd factor (column wise) is  $\frac{MSC}{MSE}$

If cal value of F is less than the table value of F with  $(R-1, n-R-C+1)$  degrees of freedom at 5% or 1% level of significance then we accept our first factor null hypothesis  
 If cal values of F is less than the table value of F with  $(C-1, n-R-C+1)$  degrees of freedom at 5% or 1% level of significance then we accept our 2nd factor null hypothesis

Ex: 1. Use 2-way classification test to the significance of students of  
Subs

	A	B	C	
BIO	157	250	156	$= 563.22$
French	148	160	158	$= 466$
English	150	150	145	$> 445$

Maths	140	160	151	161	161
Science	126	135	162	162	162
Computer	167	200	161	161	161
	<u>122</u>	<u>165</u>	<u>161</u>	<u>161</u>	<u>161</u>
	922	1055	937	937	937
					2924

6 rows represent subjects and column represent student

Rows = 6

Cols = 6

No. of total obs = 18

H<sub>0</sub>: There is no significant diff among the subs

H<sub>1</sub>: There is significant diff among the subs

H<sub>0</sub>: There is no significant diff among the subs

H<sub>1</sub>: There is significant diff among the subs

Also H<sub>0</sub> is true if  $\sum \sum \frac{X_{ij}^2}{n_i} - \frac{\sum X_{ij}^2}{n}$  is small

$$\text{Sum of } \sum \frac{X_{ij}^2}{n_i} = 487170$$

$$\text{Sum of } \sum X_{ij}^2 = 12182$$

$$\text{Sum of } \sum \frac{X_{ij}}{n_i} = 2924$$

$$CF = \frac{42}{18} = 2.33$$

$$TSS = \sum \sum X_{ij}^2 - CF = 12182.4444$$

$$SSTotal = \sum \sum \frac{X_{ij}^2}{n_i} - CF$$

$$= \frac{(163)^2}{3} + \frac{(146)^2}{3} + \frac{(165)^2}{3} + \frac{(161)^2}{3} + \frac{(161)^2}{3} - CF$$

$$= 4671.1111$$

$$SSC = \sum \frac{X_{ij}^2}{n_i} - CF$$

$$= \frac{(122)^2}{6} + \frac{(1055)^2}{6} + \frac{(937)^2}{6} - CF = 2088.7777$$

$$SSE = TSS - SSR - SSC = 5252.5556$$

$$MSSR = \frac{SSR}{R-1} = 974.2222$$

$$MSSC = \frac{SSC}{C-1} = 1029.3889$$

$$MSSE = \frac{SSE}{n-R-C+1} = 525.2556$$

$$F_R = \frac{MSSR}{MSSE} = 1.8568$$

$$F_C = \frac{MSSC}{MSSE} = 1.9598$$

The table value of  $F_R$  with  $(R-1, n-R-C+1) = (5, 10)$  at 5% level of significance is 3.35.

1. The cal value of  $F_R$  is less than tables of  $F_R$  so we accept our null hypothesis ①

The table value of  $F_C$  with  $(C-1, n-R-C+1) = (2, 10)$ , degrees of freedom at 5% level of significance is 4.10

2. The cal value of  $F_C$  is less than the table value of  $F_C$  so we accept our null hypothesis ②

The following data consists of breaking strength of 8 widgets were made from each of 4 materials and the resulting strength are shown below. Test the significant difference on the widgets of materials & breaking strength of widgets.

widgets	A	B	C	D	
1.	39	27	13	15	= 94
2.	57	43	32	28	= 160
3.	42	25	19	11	= 97
4.	32	11	11	12	= 66
5.	43	39	30	34	= 146
6.	50	49	30	34	= 163
7.	31	27	16	13	= 87
8.	51	34	28	23	= 136
	345	255	179	170	949

14 rows represents the widgets and columns represents

the material

$$L = \text{no of rows} = 8$$

$$n = \text{no of col} = 21$$

$$n = 32$$

(1)

~~H<sub>01</sub>~~ = There is no significant diff on the breaking strength of widgets

H<sub>02</sub> = There is no significant diff on the strength of material

H<sub>11</sub> = There is significant diff on the breaking strength of widgets

for hypothesis H<sub>01</sub> we have to calculate the mean of each row  
 for hypothesis H<sub>02</sub> we have to calculate the mean of each column  
 for hypothesis H<sub>11</sub> we have to calculate the mean of each row and column

$$\sum \Sigma x_{ij}^2 = 33399$$

$$G = 949$$

$$CF = \frac{G^2}{n} = 28143.7813$$

$$TSS = \sum \Sigma x_{ij}^2 - CF = 5255.2187$$

$$SSR = \left( \frac{R_1^2}{4} + \frac{R_2^2}{4} \dots \right) - CF = 2393.9687$$

$$SSC = \left( \frac{C_1^2}{8} + \frac{C_2^2}{8} \dots \right) - CF = 2480.0937$$

$$SSE = TSS - SSR - SSC = 381.1563$$

$$MSSR = \frac{SSR}{R-1} = \frac{2393.9687}{7} = 341.9955$$

$$MSSC = \frac{SSC}{C-1} = \frac{2480.0937}{6+3} = 826.6979$$

$$MSSE = \frac{SSE}{n-R-C+1} = \frac{381.1563}{8-6-1+1(2)} = 18.1503$$

$$F_R = \frac{MSSR}{MSSE} = 18.8424$$

$$F_C = \frac{MSSC}{MSSE} = 45.5473$$

The table value of  $F_R$  with  $(R-1, n-R-C+1) = (7, 2)$  is 2.49

at 5% is 2.49

The cal value of  $F$  is  $>$  table value of  $F$  so we reject our null hypothesis

The table value of  $F_c$  with  $(C-1, n-R-C+1)$  (3, 21) at 5% 3.07 the  
cal value of  $F$  is  $>$  than the table value of  $F$  so we  
reject our null hypothesis.