

Correlation:- Relationship between two variables such that a change in one variable results in a positive (or) negative change in the other variable.

There are three types of Correlation:-

1. positive Correlation
2. negative Correlation
3. perfect Correlation

1. positive Correlation:- If two variables deviate in the same direction the correlation is said to be positive (or) direct Correlation
 * One variable increase the other Variable also increase ↑ (or) ↓
 Example:- pressure Cooker.

2. Negative Correlation:- If the variables deviate in the opposite directions the correlation is said to be inverse (or) negative correlation.
 * If one variable increases the other decrease ↑ (or) ↓
 Example:- Boyle's law, exercise & Calories.

3. perfect Correlation:- If the change in one variable corresponds to a proportion change in the other variable, then the correlation is said perfect.

Karl Pearson's coefficient of correlation:- ⑤

The numerical measure of linear relationship between the random variables (rv) x, y is called coefficient of correlation between the variables and we write as R_{xy} (OR)

$$r = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \quad (\text{covariance})$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

where

$x = x - \bar{x}$, $y = y - \bar{y}$ (Observations are x and y + mean \bar{x}, \bar{y})

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$
$$\left[\frac{1}{n} \sum (x_i - \bar{x})^2 * \sum (y_i - \bar{y})^2 \right]^{1/2}$$

Note:- 1. 'r' lies between -1 and +1

2. $r=0$, variables are not correlated

(or)
Variables are Uncorrelated

3. $r=+1$, Variables positively perfectly correlated

4. $r=-1$, Variables negatively perfectly correlated.

Problem:- calculate the coefficient of correlation for the heights of fathers & their sons.

X	65	66	67	67	68	69	70	72
Y	67	68	65	68	72	72	69	71

X	Y	$X = x - \bar{x}$	$Y = y - \bar{y}$	X^2	Y^2	XY
65	67	65 - 68 = -3	67 - 69 = -2	9	4	6
66	68	-2	-1	4	1	2
67	65	-1	-4	1	16	4
67	68	-1	-1	1	1	-1
68	72	0	3	0	9	0
69	72	1	3	1	9	3
70	69	2	0	4	0	0
72	71	4	2	16	4	8

$$\frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{\sum xy}{\sqrt{44} \sqrt{44}} = \frac{\sum xy}{24} = 24$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

(3)

Problem:- Marks obtained by the students in maths and statistics by 10 students is given below. Find the correlation between two subjects

Marks in Maths	75	30	60	80	53	35	15	40	38	48
Marks in Statistics	85	45	54	91	58	63	35	43	45	44

X	Y	$X - \bar{X}$	$Y = Y - \bar{Y}$	X^2	Y^2	XY		$\sum X^2$	$\sum Y^2$	$\sum XY$
75	85	75 - 47.4 = 27.6	85 - 56.3 = 28.7	761.76	823.69	792.12				
30	45	-17.4	-11.3	302.76	127.69	196.62				
60	54	12.6	-2.3	158.76	5.29	-28.98				
80	91	32.6	34.7	1062.76	1204.09	1131.22				
53	58	5.6	1.7	31.36	2.89	9.52				
35	63	-12.4	6.7	153.76	44.89	-83.08				
15	35	-32.4	-21.3	104.976	453.69	690.12				
40	43	-7.4	-13.3	54.76	176.89	98.42				
38	45	-9.4	-11.3	88.36	127.69	106.22				
48	44	0.6	-12.3	0.36	151.29	-7.38				
				$\sum X^2 = 474$	$\sum Y^2 = 563$	$\sum XY = 3118.3$				
				$\bar{x} = 47.4$	$\bar{y} = 56.3$	$\bar{xy} = 56.3$				

$$r = \frac{\sqrt{(\sum xy) * (\sum xy)}}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

$$= \frac{\sqrt{(3118.3) * (3118.3)}}{\sqrt{474} \sqrt{563}}$$

$$= \frac{2904.8}{2904.8}$$

- * 5th decimal place ≥ 5 add 1 to previous value (5)
 * 5th decimal place < 5 the previous number is same

Correlation coefficient - another method

Direct method:-

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

where n = no of observations

Problem:- Find 'r' for the following data by direct method.

x	65	66	67	67	68	69	70	72
y	67	68	65	68	72	72	69	71

x	y	xy	x^2	y^2
65	67	4355	4225	4489
66	68	4488	4356	4624
67	65	4355	4489	4225
67	68	4556	4489	4624
68	72	4896	4624	5184
69	72	4968	4761	5184
69	69	4830	4900	4761
70		5184	5184	5041
72	71	5112	5184	5041
$\sum x = 544$	$\sum y = 552$	$\sum xy = 37560$	$\sum x^2 = 37028$	$\sum y^2 = 38132$

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

$$r = \frac{18 * (37560) - (544)(552)}{\sqrt{8 * (37028) - (544)^2} * \sqrt{8 * (38262) - (552)^2}}$$

$$= \frac{300480 - 300288}{\sqrt{296224 - 295936} * \sqrt{306096 - 304704}}$$

$$= \frac{192}{\sqrt{288} * \sqrt{1392352}}$$

$$= \frac{192}{\sqrt{400896}}$$

101376

$$= \frac{192}{\cancel{633.1635}} = 318 / 318.3960$$

$r = 0.3034$	$r = 0.6030$
--------------	--------------

② Problem:- Find the Correlation coefficient for the following data by direct method

X	20	30	40	80	60	70
Y	23	18	30	20	15	8

x	y	xy	x^2	y^2	(7)
20	23	460	400	529	
30	18	540	900	324	
40	30	1200	1600	900	
80	20	1600	6400	400	
60	15	900	3600	225	
70	8	560	4900	64	
$\sum x = 300$	$\sum y = 114$	$\sum xy = 5260$	$\sum x^2 = 17800$	$\sum y^2 = 2442$	

$$\begin{aligned}
 r &= \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} * \sqrt{n \sum y^2 - (\sum y)^2}} \\
 &= \frac{6 * (5260) - (300)(114)}{\sqrt{6 * 17800 - (300)^2} \sqrt{6 * 2442 - (114)^2}} \\
 &= \frac{31560 - 34200}{\sqrt{106800 - 90000} * \sqrt{14652 - 12996}} \\
 &= \frac{-2640}{\sqrt{16800} \sqrt{1656}} \\
 &= \frac{-2640}{10 \sqrt{168} * \sqrt{1656}} \\
 &= \frac{-264}{\sqrt{278208}}
 \end{aligned}$$

$$r = 0.5005$$

Regression:- Regression means stepping back (8)
towards the normal

There are two types of Regression lines

1. Regression line y on x
2. Regression line x on y

1. Regression line y on x :-

$$y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

where $r \frac{\sigma_y}{\sigma_x} = b_{yx}$ is the regression coefficient of y on x

σ_x = standard deviation of x

σ_y = standard deviation of y

\bar{x} = mean of x

\bar{y} = mean of y

2. Regression line x on y :-

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

where $r \frac{\sigma_x}{\sigma_y} = b_{xy}$ is the regression coefficient of x on y

\bar{x} = mean of x

\bar{y} = mean of y

σ_x = SD of x

σ_y = SD of y

Note:- $r^2 = b_{yx} * b_{xy} (01)$

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

Problem:- Calculate the coefficient of correlation from the following data. Also obtain the equation of the lines of regression & obtain an estimate of y which should correspond on the average of $x = 6.2$ ⑨

x	1	2	3	4	5	6	7	8	9
y	9	8	10	12	11	13	14	16	15

Here $n = 9$

x	y	$x = x - \bar{x}$	$y = y - \bar{y}$	x^2	y^2	xy
1	9	-4	-3	16	9	12
2	8	-3	-4	9	16	12
3	10	-2	-2	4	4	4
4	12	-1	0	1	0	6
5	11	0	-1	0	1	0
6	13	1	2	1	1	1
7	14	2	2	4	4	4
8	16	3	4	9	16	12
9	15	4	3	16	9	12
$\sum x = 45$	$\sum y = 108$			$\sum x^2 = 60$	$\sum y^2 = 60$	$\sum xy = 57$

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}} = \frac{57}{\sqrt{60 \times 60}} = \frac{57}{60} = 0.95$$

(10)

Re
vva

$$\sigma_x = \sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2}$$

$$= \sqrt{\frac{1}{n} \sum x^2}$$

$$= \sqrt{\frac{1}{n} * 60}$$

$$= \sqrt{\frac{1}{3} * 66.20}$$

$$= \sqrt{\frac{20}{3}}$$

$$\boxed{\sigma_x = 2.5819}$$

6.667

$$\sigma_y = \sqrt{\frac{1}{n} \sum (y_i - \bar{y})^2}$$

$$= \sqrt{\frac{1}{3} \sum y^2}$$

$$= \frac{1}{3} \sqrt{60}$$

$$= \frac{1}{3} * 7.7450$$

$$\boxed{\sigma_y = 2.5819}$$

Regression line y on x :-

(ii)

$$y - \bar{y} = 2 \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$y - 12 = 0.95 \left(\frac{2.5819}{2.5819} \right) (x - 5)$$

$$y - 12 = 0.95 (x - 5)$$

$$y - 12 = 0.95x - 4.75$$

$$\boxed{y = 0.95x + 7.25}$$

Regression line x on y :-

$$x - \bar{x} = 2 \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$x - 5 = 0.95 * \left(\frac{2.5819}{2.5819} \right) (y - 12)$$

$$x - 5 = 0.95 (y - 12)$$

$$x - 5 = 0.95y - 11.4$$

$$\boxed{x = 0.95y - 6.4}$$

Problem:- The following data relate to the marks of 10 students in the internal test and the university examination for the maximum of 50 each.

Internal Exam	25	28	30	32	35	36	38	39	42	45
University marks	20	26	29	30	25	18	26	35	35	46

1. The most likely internal mark for university
mark of 25 (12)
2. The most likely university mark for internal
mark of 30.

x	y	$x = x - \bar{x}$	$y = y - \bar{y}$	x^2	y^2	xy
25	20	-10	-9	100	81	90
28	26	-7	-3	49	9	21
30	29	-5	0	25	0	0
32	30	-3	1	9	1	-3
35	35	0	-4	0	16	0
36	18	-1	-11	1	121	-11
38	26	3	-3	9	9	-9
39	35	4	6	16	36	24
42	35	7	6	49	36	42
45	46	10	17	100	289	170
$\sum x = 350$	$\sum y = 290$			$\sum x^2 = 358$	$\sum y^2 = 598$	$\sum xy = 27432$
$\frac{350}{10} = 35$	$\frac{290}{10} = 29$					

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

$$= \frac{324}{\sqrt{358 * 598}}$$

$$= \frac{324}{\sqrt{190814214084}}$$

$$r = 0.7002$$

$$\sigma_x = \sqrt{\frac{\sum x^2}{n}} \quad \sigma_y = \sqrt{\frac{\sum y^2}{n}}$$

$$= \sqrt{\frac{358}{10}} \quad = \sqrt{\frac{598}{10}}$$

$$= \sqrt{35.8} \quad \sigma_y = 7.7330$$

$$\sigma_x = 5.9833$$

(13)

Regression y on x :

$$y - \bar{y} = 2 \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$y - 29 = 0.7002 \left(\frac{7.7330}{5.9833} \right) (x - 35)$$

$$y - 29 = 0.7002 (1.2924) (x - 35)$$

$$y - 29 = 0.9050 (x - 35)$$

$$y - 29 = 0.9050 x - 31.6736$$

$$y = 0.9050 x - 2.6736$$

Regression x on y :

$$x - \bar{x} = 2 \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$x - 35 = 0.7002 * \left(\frac{5.9833}{7.7330} \right) (y - 29)$$

$$x - 35 = 0.7002 * (0.7737) (y - 29)$$

$$x - 35 = 0.5418 (y - 29)$$

$$x - 35 = 0.5418 y - 15.712$$

$$x = 0.5418 y + \frac{(17.2878)}{19.2878}$$

(19)

Estimation:-

$\tilde{x} \rightarrow$ internal mark
 $y \rightarrow$ university mark

1. x when $y=25$ use x on y Regression Line

$$x = 0.5418(25) + 19.2878$$

$$x = 13.545 + 19.2878$$

$$\boxed{x = 32.8328}$$

2. y when $x=30$ use y on x Regression line

$$y = 0.9050(30) - 2.6736$$

$$= 27.15 - 2.6736$$

$$\boxed{y = 24.4764}$$

Problem:- 10 observations on price (x) and supply (y)
 \tilde{x} were obtained below

$$\sum x = 130, \sum x^2 = 2288, \sum xy = 3467$$

$$\sum y = 220, \sum y^2 = 5506$$

obtained the line of regression of y on x &
estimate the supply when the price is
16 units & also find the standard error of
the estimate.

$$\boxed{r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}}$$

Here, $n=10$

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

$$\sqrt{n \sum x^2 - (\sum x)^2} \quad \sqrt{n \sum y^2 - (\sum y)^2}$$

$$\lambda = \frac{10(3467) - (130)(220)}{\sqrt{10(2288) - (130)^2} \sqrt{10(5506) - (220)^2}}$$

$$= \frac{34670 - 28600}{\sqrt{22880 - 16900} \sqrt{55060 - 48400}}$$

$$= \frac{6070}{\sqrt{5980} \sqrt{6660}}$$

$$= \frac{6070}{\sqrt{39826800}}$$

$$= \frac{6070}{6310.8478}$$

$$\boxed{\lambda = 0.9618}$$

$$\bar{x} = \frac{\sum x}{n} = \frac{130}{10} = 13, \quad \bar{y} = \frac{\sum y}{n} = \frac{220}{10} = 13$$

$$\sigma_x = \sqrt{\frac{1}{n} \sum (x - \bar{x})^2}$$

$$\sigma_x = \sqrt{\frac{1}{n} \left[\sum x^2 - \left(\frac{\sum x}{n} \right)^2 \right]}$$

$$= \frac{1}{n} \sum (x - \bar{x})^2$$

$$= \frac{1}{n} \left[\sum (x)^2 + \sum (\bar{x})^2 - 2 \sum x \bar{x} \right]$$

$$= \frac{1}{n} \left[\sum x^2 + n(\bar{x})^2 - 2(n)(\bar{x}) \right]$$

$$= \frac{1}{n} \left[\sum x^2 - 2n\bar{x}^2 + n(\bar{x})^2 \right]$$

(15)

(16)

$$= \frac{1}{n} \left[\sum x^2 - n(\bar{x})^2 \right]$$

$$= \frac{1}{n} \left[\sum x^2 - n\left(\frac{\sum x}{n}\right)^2 \right]$$

$$\sigma_x = \frac{1}{n} \left[\sum x^2 - (\bar{x})^2 \right]$$

$$\sigma_x = \sqrt{\frac{1}{10} \left[2288 - \left[\frac{130}{10} \right]^2 \right]}$$

$$= \sqrt{\frac{1}{10} \left[2288 - \left[\frac{1690}{100} \right] \right]}$$

$$= \sqrt{\frac{1}{10} \left[[2288] - 1690 \right]}$$

$$= \sqrt{\frac{1}{10} \left[2119 \right]} = \sqrt{59.8}$$

$$\sigma_x = \sqrt{59.8}$$

 ~~σ_x~~

$$\sigma_x = 7.7330$$

$$\sigma_y = \sqrt{\frac{1}{n} \left[\sum y^2 - \frac{(\sum y)^2}{n} \right]}$$

$$= \sqrt{\frac{1}{n} \left[5506 - \frac{(220)^2}{10} \right]}$$

$$= \sqrt{\frac{1}{10} \left[5506 - \frac{48400}{10} \right]}$$

$$= \sqrt{\frac{1}{10} \left[5506 - 4840 \right]}$$

$$= \sqrt{666/10} \quad \boxed{\sigma_y = 8.1608}$$

Estimation:-

(7)

y when $x = 16$

$$y = 1.0150 * 16 + 8.8049$$

$$y = 16.24 + 8.8049$$

$$y = 25.0449$$

Standard Error:-

$$\boxed{S.E = \frac{1 - \alpha^2}{\sqrt{n}}}$$

$$\begin{aligned}
 S.E &= \frac{1 - (0.9618)^2}{\sqrt{10}} \\
 &= \frac{1 - 0.92505}{\sqrt{10}} \\
 &= \frac{1 - 0.9250}{\sqrt{10}} \\
 &= 1 - 0.2925
 \end{aligned}$$

$$\boxed{S.E = 0.0237}$$

~~***~~

Problem:- The eq'n of two regression lines obtained in a correlation analysis

$$3x + 12y = 19, 3y + 9x = 46$$

a) Find coefficient of correlation

b) Mean values of x & y

c) Ratio of coefficients of variability.

(18)

Ques:- Let

$$3x+12y = 19 \quad \text{--- (1)}$$

$$3y+9x = 46 \quad \text{--- (2)}$$

nothing is specified as y on x or x on y
 than consider (1) eqn y on x and (2) eqn
 x on y

$$\text{i.e., } 3x+12y = 19$$

$$12y = -3x + 19$$

$$y = \frac{-3x + 19}{12} \quad (\because y = mx + c)$$

$$y = -\frac{x}{4} + \frac{19}{12} \quad \text{eqn of line}$$

$m \rightarrow$ slope is an regression coefficient

$$b_{yx} = -\frac{3}{12} = -\frac{1}{4}$$

Let suppose eqn (2) to be x on y regression
 Line

$$\text{i.e., } 3y+9x = 46$$

$$9x = -3y + 46$$

$$x = \frac{-3y}{9} + \frac{46}{9}$$

$$x = -\frac{y}{3} + \frac{46}{3}$$

$b_{xy} = -\frac{1}{3}$ is the regression coefficient

$$\rho = \sqrt{b_{yx} * b_{xy}}$$

$$= \pm \sqrt{(-\frac{1}{4})(-\frac{1}{3})}$$

$$\rho = -0.2886$$

when ever b_{yx} and b_{xy} are -ve 'x' is -ive.
 If b_{yx} & b_{xy} are +ve than 'x' is +ve
 Range of 'x' is -1 to +1.

b) Mean Value:-

Mean value of $x + y$

solving eqn ① + ② we get

$$3x + 12y = 19 \quad \text{--- } ①$$

$$9x + 3y = 48 \quad \text{--- } ②$$

$$\begin{array}{r} 9x + 36y = 57 \\ - 9x - 3y = 46 \\ \hline 33y = 11 \end{array}$$

$$y = 11/33$$

$$y = 1/3$$

$$\begin{array}{r} 3x + 12y = 19 \\ - 36x - 12y = - 184 \\ \hline - 33x = - 165 \end{array}$$

$$x = \frac{-165}{33}$$

$$x = 5$$

Intersection of point of two regression lines
 is the mean value of $x + y$.

$$(5, 1/3) = (\bar{x}, \bar{y}) = (5, 0.3333)$$

3. coefficient of variability:-

$$CV = \frac{\sigma_x^2}{\bar{y}^2}$$

$$\text{we have } b_{xy} = \frac{\sigma_x^2}{\sigma_y^2}$$

$$(b_{xy})^2 = \frac{\sigma_x^2}{\sigma_y^2}$$

$$\frac{\sigma_x^2}{\sigma_y^2} = \left(\frac{b_{xy}}{2}\right)^2 = \frac{(-1/3)^2}{1/12} = \frac{1}{9} \times 12 = \frac{4}{3}$$

$$\begin{cases} \bar{x} = 4/3 \\ \bar{y} = 1.3333 \\ \text{Coff of variability} \\ = 1.3333 \end{cases}$$

19/10/2020
Monday

(20)

** Rank Correlation:- Suppose x_1, x_2, \dots, x_n and y_1, y_2, \dots, y_n be the ranks of two variables x & y in the order of merit with respect to some property. The correlation b/w these in pairs of ranks is called rank correlation.

Rank Correlation is measured by Spearman's rank correlation coefficient.

$$P = 1 - \left[\frac{6 * \left(\sum_{i=1}^n d_i^2 \right)}{n(n^2-1)} \right]$$

where n = no. of observation

x_i = Rank of X

y_i = Rank of Y

$d_i = x_i - y_i$

Note:- The rank correlation coefficient (P) lies b/w $-1 \neq 1$

Problem:- obtain the rank correlation coefficient for the following data

(Q1)

X	10	15	12	17	13	16	24	14	22	20
Y	30	42	45	46	33	34	40	35	39	38

X	Y	Rank of X $= x_i$	Rank of Y $= y_i$	$d_i = x_i - y_i$	d_i^2
10	30	10	10	0	0
15	42	6	3	3	9
12	45	9	2	7	49
17	46	4	1	3	9
13	33	8	9	-1	1
16	34	5	8	-3	9
24	40	1	4	-3	9
14	35	7	7	0	0
22	39	2	5	-3	9
20	38	3	6	-3	9

$$\sum x = 163 \quad \sum y = 382 \quad \sum x_i = 55 \quad \sum y_i = 55 \quad \sum d_i^2 = 104$$

$$P = \frac{1 + \beta * \sqrt{\sum d_i^2}}{n}$$

$$P = 1 - \left[\frac{6 * (\sum d_i^2)}{n(n^2 - 1)} \right]$$

$$= 1 - \left[\frac{6 * 104}{10(10^2 - 1)} \right]$$

$$= 1 - \left[\frac{6 * 104}{10(99)} \right]$$

$$P = \frac{624}{990}$$

$$= 1 - 6303$$

$$P = 0.3697$$

Holds good for Unique ranks.

Problem:- Calculate the rank correlation coefficient for the following data.

X	124	100	105	112	102	93	99	115	123	104	92
	113	121	103	101							
Y	80	100	102	91	92	111	109	98	89	104	115
	105	97	106	123							

X	Y	Rank of X = x_i	Rank of Y = y_i	$x_i - y_i$	d_i^2
124	80	1	15	-14	196
100	100	12	9	3	9
105	102	7	8	-1	1
112	91	6	13	-7	49
102	92	10	12	-2	4
93	111	15	3	12	144
99	109	13	4	9	81
115	98	4	10	-6	36
123	89	2	14	-12	144
104	104	8	7	1	1
98	113	14	2	12	144
113	105	5	6	-1	1
121	97	3	11	-8	64
103	106	9	5	4	16
101	123	11	1	10	100

$$\sum d_i^2 = 990$$

(23)

$$P = 1 - \left[\frac{6 * \sum d_i^2}{n(n^2-1)} \right]$$

$$= 1 - \left[\frac{6 * 990}{15(15^2-1)} \right]$$

$$= 1 - \left[\frac{6 * 990}{15(224)} \right]$$

$$= 1 - \left[\frac{5940}{3360} \right]$$

$$= 1 - [1.7679]$$

$$P = 0.7679$$

Repetition of ranks :-

$$P = 1 - \left[\frac{6 * [\sum d_i^2 + C \cdot F]}{n(n^2-1)} \right]$$

CF = Correlation Factor

CF = CF in X series + CF in Y series

Formula for

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

where m is the no. of times the rank
is repeated

Problem:- Calculate the rank correlation coefficient for the following data.

(24)

X	Y	x_i	y_i	$x_i - y_i$	d_i^2
68	62	4	5	-1	1
64	58	6	7	-1	1
75	68	2.5	3.5	-1	1
50	45	9	10	-1	1
64	81	6	1	5	1
80	60	1	6	-5	25
75	68	2.5	3.5	-1	25
40	48	10	9	1	1
55	50	8	8	0	1
64	70	6	2	4	0
					16

$$\sum d_i^2 = 72$$

$$P = 1 - \left[\frac{6 * [\sum d_i^2 + CF]}{n(n^2 - 1)} \right]$$

CF = CF in X series + CF in Y series

(CF in X series) - 75 is repeated 2 times
 64 is repeated 3 times

$$\begin{aligned}
 CF &= \frac{(2^2 - 1)}{12 \cdot 6} + \frac{(3^2 - 1)}{12 \cdot 4} \\
 &= \frac{4 - 1}{12 \cdot 6} + \frac{8 - 1}{12 \cdot 4} \\
 &= \frac{3}{12} + \frac{7}{12} \\
 CF &= 1/2 + 2 \\
 &= 2.5
 \end{aligned}$$

C
2

CF in Y series:-

(25)

$$CF \text{ in } Y \text{ series} = \frac{\alpha(2^d - 1)}{2}$$
$$= 0.5$$

$CF = CF \text{ in } X \text{ series} + CF \text{ in } Y \text{ series}$

$$= 2.5 + 0.5$$

$$CF = 3.0$$

$$\rho = 1 - \left[\frac{6 * (72+3)}{10(10^2-1)} \right]$$
$$= 1 - \left[\frac{6 * (75)}{10 * 99} \right]$$
$$= 1 - \left[\frac{450}{990} \right]$$
$$= 1 - \left[\frac{45}{99} \right]$$
$$= 1 - 0.4545$$

$P = 0.5455$

Problem:- A sample of 12 fathers and their sons gave the following data about their heights. Calculate their rank correlation.

x	y	x_i	y_i	$x_i - y_i$	$(x_i - y_i)^2$
65	68	9	5.5	3.5	12.25
63	66	11	9.5	1.5	2.25
67	68	6.5	5.5	1	1
64	65	10	7.5	-1.5	2.25
68	69	4.5	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	11.5	-3.5	12.25
68	71	4.5	1	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
					$\sum d_i^2 = 72.5$

CF in X series:- 67 - 2 times
 ~~~~~ 68 - 2 times

$$CF_x = \left( \frac{2(2^2-1)}{12} \right) + 2 \left( \frac{2^2-1}{12} \right)$$

$$= 0.5 + 0.5$$

$$= 1$$

CF in Y series:- 68 - 4 times

65 - 2 times

66 - 2 times

$$CF_y = \frac{4(4^2-1)}{12} + \frac{2(2^2-1)}{12} + \frac{2(2^2-1)}{12}$$

$$= \frac{15.5}{8} + \frac{2}{8} + \frac{2}{8}$$

$$CF_y = 5 + 0.5 + 0.5$$

(27)

$$CF_y = 6$$

$$CF = CF_x + CF_y$$

$$= 1 + 6$$

$$\boxed{CF = 7}$$

$$\rho = 1 - \left[ \frac{6(-72.5 + 7)}{12(144 - 1)} \right]$$

$$= 1 - \left[ \frac{6(-65.5)}{12(143)} \right]$$

$$= 1 - \left[ \frac{\frac{2}{3}(-65.5)}{286} \right]$$

$$\left. \begin{aligned} &= 1 - [0.1512] \\ &\rho = 0.8487 \\ &\rho = 1 - 0.27189 \end{aligned} \right\}$$

$$\rho = \left[ 1 - \frac{79.5}{286} \right]$$

$$= 1 - 0.27197$$

$$\boxed{\rho = 0.7221}$$

Limits of Correlation:-

(28)

$$r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$$r = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$\text{Let } a_i = x_i - \bar{x}$$

$$b_i = y_i - \bar{y}$$

$$r^2_{x,y} = \left( \frac{1}{n} \right)^2 \left[ \sum_{i=1}^n (a_i)(b_i) \right]^2$$
$$= \frac{1}{n^2} \frac{\left( \sum_{i=1}^n a_i b_i \right)^2}{\left( \frac{1}{n} \sum_{i=1}^n a_i^2 \right) \left( \frac{1}{n} \sum_{i=1}^n b_i^2 \right)}$$

$$r^2_{x,y} = \frac{\sum_{i=1}^n (a_i b_i)^2}{\sum a_i^2 \sum b_i^2} \rightarrow ①$$

We have Cauchy-Schwarz inequality which states that if  $a_i b_i, i = 1, 2, \dots, n$  are real quantities then

$$(\sum a_i b_i)^2 \leq (\sum a_i^2)(\sum b_i^2)$$

The sign of equality holding if

$$\frac{a_1}{b_1} = \frac{a_2}{b_2} = \dots = \frac{a_n}{b_n}$$

(29)

Using Cauchy Schwartz inequality in ① we have  $\sigma^2(x, y) \leq 1$

$$\Rightarrow |\rho(x, y)| \leq 1$$

$$-1 \leq \rho \leq 1$$

Properties of Correlation Coefficient:-

1. Limits for Correlation coefficient are  $(-1, +1)$
2. Correlation coefficient is independent of change of origin & scale.
3. Two independent variables are Uncorrelated

Properties of Regression Coefficient:-

1. Correlation coefficient is the geometric mean between the regression coefficients

Proof:- We have  $b_{yx} * b_{xy} = \frac{\partial f}{\partial x} * \frac{\partial f}{\partial y}$   
 $= \rho^2$

$$b_{yx} * b_{xy} = \rho^2 = \rho^2$$

$$\rho = \pm \sqrt{b_{yx} * b_{xy}}$$

Note:- Sign of  $\rho$  is same as that of regression coefficients

( $\because$  the sign of each regression coefficient depends on the covariance term)

If both the regression coefficients are +ve  
 'r' is +ve  
 If both the regression coefficients are -ve  
 'r' is -ve.

2. If one of the regression coefficients is greater than Unity the other is less than Unity

Proof:- let  $b_{yx}$  is greater than Unity we have to show that  $b_{xy}$  is less than Unity

$$b_{xy} > 1$$

$$\frac{1}{b_{xy}} < 1 \quad \textcircled{1}$$

$$\text{Also } r^2 \leq 1 \Rightarrow b_{yx} b_{xy} \leq 1$$

$$\Rightarrow b_{xy} \leq \frac{1}{b_{yx}}$$

$$\Rightarrow b_{xy} < 1 \quad (\text{from eqn 1})$$

$\therefore$  If  $b_{yx} > 1$  then  $b_{xy} < 1$  & Vice Versa.

3. The modulus of the arithmetic mean of the regression coefficient is not less than the modulus value of the correlation coefficient 'r'.

$$\text{i.e., } \left| \frac{b_{yx} + b_{xy}}{2} \right| > |r|$$

Proof:- To prove that

$$\left| \frac{1}{2} (b_{yx} + b_{xy}) \right| > |r|$$

$$\Rightarrow \left| \frac{1}{2} \left( \frac{\partial \bar{y}}{\partial x} + \frac{\partial \bar{x}}{\partial y} \right) \right| > |\gamma| \quad (31)$$

$$\Rightarrow \left| \gamma \left( \frac{\partial \bar{y}}{\partial x} \right) + \gamma \left( \frac{\partial \bar{x}}{\partial y} \right) \right| > 2|\gamma|$$

$$\Rightarrow \left| \frac{\partial \bar{y}}{\partial x} + \frac{\partial \bar{x}}{\partial y} \right| > 2 \quad (\because |\gamma| > 0)$$

$$\Rightarrow \frac{\sigma_y^2 + \sigma_x^2}{\sigma_x \sigma_y} > 2 \quad (\because \sigma_x + \sigma_y > 0 \\ \text{modulus removed})$$

$$\Rightarrow \sigma_y^2 + \sigma_x^2 > 2\sigma_x \sigma_y$$

$$\Rightarrow \sigma_y^2 + \sigma_x^2 - 2\sigma_x \sigma_y > 0$$

$$\Rightarrow (\sigma_y - \sigma_x)^2 > 0$$

which is always true since the square of a real quantity is  $> 0$ .

4. The regression Coefficients are independent of change of Origin but not scale.

Proof :- let  $U = \frac{x-a}{h}$  &  $V = \frac{y-b}{k}$  where  $a, b, h, k$  are constants.

$$\Rightarrow x = a + uh ; y = b + vk$$

$$\begin{aligned} \text{cov}(x, y) &= E \left[ \{x - E(x)\} \{y - E(y)\} \right] \\ &= E \left[ \{(a+uh) - E(a+uh)\} \right] \\ &\quad * \left[ \{(b+vk) - E(b+vk)\} \right] \end{aligned}$$

$$\begin{aligned}\text{cov}(X, Y) &= E \left[ h \{U - E(U)\} \{k \{V - E(V)\}\} \right] \\ &= hk E \left[ \{U - E(U)\} \{V - E(V)\} \right]\end{aligned}$$

$$\text{cov}(X, Y) = hk \text{cov}(U, V) \quad \textcircled{1}$$

$$\begin{aligned}\text{Var}(X) &= \sigma_x^2 = V(X) = \text{Var}(a + uh) \\ &= h^2 V(u) \rightarrow \textcircled{2}\end{aligned}$$

Similarly

$$\begin{aligned}\sigma_y^2 &= V(Y) = \text{Var}(b + kv) \\ &= k^2 \text{Var}(v)\end{aligned}$$

$$\text{we have } b_{yx} = \frac{\text{cov}(X, Y)}{\sigma_x^2}$$

$$r = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y}$$

$$b_{yx} = r \frac{\sigma_y}{\sigma_x}$$

$$b_{yx} = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} * \frac{\sigma_y}{\sigma_x}$$

|                                                |
|------------------------------------------------|
| $b_{yx} = \frac{\text{cov}(X, Y)}{\sigma_x^2}$ |
|------------------------------------------------|

$$b_{yx} = \frac{\text{cov}(X, Y)}{\sigma_x^2} = \frac{hk \text{cov}(U, V)}{h^2 V(u)} \quad (\text{from } \textcircled{1} + \textcircled{2})$$

$$= \frac{k}{h} \frac{\text{cov}(U, V)}{\text{Var}(U)}$$

|                               |
|-------------------------------|
| $b_{yx} = \frac{k}{h} b_{vu}$ |
|-------------------------------|

Similarly we can prove that  $b_{xy} = \frac{h}{k} b_{uv}$  (33)

∴ Regression coefficient are independent of change of origin but not scale

Angle between two regression lines:-

$$\text{Let } y - \bar{y} = \alpha \frac{\partial y}{\partial x} (x - \bar{x}) \rightarrow ①$$

be the regression line  $y$  on  $x$

$$\text{Let } x - \bar{x} = \alpha \frac{\partial x}{\partial y} (y - \bar{y}) \rightarrow ② \text{ be the regression line } x \text{ on } y.$$

$$② \Rightarrow y - \bar{y} = \frac{\partial y}{\partial x} (x - \bar{x})$$

Angle between the lines of regression:-

$$\begin{aligned} \tan \theta &= \left| \frac{\frac{\partial y}{\partial x} - \frac{\partial y}{\partial x}}{1 + \left( \frac{\partial y}{\partial x} \right) \left( \frac{\partial y}{\partial x} \right)} \right| \\ &= \left| \frac{\alpha^2 \partial y - \partial y}{\alpha \partial x} * \frac{\partial x^2}{\partial x^2 + \partial y^2} \right| \\ &= \left| \frac{\partial y (\alpha^2 - 1)}{\alpha \partial x} * \frac{\partial x^2}{\partial x^2 + \partial y^2} \right| \end{aligned}$$

$$\tan \theta = \frac{\partial x \partial y (\alpha^2 - 1)}{\alpha (\partial x^2 + \partial y^2)}$$

$$\tan \theta = \left| \frac{(r^2 - 1)}{r} \right| \left| \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right|$$

(34)

$$\tan \theta = \left| \frac{1 - r^2}{r} \right| \left( \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \quad (r^2 \leq 1)$$

$$\boxed{\theta = \tan^{-1} \left\{ \left| \frac{1 - r^2}{r} \right| \left( \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \right\}}$$

Case 1 :- If  $r = 0$ ,  $\tan \theta = \infty \Rightarrow \theta = \pi/2$

Thus if the two variables are uncorrelated, the lines of regression becomes  $\perp$  to each other.

Case 2 :- If  $r = \pm 1$ ,  $\tan \theta = 0 \Rightarrow \theta = 0$  or  $\pi$

In this case, the two lines of regression either coincide (or) they are parallel to each other. Since lines of regression pass through  $(\bar{x}, \bar{y})$ , they can't be parallel.

So, in case of perfect correlation +ve (or) -ve, the two lines of regression coincide.

- Problem :- The data given below relates to marks in two subjects mathematics and statistics of B.Tech students. The correlation coefficient b/w the marks in two subjects is 0.42
1. Estimate the marks in statistics if the marks in maths is 52
  2. Find the angle b/w two regression lines

|                        | Maths x | statistics y |
|------------------------|---------|--------------|
| Average marks<br>means | 39.5    | 49.5         |
| standard deviation     | 10.8    | 16.8         |

35

The Regression line y on x:-

$$y - \bar{y} = \frac{\sigma_{xy}}{\sigma_x} (x - \bar{x})$$

$$y - 49.5 = (0.42) \left( \frac{16.8}{10.8} \right) (x - 39.5)$$

$$y - 49.5 = (0.42)(1.5556) (x - 39.5)$$

$$y - 49.5 = 0.6533 (x - 39.5)$$

$$y - 49.5 = 0.6533x - 25.8074$$

$$y = 0.6533x - 25.8074 + 49.5$$

$$y = 0.6533x + 23.6933$$

when  $x = 52$

$$y = 0.6533 * 52 + 23.6933$$

$$y = 33.9716 + 23.6933$$

$$y = 57.6649$$

$$\begin{aligned} \theta &= \tan^{-1} \left\{ \left( \frac{1 - r^2}{1 + r^2} \right) \left( \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \right\} \\ &= \tan^{-1} \left\{ \left( \frac{1 - (0.42)^2}{0.42} \right) \left( \frac{10.8 * 16.8}{(10.8)^2 + (16.8)^2} \right) \right\} \\ &= \tan^{-1} \left\{ \left( \frac{1 - 0.1764}{0.42} \right) \left( \frac{10.8 * 16.8}{116.64 + 282.24} \right) \right\} \\ &= \tan^{-1} \left\{ \left( \frac{0.8236}{0.42} \right) \left( \frac{181.44}{398.88} \right) \right\} \end{aligned}$$

$$\theta = \tan^{-1} \left( 1.9610 * 0.4549 \right)$$

(36)

$$= \tan^{-1} (0.8920)$$

$$\boxed{\theta = 41.7330^\circ}$$

Problem:- calculate the Correlation coefficient for the following data

|   |    |    |    |    |    |    |    |    |
|---|----|----|----|----|----|----|----|----|
| x | 65 | 66 | 67 | 67 | 68 | 69 | 70 | 72 |
| y | 67 | 68 | 65 | 68 | 72 | 72 | 69 | 71 |

$$x - x \quad U = x - 68$$

$$v = y - 72$$

| $U-68$ | $v-72$ | $U^1 = U - U'$ | $V^1 = V - V'$ | $(U^1)^2$       | $(V^1)^2$       | $U^1 V^1$ | $U'$ |
|--------|--------|----------------|----------------|-----------------|-----------------|-----------|------|
| -3     | -5     | -3             | -2             | 9               | 4               | 6         |      |
| -2     | -4     | -2             | -1             | 4               | 1               | 2         |      |
| -1     | -7     | -1             | -4             | 1               | 16              | 1         |      |
| -1     | -4     | -1             | -1             | 0               | 9               | 0         |      |
| 0      | 0      | 0              | 3              | 1               | 9               | 1         |      |
| 1      | 0      | 1              | 3              | 4               | 0               | 0         |      |
| 2      | -3     | 2              | 0              | 16              | 4               | 1         |      |
| 4      | -1     | 4              | 2              |                 |                 |           |      |
|        |        |                |                | $\sum U^1 = 36$ | $\sum V^1 = 55$ |           |      |
|        |        |                |                | $\frac{36}{8}$  | $\frac{55}{8}$  |           |      |

$$\sum U = 0$$

$$\sum V = 24$$

$$4.375$$

(37)

$$\lambda = \frac{\sum U'V'}{\sqrt{\sum (U')^2 \sum (V')^2}}$$

$$= \frac{24}{\sqrt{36 + 44}}$$

$$\lambda = \frac{24}{\sqrt{44}}$$

$$\lambda = \frac{4^2}{2\sqrt{11}}$$

$$\lambda = 2/\sqrt{11}$$

$$\boxed{\lambda = 0.6030}$$

Direct method:-  
 m m m

| $U = x - 68$ | $V = y - 69$ | $U^2$ | $V^2$ | $UV$ |
|--------------|--------------|-------|-------|------|
| -3           | -2           | 9     | 4     | 6    |
| -2           | -1           | 4     | 1     | 2    |
| -1           | -4           | 1     | 16    | 4    |
| -1           | -1           | 1     | 1     | 1    |
| 0            | 3            | 0     | 9     | 0    |
| 0            | 3            | 1     | 9     | 3    |
| 1            | 0            | 4     | 0     | 0    |
| 2            | 0            | 16    | 4     | 8    |
| 4            | 2            |       |       |      |

$$\sum U^2 = 36 \quad \sum V^2 = 44 \quad \sum UV = 24$$

(38)

$$\lambda = n * \sum_{UV} \epsilon_U \epsilon_V$$

$$\sqrt{n \sum_U \epsilon_U^2 - (\sum_U \epsilon_U)^2} \quad \sqrt{n * \sum_V \epsilon_V^2 - (\sum_V \epsilon_V)^2}$$

$$= \cancel{8 * 24} - 0$$

$$\sqrt{8(36) - 0} \quad \sqrt{8(44) - 0}$$

$$= \frac{242}{6 \times 2 \sqrt{11}}$$

$$= 2 / \sqrt{11}$$

$$\boxed{\lambda = 0.6030}$$

Rank Correlation :-

(39)

problem:- Ten competitors in a music test were ranked by the three judges A, B & C in the following order

|            |   |   |   |    |   |    |   |    |   |   |
|------------|---|---|---|----|---|----|---|----|---|---|
| Ranks of A | 1 | 6 | 5 | 10 | 3 | 2  | 4 | 9  | 7 | 8 |
| Ranks of B | 3 | 5 | 8 | 4  | 7 | 10 | 2 | 1  | 6 | 9 |
| Ranks of C | 6 | 4 | 9 | 8  | 1 | 2  | 3 | 10 | 5 | 7 |

using the rank correlation method discuss which pair of judges has the nearest approach to common licks in music

| Ranks A<br>X | Ranks B<br>Y | Ranks C<br>Z | $d_1 = x - y$ | $d_2 = x - z$ |
|--------------|--------------|--------------|---------------|---------------|
| 1            | 3            | 6            | -2            | -5            |
| 6            | 5            | 4            | 1             | 2             |
| 5            | 8            | 9            | -3            | -4            |
| 10           | 4            | 8            | 6             | 2             |
| 3            | 7            | 1            | -4            | 2             |
| 2            | 10           | 2            | -6            | 0             |
| 4            | 2            | 3            | 2             | 1             |
| 9            | 1            | 10           | 8             | -1            |
| 7            | 6            | 5            | 1             | 2             |
| 8            | 9            | 7            | -1            | 1             |

| $d_3 = y - z$    | $d_1^2$ | $d_2^2$         | $d_3^2$          |
|------------------|---------|-----------------|------------------|
| -3               | 4       | 25              | 9                |
| 1                | 1       | 4               | 1                |
| -1               | 9       | 16              | 1                |
| -4               | 36      | 4               | 16               |
| 6                | 16      | 4               | 36               |
| 8                | 36      | 0               | 64               |
| -1               | 4       | 1               | 1                |
| 9                | 64      | 1               | 81               |
| 1                | 1       | 4               | 1                |
| 2                | 1       | 1               | 4                |
| $\sum d_1 = 200$ |         | $\sum d_2 = 60$ | $\sum d_3 = 214$ |

(40)

$$P(x, y) = 1 - \left\{ \frac{6 * \sum d_1^2}{10(99)} \right\}$$

$$= 1 - \left\{ \frac{6 * 200}{10 * 99} \right\}$$

$$= 1 - \frac{40}{33}$$

$$P(x, y) = -0.2121$$

$$P(x, z) = 1 - \left\{ \frac{6 * \sum d_2^2}{10 * 99 / 33} \right\}$$

$$= 1 - \left\{ \frac{2 * 60}{33} \right\}$$

$$\neq 1 \neq 0.0606$$

≠

$$= 1 - \frac{4}{11}$$

$$= 11 - 4$$

$$P(x, z) = \frac{11}{11} = 0.6364$$

$$P(Y_1, Z) = 1 - \frac{8 * 214}{165 * 9933} \quad (41)$$

$$= 1 - \frac{214}{165}$$

$$= \frac{165 - 214}{165}$$

$$P(Y_1, Z) = -0.2970$$

$\therefore$  Inference  $P(X_1, Z)$  is maximum we concluded that pairs of judges A and C has the nearest approach common likings in music.

problem:- calculate the correlation coefficient b/w  $x$  &  $y$  for the following data

|     |   |   |   |    |    |    |    |
|-----|---|---|---|----|----|----|----|
| $x$ | 1 | 3 | 4 | 5  | 7  | 8  | 10 |
| $y$ | 2 | 6 | 8 | 10 | 14 | 16 | 20 |

| $x$                                | $y$ | $x = x - \bar{x}$ | $y = y - \bar{y}$ | $x^2$                  | $y^2$                  | $xy$                  |
|------------------------------------|-----|-------------------|-------------------|------------------------|------------------------|-----------------------|
| 1                                  | 2   | -4.4286           | -8.8571           | 19.6125                | 78.4482                | 39.224                |
| 3                                  | 6   | -2.4286           | -4.8571           | 5.8981                 | 23.5914                | 11.7961               |
| 4                                  | 8   | -1.4286           | -2.8571           | 2.0409                 | 8.1630                 | 4.081                 |
| 5                                  | 10  | -0.4286           | -0.8571           | 0.1837                 | 0.7346                 | 0.3674                |
| 7                                  | 14  | 1.5714            | 3.1429            | 2.4693                 | 9.8778                 | 4.9388                |
| 8                                  | 16  | 2.5714            | 5.1429            | 6.6121                 | 26.4494                | 13.2245               |
| 10                                 | 20  | 4.5714            | 9.1429            | 20.8977                | 83.5926                | 41.795                |
| $\Sigma x = 38$                    |     | $\Sigma y = 76$   |                   | $\Sigma x^2 = 57.7143$ | $\Sigma y^2 = 280.857$ | $\Sigma xy = 115.428$ |
| $\bar{x} = \frac{38}{7} = 5.4286$  |     |                   |                   |                        |                        |                       |
| $\bar{y} = \frac{76}{7} = 10.8571$ |     |                   |                   |                        |                        |                       |

(42)

$$\rho = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

$$= \frac{115.4289}{\sqrt{57.7143 * 230.857}}$$

$$\rho = \frac{115.4289}{\sqrt{13323.7502}} = \frac{115.4289}{115.4286}$$

$$\boxed{\rho = 1.0000}$$

- ② The following table gives the no of blind per lakh of population in different age groups  
find out the correlation b/w age & blindness.

| Age in years          | 0-10 | 10-20 | 20-30 | 30-40 | 40-50 | 50-60 | 60-70 | 70-80 |
|-----------------------|------|-------|-------|-------|-------|-------|-------|-------|
| no. of blind per lakh | 55   | 67    | 100   | 111   | 150   |       |       |       |
|                       | 200  | 300   | 500   |       |       |       |       |       |

| $x$   | $y$ | $\gamma$ | $x = \bar{x} - \bar{y}$ | $y = \bar{y} - \bar{x}$ | $x^2$                  | $y^2$                   | $xy$                    | $\sum x^2$     | $\sum y^2$ | $\sum xy$  |
|-------|-----|----------|-------------------------|-------------------------|------------------------|-------------------------|-------------------------|----------------|------------|------------|
| 0-10  | 10  | 65       | -30.625                 | -130.375                | 130.375                | 6997.6406               | 130.375                 | 130.375        | 3033.9844  | 3912.1344  |
| 10-20 | 15  | 67       | -25.625                 | -118.375                | 118.375                | 656.6406                | 14012.6406              | 14012.6406     | 1333.9844  | 3033.9844  |
| 20-30 | 25  | 100      | -15.625                 | -85.375                 | 85.375                 | 244.1406                | 7288.8906               | 7288.8906      | 418.3594   | 418.3594   |
| 30-40 | 35  | 111      | -5.625                  | -74.375                 | 74.375                 | 31.6406                 | 1251.3906               | 1251.3906      | -154.7656  | -154.7656  |
| 40-50 | 45  | 150      | 4.375                   | 41.375                  | 41.375                 | 19.1406                 | 213.8906                | 213.8906       | 210.2344   | 210.2344   |
| 50-60 | 55  | 200      | 14.375                  | 14.375                  | 14.375                 | 206.6406                | 13138.8906              | 13138.8906     | 279.844    | 279.844    |
| 60-70 | 65  | 300      | 24.375                  | 114.625                 | 114.625                | 594.1406                | 98988.8906              | 98988.8906     | 10815.2344 | 10815.2344 |
| 70-80 | 75  | 500      | 34.375                  | 314.625                 | 314.625                | 1181.6406               | 157423.8748             | 157423.8748    | $\leq y^2$ | $\leq y^2$ |
|       |     |          | $\sum x = 325$          | $\sum y = 1483$         | $\sum x^2 = 3871.8748$ | $\sum y^2 = 19929.1252$ | $\sum xy = 157423.8748$ | $= 19929.1252$ | $= 0.8072$ | $= 0.8072$ |

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}} = \frac{19929.1252}{\sqrt{3871.8748 * 157423.8748}} = \frac{19929.1252}{94688.5709} = 0.8012$$

2/11/2020  
Monday

(44)

### Partial Coefficient:-

#### Correlation

Sometimes the correlation b/w two variables  $x_1$  &  $x_2$  may be partly due to the

Coefficient of Partial Correlation:- sometimes the correlation b/w two variables  $x_1$  and  $x_2$  may be partly due to the correlation of the third variable,  $x_3$  with both  $x_1$  and  $x_2$ .

In such a situation, we may want to know, what the correlation b/w  $x_1$  and  $x_2$  would be if the effect of  $x_3$  on each of  $x_1$  and  $x_2$  would be if the effect of  $x_3$  on each of  $x_1$  and  $x_2$  were eliminated. This correlation is called partial correlation and the correlation coefficient b/w  $x_1$  and  $x_2$  after the linear effect of  $x_3$  on each of them has been eliminated is called the partial correlation coefficient.

Thus the partial correlation coefficient b/w  $x_1$  and  $x_2$  usually denoted by

$$\rho_{12.3} = \frac{\text{Cov}(x_{1.3}, x_{2.3})}{\sqrt{\text{Var}(x_{1.3}) \text{Var}(x_{2.3})}}$$

$$\boxed{\rho_{12.3} = \frac{\text{Cov}(x_{1.3}, x_{2.3})}{\sqrt{\text{Var}(x_{1.3}) \text{Var}(x_{2.3})}}}$$

$$\rho_{12 \cdot 3} = \frac{\rho_{12} - \rho_{13} * \rho_{23}}{\sqrt{(1 - \rho_{13}^2)(1 - \rho_{23}^2)}}$$

Note:- ① Similarly

$$\rho_{13 \cdot 2} = \frac{\rho_{13} - \rho_{12} * \rho_{32}}{\sqrt{(1 - \rho_{12}^2)(1 - \rho_{32}^2)}}$$

+

$$\rho_{23 \cdot 1} = \frac{\rho_{23} - \rho_{21} * \rho_{31}}{\sqrt{(1 - \rho_{21}^2)(1 - \rho_{31}^2)}}$$

- ② If  $\rho_{12 \cdot 3} = 0$ , we have then  $\rho_{12} = \rho_{13} * \rho_{23}$   
 If means that  $\rho_{12}$  will not be zero  
 If  $x_3$  is correlated with both  $x_1$  and  $x_2$   
 Thus although  $x_1$  and  $x_2$  may be uncorrelated  
 when effect of  $x_3$  is eliminated, yet  $x_1$   
 +  $x_2$  may appear to be correlated because  
 they carry the effect of  $x_3$  on them.
- ③ Partial correlation coefficient helps in  
 deciding whether to include (or) not an  
 additional independent variable in  
 regression analysis

Problems:- Find the partial correlation coefficient (76) from the following data

$$r_{12} = 0.77; r_{13} = 0.72, r_{23} = 0.52$$

soln:- given  $r_{12} = 0.77$

$$r_{13} = 0.72$$

$$r_{23} = 0.52$$

$$r_{12 \cdot 3} = r_{12} - r_{13} r_{23}$$

$$\frac{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}}{\sqrt{(1 - 0.72^2)(1 - 0.52^2)}}$$

$$= 0.77 - (0.72)(0.52)$$

$$\frac{\sqrt{(1 - 0.72^2)(1 - 0.52^2)}}{\sqrt{(0.4816)(0.7296)}}$$

$$= 0.77 - 0.3744$$

$$\frac{0.77 - 0.3744}{0.3514}$$

$$= \frac{0.3956}{\sqrt{0.3514}} = \frac{0.3956}{\sqrt{0.5988}}$$

$$\boxed{r_{12 \cdot 3} = 0.6674}$$

(47)

In a trivariant distribution it is found  
 $\tau_{12} = 0.7$ ,  $\tau_{13} = 0.61$ ,  $\tau_{23} = 0.4$ , Find the  
 Value of Partial correlation coefficient of  
 $\rho_{23.1}$ ,  $\rho_{13.2}$ ,  $\rho_{12.3}$ .

$$\begin{aligned}\rho_{23.1} &= \frac{\rho_{23} - \rho_{21} * \rho_{31}}{\sqrt{(1-\rho_{21}^2)(1-\rho_{31}^2)}} \\ &= \frac{0.4 - (0.7)(0.61)}{\sqrt{(1-0.7^2)(1-0.61^2)}} \\ &= -0.0477\end{aligned}$$

$$\begin{aligned}\rho_{13.2} &= \frac{\rho_{13} - \rho_{12} \rho_{23}}{\sqrt{(1-\rho_{12}^2)(1-\rho_{23}^2)}} \\ &= \frac{0.61 - (0.7)(0.4)}{\sqrt{(1-0.7^2)(1-0.4^2)}} \quad \frac{0.61 - (0.7)(0.4)}{\sqrt{(1-0.7^2)(1-0.4^2)}}\end{aligned}$$

$$\rho_{13.2} = -0.5041$$

$$\begin{aligned}\rho_{12.3} &= \frac{\rho_{12} - \rho_{13} \rho_{23}}{\sqrt{(1-\rho_{13}^2)(1-\rho_{23}^2)}} \\ &= \frac{(0.7) - (0.61)(0.4)}{\sqrt{(1-0.61^2)(1-0.4^2)}} = 0.6279\end{aligned}$$

<sup>(48)</sup>  
Multiple correlation :- In studying the joint effect of a group of variables upon a variable not included in the group such type of study is called multiple correlation.

For example, the yield of the crop per acre say  $x_1$ , depends on the quality of seeds  $x_2$ , fertility of soil  $x_3$ , fertilizers used  $x_4$ , irrigation say  $x_5$ , weather conditions say  $x_6$  and so on.

Whenever we are interested in studying the joint effect of a group of variables upon a variable not included in the group. One study is that of multiple correlation.

Coefficient of multiple correlation :- In trivariate distribution in which each of the variables  $x_1, x_2$  &  $x_3$  has  $N$  observations, the multiple correlation coefficient of  $x_1$  on  $x_2$  &  $x_3$  usually denoted by  $R_{1.23}$  is the simple correlation b/w  $x_1$  and the joint effect of  $x_2$  and  $x_3$  on  $x_1$ . In other words  $R_{1.23}$  is the correlation coefficient b/w  $x_1$  & its estimated value as given by the plane of regression of  $x_1$  on  $x_2 + x_3$ .

$$R_{1.23}^2 = \frac{r_{12}^2 + r_{23}^2 - 2r_{12}r_{13}r_{23}}{(1 - r_{13}^2)}$$

Similarly

(49)

$$R_{3,12}^2 = \frac{\lambda_{12}^2 + \lambda_{23}^2 - 2\lambda_{12}\lambda_{13}\lambda_{23}}{(1 - \lambda_{12}^2)}$$

$$R_{1,23}^2 = \frac{\lambda_{12}^2 + \lambda_{13}^2 - 2\lambda_{12}\lambda_{13}\lambda_{23}}{1 - \lambda_{23}^2}$$

$$R_{1,23}^2 = \frac{\lambda_{12}^2 + \lambda_{23}^2 - 2\lambda_{12}\lambda_{13}\lambda_{23}}{1 - \lambda_{13}^2}$$

06/11/2020

Friday

$$R_{1,23}^2 = \frac{\lambda_{12}^2 + \lambda_{13}^2 - 2\lambda_{12}\lambda_{13}\lambda_{23}}{1 - \lambda_{23}^2}$$

$$R_{3,12}^2 = \frac{\lambda_{13}^2 + \lambda_{23}^2 - 2\lambda_{12}\lambda_{13}\lambda_{23}}{1 - \lambda_{12}^2}$$

$$R_{2,13}^2 = \frac{\lambda_{12}^2 + \lambda_{13}^2 - 2\lambda_{12}\lambda_{13}\lambda_{23}}{1 - \lambda_{13}^2}$$

Properties of multiple correlation:-

(50)

- ① Multiple correlation coefficient measures the closeness of the association b/w the observed values of a variable obtained from the multiple linear regression of that variable on other variables.
- ② McC b/w observed values and the Expected Values, when the Expected Values are calculated from the linear relation of the Variable determined by the method of least squares is always greater than that where Expected values are calculated from any other linear combination of the variables.
- ③  $0 \leq R_{1.23} \leq 1$
- ④ If  $R_{1.23} = 0$ , then total f partial Correlations involving  $x_1$  are zero.
- ⑤ If  $R_{1.23} = 1$ , then association is perfect
- ⑥  $R_{1.23} \geq r_{12}, r_{13}, r_{23}$

Problem) - calculate  $R_{1.23}, R_{3.12}, R_{2.13}$  for the following data

$$r_{12} = 0.6, r_{13} = 0.7, r_{23} = 0.65$$

Soln given that

$$r_{12} = 0.6, r_{13} = 0.7, r_{23} = 0.65$$

$$R_{1.23}^2 = \frac{r_{12}^2 + r_{13}^2 - 2 * r_{12} * r_{13} * r_{23}}{1 - r_{23}^2}$$

$\frac{608}{1155}$

$$R_{1.23}^2 = (0.6)^2 + (0.7)^2 - 2 * (0.6)(0.7) / 0.65$$

$$R_{1\cdot 23} = 0.7255$$

(51)

$$\begin{aligned} R_{3\cdot 12}^2 &= \frac{r_{13}^2 + r_{23}^2 - 2r_{12}r_{13}r_{23}}{1 - r_{12}^2} \\ &= \frac{0.7^2 + 0.65^2 - 2(0.6 \times 0.7 \times 0.65)}{1 - 0.6^2} \\ &= \frac{733}{1280} \end{aligned}$$

$$R_{3\cdot 12} = 0.5727$$

$$R_{3\cdot 12} = \sqrt{0.5727}$$

$$R_{3\cdot 12} = 0.7568$$

$$\begin{aligned} R_{2\cdot 13}^2 &= \frac{r_{12}^2 + r_{23}^2 - 2r_{12}r_{13}r_{23}}{1 - r_{13}^2} \\ &= \frac{0.6^2 + 0.65^2 - 2 \times 0.6 \times 0.7 \times 0.65}{1 - 0.7^2} \\ &= 0.4637 \end{aligned}$$

$$R_{2\cdot 13} = \sqrt{0.4637}$$

$$R_{2\cdot 13} = 0.6810$$

problem:-

(52)

$$^2 \text{ If } r_{12} = 0.9, r_{13} = 0.75 \text{ and } r_{23} = 0.7$$

find  $R_{1.23}, R_{3.12}, R_{2.13}$

$$R_{1.23}^2 = r_{12}^2 + r_{13}^2 - 2 * r_{12} * r_{13} * r_{23}$$

$$\frac{1 - r_{23}^2}{1 - r_{23}^2}$$

$$= 0.9^2 + 0.75^2 - 2 * 0.9 * 0.75 * 0.7$$
$$\frac{1 - 0.7^2}{1 - 0.7^2}$$

$$= \frac{57}{68}$$

$$R_{1.23}^2 = \frac{57}{68}$$

$$R_{1.23} = \sqrt{\frac{57}{68}}$$

$$R_{1.23} = 0.9156$$

$$R_{3.12}^2 = r_{13}^2 + r_{12}^2 - 2 * r_{12} * r_{13} * r_{23}$$

$$\frac{1 - r_{12}^2}{1 - r_{12}^2}$$

$$= 0.75^2 + 0.9^2 - 2 * 0.9 * 0.75 * 0.7$$
$$\frac{1 - 0.9^2}{1 - 0.9^2}$$

$$R_{3.12}^2 = 2.25$$

$$R_{3.12} = \sqrt{2.25}$$

$$R_{3.12} = 1.5$$

(53)

$$R_{2 \cdot 13}^2 = \frac{\lambda_{12}^2 + \lambda_{13}^2 - 2\lambda_{12}\lambda_{13}\lambda_{23}}{1 - \lambda_{13}^2}$$

$$= \frac{0.9^2 + 0.75^2 - (2 * 0.9 * 0.75 * 0.7)}{1 - 0.75^2}$$

$$= 0.9771$$

$$R_{2 \cdot 13} = \sqrt{0.9771}$$

$$R_{2 \cdot 13} = 0.9885$$

given that  $\tau_{12} = 0.4$ ,  $\tau_{13} = 0.5$ ,  $\tau_{23} = 0.6$

find  $\tau_{12 \cdot 3}$  &  $R_{12 \cdot 3}$

$$\tau_{12 \cdot 3} = \frac{\lambda_{12} - \lambda_{13}\lambda_{23}}{1 - \lambda_{13}^2}$$

$$\sqrt{(1 - \lambda_{13}^2)(1 - \lambda_{23}^2)}$$

$$= \frac{0.4 - (0.5)(0.6)}{\sqrt{(1 - 0.5^2)(1 - 0.6^2)}}$$

$$= \frac{\sqrt{3}}{12}$$

$$\boxed{\lambda_{12 \cdot 3} = 0.1443}$$

$$R_{1 \cdot 23}^2 = \frac{0.4^2 + 0.5^2 - 2 * 0.4 * 0.5 * 0.6}{1 - 0.6^2}$$

$$R_{1 \cdot 23}^2 = 0.2656$$

$$R_{1 \cdot 23} = 0.5154$$

(54)

Given that  $r_{12} = 0.6$ ,  $r_{13} = 0.72$  &  $r_{23} = 0.65$

Find  $\alpha_{123}$ ,  $R_{123}$  &  $\alpha_{12 \cdot 3}$

$$R_{1 \cdot 2 \cdot 3}^2 = \frac{\alpha_{12}^2 + \alpha_{13}^2 - 2 * \alpha_{12} * \alpha_{13} * \alpha_{23}}{1 - \alpha_{23}^2}$$

$$= \frac{0.6^2 + 0.72^2 - 2 * 0.6 * 0.72 * 0.65}{1 - 0.65^2}$$

$$R_{1 \cdot 2 \cdot 3}^2 = 0.5486$$

$$R_{1 \cdot 2 \cdot 3} = 0.7407$$

$$\alpha_{12 \cdot 3} = \frac{\alpha_{12} - \alpha_{13} \alpha_{23}}{\sqrt{(1 - \alpha_{13}^2)(1 - \alpha_{23}^2)}}$$

$$= \frac{0.6 - 0.72 * 0.65}{\sqrt{(1 - 0.72^2)(1 - 0.65^2)}}$$

$$\boxed{\alpha_{12 \cdot 3} = 0.2432}$$

Correlation ratio:-

Correlation ratio ' $\eta$ ' is the appropriate measure of curvilinear relationship b/w the two variables

just as  $\alpha$  measures the concentration of points about the curve of best fit.  $\eta$  measures the concentration of points about the curve of best fit.

Note:- If Regression is Linear,  $\eta = e^{55}$  otherwise  
 $|\eta| \geq 10$

Measure of Correlation ratio:-

In Correlation we observe only one value of  $y$  corresponding to  $\underline{\text{that}}$  value of  $x$ . There may be some cases where variable  $y$  takes the value  $y_{ij}$  with frequencies  $f_{ij}$  where  $j = 1, 2, \dots, n$  corresponding to the value of  $x = x_i$  where  $i = 1, 2, \dots, m$ . These values can be shown in the following table.

| $x_1$    | $x_2$            | $\dots$          | $x_i$            | $\dots$          | $x_m$    | $n_j$    |
|----------|------------------|------------------|------------------|------------------|----------|----------|
| $y_1$    | $f_{11}(y_{11})$ | $f_{21}(y_{21})$ | $f_{i1}(y_{i1})$ | $f_{m1}(y_{m1})$ | $n_1$    |          |
| $y_2$    | $f_{12}(y_{12})$ | $f_{22}(y_{22})$ | $f_{i2}(y_{i2})$ | $f_{m2}(y_{m2})$ | $n_2$    |          |
| $\vdots$ | $\vdots$         | $\vdots$         | $\vdots$         | $\vdots$         | $\vdots$ | $\vdots$ |
| $y_j$    | $f_{1j}(y_{1j})$ | $y_{2j}(y_{2j})$ | $f_{ij}(y_{ij})$ | $f_{mj}(y_{mj})$ | $\vdots$ | $n_j$    |
| $\vdots$ | $\vdots$         | $\vdots$         | $\vdots$         | $\vdots$         | $\vdots$ | $\vdots$ |
| $y_n$    | $f_{1n}(y_{1n})$ | $f_{2n}(y_{2n})$ | $f_{in}(y_{in})$ | $f_{mn}(y_{mn})$ | $\vdots$ | $n_n$    |
| $n_1$    | $\dots$          | $n_2$            | $\dots$          | $n_i$            | $\dots$  | $n_m$    |
| $N$      |                  |                  |                  |                  |          |          |

Here pair of observations  $(x_i, y_{ij})$   
 with the frequencies  $f_{ij}$  so that 56

$$N = \sum_{i=1}^m \sum_{j=1}^n f_{ij}$$

$$n_i = \sum_{j=1}^n f_{ij} \quad j = 1, 2, \dots, n$$

$$n_j = \sum_{i=1}^m f_{ij}$$

$$T_i = \sum_{j=1}^n f_{ij} y_{ij}$$

$$T = \sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij} = \sum_{i=1}^m T_i$$

Let us denote  $\bar{y}_i, \bar{y}$  be the means of  
 i<sup>th</sup> array and overall mean respectively

$$\bar{y}_i = \frac{\sum_{j=1}^n f_{ij} y_{ij}}{\sum_{j=1}^n f_{ij}} = \frac{T_i}{n_i}$$

$$\boxed{\bar{y}_i = \frac{T_i}{n_i}}$$

$$\bar{y} = \frac{\sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij}}{\sum_{i=1}^m \sum_{j=1}^n f_{ij}} = \frac{T}{N}$$

$$\bar{y} = \frac{T}{N}$$

where

$$T = \bar{y} \cdot N$$

The correlation ratio of  $y$  on  $x$  is denoted by  $\eta_{yx}^2$  defined as

$$\eta_{yx}^2 = \frac{\sum_{i=1}^m \left( \frac{T_i^2}{n_i} \right) - \frac{T^2}{N}}{\sum_{i=1}^m \sum_{j=1}^n f_{ij} (y_{ij})^2 - \left( \frac{T^2}{N} \right)}$$

Similarly the correlation ratio of  $x$  on  $y$  is denoted by  $\eta_{xy}^2$  defined as

$$\eta_{xy}^2 = \frac{\sum_{j=1}^n \frac{T_j^2}{n_j} - \left( \frac{T^2}{N} \right)}{\sum_{i=1}^m \sum_{j=1}^n f_{ij} (y_{ij})^2 - \left( \frac{T^2}{N} \right)}$$

$$N = \sum_{i=1}^m \sum_{j=1}^n f_{ij} ; \quad n_i = \sum_{j=1}^n f_{ij} ; \quad n_j = \sum_{i=1}^m f_{ij}$$

$$T_i = \sum_{j=1}^n f_{ij} y_{ij} ; \quad T = \sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij} = \sum_{i=1}^m T_i$$

Properties of correlation ratio:-

1. Correlation ratio is independent of change of origin & scale

i.e,

$$\eta_{yx}^2 = \eta_{xy}^2 * \eta_{uv}^2$$

2. limits of correlation ratio are 0 to 1

$$\text{i.e } 0 \leq \eta_{yx}^2 \leq 1$$

Also

$$0 \leq \eta_{yx} \leq 1$$

9/11/2020

obtain correlation ratio of y on x to the following data

| $y \backslash x$ | 10 | 50 | 20 | 25 |
|------------------|----|----|----|----|
| 7                | 3  | 2  | -  | -  |
| 9                | -  | 1  | 4  | 6  |
| 11               | -  | 3  | 4  | 2  |
| 13               | 2  | 1  | 5  | -  |
| 15               | -  | 6  | -  | 1  |

The  $\gamma$  on  $x$  correlation is

(59)

$$\eta_{yx}^2 = \frac{\sum_{i=1}^m \left( \frac{T_i^2}{n_i} \right) - \left( \frac{T^2}{N} \right)}{\sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij}^2 - \frac{T^2}{N}}$$

where  $N = \sum_{i=1}^m \sum_{j=1}^n f_{ij}$

$$n_i = \sum_{j=1}^n f_{ij} \quad T = \sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij} = \sum_{i=1}^m T_i$$

$$n_j = \sum_{i=1}^m f_{ij}$$

$$\bar{T}_i = \frac{\sum_{j=1}^n f_{ij} y_{ij}}{n_j}$$

$$\bar{y}_i = \frac{\sum_{j=1}^n f_{ij} y_{ij}}{\sum_{j=1}^n f_{ij}} = \frac{\bar{T}_i}{n_i}$$

$$\boxed{\bar{T}_i = n_i \cdot \bar{y}_i}$$

$$\bar{y} = \frac{\sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij}}{\sum_{i=1}^m \sum_{j=1}^n f_{ij}} = \frac{T}{N}$$

$$\boxed{T = N \cdot \bar{y}}$$

13/11/2020

Fri day

| $f_{ij} y_{ij}$              | 10           | 15    | 20    | 25      | $\eta_j$ | $f_{ij} y_{ij}$               | $\sum_{i=1}^m f_{ij} y_{ij}$                                          |
|------------------------------|--------------|-------|-------|---------|----------|-------------------------------|-----------------------------------------------------------------------|
| 7                            | 3            | 2     | -     | -       | 5        | $7(3) + 7(2) = 35$            | $3 \times 7^2 + 2 \times 7^2 = 845$                                   |
| 9                            | -            | 1     | 4     | 6       | 11       | $9(1) + 9(9) + 9(6) = 99$     | $1 \times 9^2 + 4 \times 9^2 + 6 \times 9^2 = 891$                    |
| 11                           | -            | 3     | 4     | 2       | 9        | $11(3) + 11(9) + 11(2) = 99$  | $3 \times 11^2 + 4 \times 11^2 + 4 \times 11^2 = 1089$                |
| 13                           | 2            | -     | 5     | -       | 8        | $13(2) + 13(1) + 13(5) = 104$ | $2 \times 13^2 + 1 \times 13^2 + 5 \times 13^2 = 1352$                |
| 15                           | -            | 1     | -     | -       | 7        | $15(6) + 15(1) + 105$         | $6 \times 15^2 + 1 \times 15^2 = 1572$                                |
| $\sum_{i=1}^m f_{ij} y_{ij}$ | 7(3) + 13(2) | 15(9) | 145   | 91      | N = 40   | $T = \sum_i T_i = 442$        | $T = 5152 = \sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij}$                  |
| $T_i^2$                      | 47           | 2209  | 25281 | 21025   | 5676     | $\sum_{i=1}^m T_i^2 = 5676$   | $\sigma_{xy}^2 = \frac{\sum_i T_i^2}{n} - \left(\frac{T}{N}\right)^2$ |
| $\frac{T^2}{n}$              | 441.8        | 8209  | 5     | 1944.69 | 16173    | 91011                         | $= 4923.90 - 4884.1$                                                  |

$$\sigma_{xy}^2 = 0.1986 \Rightarrow \eta_{xy} = \sqrt{0.1986} \Rightarrow 0.3855$$

## \* Difference b/w Correlation & Regression:-

(6)

- |                                                                                                                                                |                                                                                                                                                                                                                                                                                                                                          |
|------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 1. Correlation means the relationship b/w two (or) more variables. It measures effect of one variable due to the change in the other variable. | 1. Regression means stepping back to the mean value. It Expresses average relationship b/w two (or) more variables.                                                                                                                                                                                                                      |
| 2. Correlation did not need not imply cause and effect relationship b/w the variables under study.                                             | 2. But Regression analysis indicates the cause & effect relationship. The variable constituting cause is taken as independent variable & the variable constituting the effect is taken as dependent variable.                                                                                                                            |
| 3. Correlation analysis measure linear relationship only b/w the variables. Hence practical applications are very less.                        | 3. Regression analysis measure the linear & non linear relationships of the variable. Hence practical applications are more.                                                                                                                                                                                                             |
| 4. Sometimes correlation may be non-sense                                                                                                      | 4. There is no non-sense regression.                                                                                                                                                                                                                                                                                                     |
| 5. Correlation measures the direction & degree of linear relationship. It is symmetric. i.e.,<br>$r_{xy} = r_{yx}$ .                           | 5. Regression analysis measures the functional relationship b/w the variables and it also used to estimate the value of dependent variable for any given independent variable. It identifies the nature of the variable i.e; which is independent and which is dependent i.e, $b_{yx} \neq b_{xy}$ . Regression coeff are not symmetric. |

H/W:-

(62)

- ① calculate rank correlation coefficient for the following data

|   |    |    |    |    |    |    |    |    |    |    |
|---|----|----|----|----|----|----|----|----|----|----|
| X | 10 | 15 | 12 | 17 | 13 | 16 | 24 | 14 | 22 | 20 |
| Y | 30 | 42 | 45 | 46 | 33 | 34 | 40 | 35 | 39 | 38 |

2. calculate the coefficient of correlation for ranks from the following

$(X, Y) : (5, 8), (10, 3), (6, 2), (3, 9), (19, 12), (5, 3), (16, 17), (12, 18), (8, 22), (2, 12), (10, 17), (19, 20)$

Analysis of Variance (ANOVA) :-

Analysis of variance (ANOVA) is the separation of variance ascribable to one group of causes from the variance ascribable to other group. ANOVA consists in the estimation of amount of variation due to each of the independent factors (or) causes separately and then comparing these estimates due to assignable factors with the estimates due to chance factors, the later being known as experimental error.

- \*\*\*
- 1. ANOVA I one way classification
  - 2. ANOVA II two way classification.

i.e.,  $b_{yx} \neq b_{xy}$

~~\*\*\*~~ ANOVA I way classification:-

Formulae:-

RSS = Raw sum of squares

$$\boxed{\text{RSS} = \sum_{i=1}^k \sum_{j=1}^n x_{ij}^2}$$

$$G = \text{Grand total} = \sum_{i=1}^k \sum_{j=1}^n x_{ij}$$

$$CF = \text{correction factor} = \frac{G^2}{N}$$

N = Total no of observations

TSS = Total sum of squares

$$\boxed{TSS = RSS - CF}$$

$$SST = \text{Treatment sum of squares} = \sum_{i=1}^k \left( \frac{\bar{T}_{ii}}{n_i} \right)^2 - CF$$

SSE = Error sum of squares (Sum of square due to error)

$$\boxed{SSE = TSS - SST}$$

ANOVA I way Table:-

| 10 |  |  |  |  |
|----|--|--|--|--|
|    |  |  |  |  |
|    |  |  |  |  |
|    |  |  |  |  |

65

| Source of variation | Variance            |                         |
|---------------------|---------------------|-------------------------|
|                     | mean sum of squares | Ratio                   |
| Treatments          | $k-1$               | $MST = \frac{SST}{k-1}$ |
| Error               | $n-k$               | $MSE = \frac{SSE}{n-k}$ |
| Total               | $n-1$               |                         |

$F = \frac{MST}{MSE} \sim F_{(k-1), (n-k)}$

Show that a significance test doesn't reject their homogeneity.

Problem:- The following shows the times (in hours) of four batches of electric bulbs in hrs

| Batches | 1    | 2    | 3    | 4    |
|---------|------|------|------|------|
| 1       | 1600 | 1610 | 1650 | 1680 |
| 2       | 1580 | 1640 | 1640 | 1700 |
| 3       | 1460 | 1550 | 1600 | 1520 |
| 4       | 1510 | 1530 | 1570 | 1600 |

| Batches | 1    | 2    | 3    | 4    |
|---------|------|------|------|------|
| 1       | 1700 | 1720 | 1800 | —    |
| 2       | 1750 | —    | —    | —    |
| 3       | 1640 | 1620 | 1640 | 1740 |
| 4       | 1600 | 1600 | 1680 | 1820 |

23/11/2020

(65)

1. From the following data, compute the coefficient of correlation b/w X & Y

| no. of items                                                                         | X-series | Y-series |
|--------------------------------------------------------------------------------------|----------|----------|
| Arithmetic mean                                                                      | 15       | 15       |
| sum of squares of deviation from mean                                                | 25       | 18       |
| summation of product of deviations of X and Y series from respective arithmetic mean | 136      | 138      |
| $\Sigma xy = 122$                                                                    |          |          |

| no. of items                          | X                                    | Y   |
|---------------------------------------|--------------------------------------|-----|
| Arithmetic mean                       | 15                                   | 15  |
|                                       | 25                                   | 18  |
| Sum of squares of deviation from mean | 136                                  | 138 |
|                                       | $x = 15, y = 15$                     |     |
|                                       | $\bar{x} = 25, \bar{y} = 18$         |     |
|                                       | $\Sigma x^2 = 136, \Sigma y^2 = 138$ |     |
|                                       | $\Sigma xy = 122$                    |     |

$$\rho = \frac{\Sigma xy}{\sqrt{\Sigma x^2 \Sigma y^2}}$$

$$= \frac{122}{\sqrt{136 * 138}}$$

$$\rho = 0.8905$$

A sample of 12 fathers and their eldest son gave the following data about their height in inches.

Father 65 63 67 64 68 62 70 66 68  
67 69 71

Sons 68 66 68 65 69 66 68 65  
61 71 67 68 70

calculate the coefficient of rank correlation.

| x  | y  | Rank of $x_i$ | Rank of $y_i$ | $x_i - y_i$ | $d_i^2$                       |
|----|----|---------------|---------------|-------------|-------------------------------|
| 65 | 68 | 9             | 5.5           | 3.5         | 12.25                         |
| 63 | 66 | 11            | 9             | 2           | 4                             |
| 67 | 68 | 6.5           | 5.5           | 1           | 1                             |
| 64 | 65 | 10            | 11.5          | -1.5        | 2.25                          |
| 68 | 69 | 4.5           | 3             | 1.5         | 2.25                          |
| 62 | 66 | 12            | 10            | 2           | 4                             |
| 70 | 68 | 2             | 5.5           | -3.5        | 12.25                         |
| 66 | 65 | 8             | 10.5          | -3.5        | 12.25                         |
| 68 | 71 | 4.5           | 1             | 3.5         | 12.25                         |
| 67 | 67 | 6.5           | 8.5           | -1.5        | 2.25                          |
| 69 | 68 | 3             | 5.5           | -2.5        | 6.25                          |
| 71 | 70 | 1             | 2             | -1          | $\frac{4+5}{2} = \frac{9}{2}$ |

$$\sum d_i^2 = 72$$

$$CF \text{ in } X \text{ series} = 67 - 2, 68 - 2 \quad \frac{4+5+6+7}{4} = \frac{6+7}{2}$$

$$2(2^2 - 1) + \frac{2(2^2 - 1)}{12} = \frac{13}{2}$$

$$CF_x = 0.5 + 0.5 = 1$$

$$CF_y = 68 - 2\sqrt{65 - 2}$$

$$CF_y = \frac{4(4^2 - 1)}{12} + \frac{2(2^2 - 1)}{12}$$

$$CF_y = \frac{4(15)}{12 \times 3} + \frac{2(3)}{12 \times 2}$$

$$= 5 + \frac{1}{2}$$

$$= \frac{11}{2}$$

$$CF_y = 5.5 - 4.5 = 5.5$$

$$CF_x + CF_y = 1 + 5.5 = 6.5$$

$$\rho = 1 - \left[ \frac{6(72 + 6.5)}{12(143)} \right]$$

$$= 1 - \left[ \frac{6 \times 78.5}{12 \times 143} \right]$$

$$= 1 - \left[ \frac{471}{1716} \right]$$

$$= 1 - 0.2745$$

$$= 0.7255 \quad (0.7220)$$

68

Ten competitors in a beauty contest are ranked by three judges as follows

| judges | 1 | 2 | 3 | 4  | 5  | 6 | 7  | 8 | 9 | 10 |
|--------|---|---|---|----|----|---|----|---|---|----|
| A      | 6 | 5 | 3 | 10 | 2  | 4 | 9  | 7 | 8 | 1  |
| B      | 5 | 8 | 4 | 7  | 10 | 2 | 1  | 6 | 9 | 3  |
| C      | 4 | 9 | 8 | 1  | 2  | 3 | 10 | 5 | 7 | 6  |

Discuss which pair has the nearest approach to common tastes of beauty.

| Rank A | Rank B | Rank C | $d_1 = x - y$ | $d_2 = x - z$ | $d_3 = y - z$ | $d_1^2$ | $d_2^2$ | $d_3^2$ |
|--------|--------|--------|---------------|---------------|---------------|---------|---------|---------|
| X      | Y      | Z      | 1             | 5             | 1             | 25      | 1       | 1       |
| 6      | 5      | 4      | -3            | -4            | -1            | 16      | 1       | 16      |
| 5      | 8      | 8      | -1            | -5            | -4            | 25      | 36      | 36      |
| 3      | 4      | 1      | 3             | 9             | 6             | 81      | 0       | 64      |
| 9      | 10     | 7      | -8            | 0             | 8             | 64      | -1      | 81      |
| 2      | 2      | 2      | 2             | 1             | -1            | 9       | 1       | 9       |
| 10     | 7      | 3      | 8             | 1             | 1             | 64      | 1       | 9       |
| 2      | 10     | 5      | -1            | 2             | 2             | 1       | 1       | 9       |
| 6      | 9      | 7      | -1            | -1            | -1            | 2       | -3      | 25      |
| 9      | 9      | 6      | 2             | 1             | 1             | -5      | 4       | 25      |
| 7      | 7      | 7      | 6             | 6             | 7             | 1       | 1       | 25      |
| 8      | 8      | 9      | 9             | 9             | 3             | -3      | 4       | 25      |
| 1      | 1      | 1      |               |               |               |         |         |         |

$$\begin{aligned}
 P(X_1Y) &= 1 - \left\{ \frac{6 * \leq d_1^2}{n(n^2-1)} \right\} \\
 &= 1 - \left\{ \frac{\frac{2}{6} * 158}{10(99) / 33} \right\} \\
 &= 1 - \left\{ \frac{158}{\cancel{5} \cancel{13}} \right\} \\
 &= 1 - \left\{ \frac{158}{165} \right\} \\
 &= 1 - \left\{ \frac{158}{165} \right\} \\
 &\approx 1 - 0.9576 \\
 &\approx 0.0424
 \end{aligned}$$

$$\begin{aligned}
 P(X_1Z) &= 1 - \left\{ \frac{6 * 179}{10(99)} \right\} \\
 &= 1 - \left\{ \frac{1074}{990} \right\} \\
 P(X_1Z) &\approx 0.0848
 \end{aligned}$$

$$\begin{aligned}
 P(Y_1Z) &= 1 - \left\{ \frac{6 * 214}{10 * 99} \right\} \\
 &= 1 - \left\{ \frac{1284}{990} \right\} \\
 &= 1 - 1.2940 \\
 P(Y_1Z) &\approx 0.2970
 \end{aligned}$$

27/11/2020

Marks of 6 students of a class in paper I and paper II

(70)

1) both Regression coefficients

2) both Regression lines

3) Re. correlation coefficient

paper I 45 55 66 75 85 100

paper II 56 55 45 65 62 71

| $x$        | $y$        | $x - \bar{x}$ | $y - \bar{y}$ | $x^2$        | $y^2$        | $xy$        |
|------------|------------|---------------|---------------|--------------|--------------|-------------|
| 45         | 56         | -11.26        | -3            | 676          | 9            | 78          |
| 55         | 55         | -6.16         | -4            | 256          | 16           | 64          |
| 66         | 45         | 21.5          | -14           | 25           | 196          | 70          |
| 75         | 65         | 16.4          | 6             | 16           | 36           | 24          |
| 85         | 62         | 23.4          | 3             | 196          | 9            | 42          |
| 100        | 71         | 29.29         | 12            | 841          | 144          | 348         |
| $\sum x =$ | $\sum y =$ |               |               | $\sum x^2 =$ | $\sum y^2 =$ | $\sum xy =$ |
| 426        | 354        |               |               | 2010         | 410          | 626         |
| 71         | 59         |               |               |              |              |             |

$$\lambda = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

$$= \frac{626}{\sqrt{2010 * 410}}$$

$$\bar{x} = \sqrt{\frac{1}{n} \sum x^2}$$

$$= \sqrt{\frac{1}{6} * 2010}$$

$$\bar{x} = 18.3030$$

$$\lambda = 0.6896$$

$$\bar{y} = \sqrt{\frac{1}{6} * 410}$$

$$\bar{y} = 8.2664$$

71

$$b_{yx} = 2 \frac{\sigma_y}{\sigma_x}$$

$$= 0.6896 * \frac{8.2664}{18.3030}$$

$$= 0.3114$$

$$b_{yx} = \frac{628}{2010} \frac{\sum xy}{\sum x^2}$$

$$\frac{= 314}{1005}$$

$$= 0.3124$$

$$b_{yx}$$

$$b_{xy} = 2 \frac{\sigma_x}{\sigma_y}$$

$$= 0.6896 * \frac{18.3030}{8.2664}$$

$$b_{xy} = 1.5269$$

$$b_{xy} = \frac{\sum xy}{\sum y^2}$$

$$\frac{= 628}{410}$$

$$b_{xy} = 1.5317$$

Regression line  $y$  on  $x$  :-

$$y - \bar{y} = 2 \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$y - 59 = 0.3124 (x - 71)$$

$$y - 59 = 0.3124 x - 22.1804$$

$$y - 59 = 0.3124 x + 36.8196$$

$$y = 0.3124 x + 36.8196$$

Regression line on  $x$  on  $y$  :-

$$x - \bar{x} = 2 \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

72

$$x - 71 = 1.5317(y - 59)$$

$$x - 71 = 1.5317y - 90.3703$$

$$x = 1.5317y + 19.3703$$

$$x = 1.5317 * 0.49 + 19.3703$$

$$\approx 0.7505 + 19.3703$$

$$\boxed{x = 20.1208}$$

$$\boxed{x = -18.6198}$$

$y$  for  $x = 0.2$

$$y = 0.3124 * 0.2 + 36.8196$$

$$\approx 0.06248 + 36.8196$$

$$\boxed{y \approx 36.8820}$$

2. regression line of  $y$  on  $x$  on  $y$  respectively  
are

$$2x - 3y = 8$$

$$5x - y = 6$$

then find

(i) mean values of  $x$  and  $y$

(ii) coefficient of correlation b/w  $x$  and  $y$

(iii) standard deviation of  $y$  for given

Variance of  $x = 5$

eqn ① as y on x & eqn ② is x on y 7.

$$2x - 3y = -8$$

$$-3y = -8 + 2x$$

$$-3y = 2x - 8$$

$$y = \frac{+2x + 8}{3} \quad \frac{3}{3}$$

$$byx = +\frac{2}{3}$$

$$2x - 3y = -8$$

$$2x = 3y + 8$$

$$x = \frac{3}{2}y + \frac{8}{2}$$

$$bxy = 3/2$$

eqn ② is x on y

$$5x - y = 6$$

$$-y = -5x + 6$$

$$5x = y + 6$$

$$y = 5x - 6$$

$$x = \frac{y}{5} + \frac{6}{5}$$

$$byx = 5$$

$$bxy = 1/5$$

$$d = \sqrt{bxy * byx} \quad d = \sqrt{3/2 * 5}$$

$$= \sqrt{(+\frac{2}{3})(\frac{1}{5})}$$

$$= \sqrt{15/2}$$

$$d = 2.7386$$

$$= 0.3651$$

$$2x - 3y = -8$$

$$x = 2$$

$$5x - y = 6$$

$$y = 4$$

(74)

$$g = \sqrt{b_{xy} * b_{yx}}$$

$$= \sqrt{\left(\frac{2}{3}\right) * \frac{1}{5}}$$

$$= \sqrt{\frac{2}{15}}$$

$$\boxed{g = 0.3651}$$

$$b_{yx} = \frac{\sigma_y}{\sigma_x}$$

$$\frac{2}{3} = 0.3651 \left( \frac{\sigma_y}{\sigma_x} \right)$$

$$\frac{2}{3} \times 5 = 0.3651 (\sigma_y)$$

$$\frac{10}{3 \times 0.3651} = \sigma_y$$

$$\frac{2}{3} = \frac{0.3651}{2.2360} (\sigma_y)$$

$$\frac{4.4721}{1.6953} = \sigma_y$$

$$\sigma_x^2 = 5$$

$$\sigma_x =$$

~~$$\sigma_y = b_{yx} * b_{xy}$$~~

$$\sigma_y = \frac{\sigma_x * b_{yx}}{r}$$

$$0.6492$$

1. calculate the Karl Pearson's coefficient of correlation b/w price and demand for the following data

|        |    |    |    |    |    |    |    |    |    |
|--------|----|----|----|----|----|----|----|----|----|
| price  | 17 | 18 | 19 | 20 | 22 | 24 | 26 | 28 | 30 |
| demand | 40 | 38 | 35 | 30 | 28 | 25 | 22 | 21 | 20 |

| x          | y                | $x - \bar{x}$ | $y - \bar{y}$ | $x^2$              | $y^2$                | $xy$ |
|------------|------------------|---------------|---------------|--------------------|----------------------|------|
| 17         | 40               | -5.6667       | 11.2222       | 32.1115            | 125.9378             |      |
| 18         | 38               | -4.6667       | 9.2222        | 21.778             | 85.0490              |      |
| 19         | 35               | -3.6667       | 6.2222        | 13.4447            | 38.7158              |      |
| 20         | 30               | -2.6667       | 1.2222        | 7.1113             | 1.4938               |      |
| 22         | 28               | -0.6667       | -0.7778       | 0.4445             | 0.6090               |      |
| 24         | 25               | 1.3333        | -3.7778       | 1.7777             | 14.2718              |      |
| 26         | 22               | 3.3333        | -6.7778       | 11.1109            | 45.9386              |      |
| 28         | 21               | 5.3333        | -7.7778       | 28.4441            | 60.4926              |      |
| 30         | 20               | 7.3333        | -8.7778       | 53.7773            | 53.77                |      |
| $\Sigma x$ | $\Sigma y = 259$ |               |               | $\Sigma x^2 = 170$ | $\Sigma y^2 = 381.0$ |      |
| = 204      | = 28.7778        |               |               |                    |                      |      |
|            |                  | = 22.6667     |               |                    |                      |      |
|            |                  |               |               |                    |                      |      |
|            |                  |               |               |                    |                      |      |

$$-63.5928 \quad \Sigma xy = -265.6663$$

$$-43.0372$$

$$-22.8149$$

$$-3.2592$$

$$0.5186$$

$$-5.0369$$

$$-22.5924$$

$$-41.4813$$

$$-4.3702$$

$$\rho = \frac{-265.6663}{\sqrt{170 * 381.0739}}$$

$$\sqrt{170 * 381.0739}$$

30/11/2020  
Monday

(76)

$H_0: \mu_1 = \mu_2 = \mu_3 = \mu_4$  (null hypothesis)

$H_1: \mu_1 \neq \mu_2 \neq \mu_3 \neq \mu_4$  (Alternative hypothesis)

Coding method:- subtract each number from 1600  
~~then~~  
 - & divide by 10.

|   |      |    |    |     |     |     |     | $\sum_{j=1}^{120}$ |
|---|------|----|----|-----|-----|-----|-----|--------------------|
| 1 | 0    | -1 | -5 | -8  | -10 | -12 | -20 | -                  |
| 2 | +102 | -4 | -4 | -10 | -15 | -   | -   | -56                |
| 3 | +154 | 5  | 0  | -2  | -4  | -6  | -14 | -22                |
| 4 | 9    | 8  | 7  | 3   | 0   | -8  | -   | -19                |
|   |      |    |    |     |     |     |     | <u>G = -97</u>     |

$x_{ij}^2$

$$N = 26$$

$$RSS = 2319$$

134

$$G = -97 \text{ (grand total)}$$

361

$$G = \sum \sum x_{ij} = -97$$

157

$$\frac{67}{= 2319} \quad \text{correction factor} = G^2 = \frac{9409}{N} = \frac{9409}{26}$$

$$\text{Correction factor} = 361.8846.$$

Total sum of squares =  $T_{SS}$

$$= RSS - CF$$

$$= 2319 - 361.8846$$

$$= 1957.1154$$

Sum of squares due to  
Total

sum of squares due to Treatments

(77)

$$= SST = \sum_{i=1}^k \left( \frac{T_i^2}{n_i} \right) - CF$$

$$= \frac{(-56)^2}{7} + \frac{(-31)^2}{5} + \frac{(-29)^2}{8} + \frac{(19)^2}{6} - CF$$

$$= \frac{96659}{120} - 361.8846$$

$$= 805.4917 - 361.8846$$

$$SST = 443.6071$$

Sum of squares due to Errors

$$SSE = TSS - SST$$

$$= 1957.1154$$

$$SSE = 1513.5083$$

ANOVA    I way table:-

| Source of Variation      | degree of freedom  | Sum of squares    | mean sum of squares                          | Variance ratio                   |
|--------------------------|--------------------|-------------------|----------------------------------------------|----------------------------------|
| Treatments (b/w batches) | $k-1 = 3$          | $SST = 443.6071$  | $MST = \frac{SST}{k-1} = \frac{443.6071}{3}$ | $F = \frac{MST}{MSE} = 147.8690$ |
| Error                    | $N-k$<br>$26-4=22$ | $SSE = 1513.5083$ | $147.8690$                                   | $\frac{68.7952}{2.1494}$         |
| Total                    | $N-1$<br>$26-1=25$ |                   | $68.7952$                                    |                                  |

\* The variance ratio the numerator should be always greater than denominator  
If the condition then reverse the fraction

Table value of F at 5% level of significance  
at (3, 22) degrees of freedom is 3.0491

calculated value is  $<$  Table value, <sup>than</sup> accept  $H_0$ .

problem:- 4 sales men were posted in different areas by a company. The no. of units of commodity X sold by them are as follows.  
Is there significant difference in the performance of salesmen?

|   |    |    |    |    |
|---|----|----|----|----|
| A | 20 | 23 | 28 | 29 |
| B | 25 | 32 | 30 | 21 |
| C | 23 | 28 | 35 | 18 |
| D | 15 | 21 | 19 | 25 |

$$H_0 \Rightarrow \mu_1 = \mu_2 = \mu_3 = \mu_4$$

$$H_1 \Rightarrow \mu_1 \neq \mu_2 \neq \mu_3 \neq \mu_4$$

Coding method:- Subtract each number from 28

|   | $x_{ij} - T_i$ | $x_{ij}^2$ |
|---|----------------|------------|
| A | 8 5 0 -1       | 12 90      |
| B | 3 -4 -2 7      | 4 78       |
| C | 5 0 -7 10      | 8 174      |
| D | 13 7 9 3       | 32 308     |
|   | G = 56         | RSS = 650  |

$$N = 16, RSS = 650, G = \sum \sum x_{ij} = 650$$

(79)

$$\text{Correction factor} = \frac{G^2}{N} = \frac{650^2}{16} - \frac{56^2}{16}$$

$$= \frac{105625}{16} - 196$$

$$\text{Correction factor} = 6601.5625$$

$$\text{Total sum of squares} = TSS = \text{PSS} - CF$$

$$= 650 - 6601.5625$$

$$= 650 - 196$$

$$= 454$$

Sum of squares due to treatments

$$= SST = \sum_{i=1}^4 \left( \frac{T_i^2}{n_i} \right) - CF$$

$$\frac{(12)^2}{4} + \frac{(4)^2}{4} + \frac{(8)^2}{4} + \frac{(32)^2}{4} - 454$$

$$SST = 312 - 454$$

$$SST = -142 \quad (116)$$

Sum of squares due to Errors:-

$$SSE = TSS - SST$$

$$= 454 - (-142) \quad (116 - 116)$$

$$= 454 + 142$$

$$SSE = 596 \quad (338)$$

Table value of F at 5% level of significance at  
3,12 df is 3.4903

ANOVA I way table :-

| source of variation       | degree of freedom         | sum of squares | mean of sum squares                                 | Variance ratio                |
|---------------------------|---------------------------|----------------|-----------------------------------------------------|-------------------------------|
| Treatments<br>( batches ) | $K-1 \Rightarrow 4-1 = 3$ | $SST = -142$   | $MST = \frac{SST}{k-1} = \frac{-142}{3} = 47.3333$  | $F = \frac{67}{47.3333}$      |
| error                     | $N-k$<br>$16-4=12$        | $SS_E = 596$   | $MSE = \frac{SS_E}{N-k} = \frac{596}{12} = 49.6667$ | $F = \frac{47.3333}{49.6667}$ |
| Total                     | $N-1$<br>$16-1=15$        | TSS            |                                                     | $F = +$<br>$-1.0493$          |

$-1.0493 < 3.4903$  then accept  $H_0$ .

3/12/2020

ANOVA II way :-

Rows  $\rightarrow k$ columns  $\rightarrow h$ Formulae:-  $TSS = RSS - CF$ 

$$CF = \frac{G^2}{N}$$

$$SST = \frac{1}{h} \sum T_i^2 - CF \quad (h\text{-columns})$$

$$SSV = \frac{1}{k} \sum T_j^2 - CF$$

$$G = \sum_{j=1}^k y_{ij}, \quad T_i = \sum_{j=1}^k y_{ij}$$

$$T_j = \sum_{i=1}^h y_{ij}$$

$$SSE = TSS - SST - SSV$$

ANOVA II way Table:-

| Source of Variance | d.f          | Sum of Squares | Mean sum of squares            | Variance Ratio                                        |
|--------------------|--------------|----------------|--------------------------------|-------------------------------------------------------|
| Treatments         | $k-1$        | SST            | $MST = \frac{SST}{k-1}$        | $F = \frac{MST}{MSE}$                                 |
| Varieties          | $h-1$        | SSV            | $MSV = \frac{SSV}{h-1}$        | $\sim F_{(k-1), h-1}$                                 |
| Errors             | $(k-1)(h-1)$ | SSE            | $MSE = \frac{SSE}{(k-1)(h-1)}$ | $F = \frac{MSV}{MSE}$<br>$\sim F_{(h-1), (k-1)(h-1)}$ |
| Total              | $hk-1$       |                |                                |                                                       |

numerator is greater than denominator

(82)

problem:- A tea company appoints 4 salesmen A, B, C and D & observes that sales in three seasons - summer, winter & monsoon of the figures (in lakhs) are given in the following table.

| seasons  | A  | B  | C  | D  | seasons total |
|----------|----|----|----|----|---------------|
| summer   | 36 | 36 | 21 | 35 | 128           |
| winter   | 28 | 29 | 31 | 32 | 120           |
| monsoon  | 26 | 28 | 29 | 29 | 112           |
| salesmen | 90 | 93 | 81 | 96 | 360           |
| Total    |    |    |    |    |               |

$$H_0: \mu_1 = \mu_2 = \mu_3$$

$$H_{11}: \mu_1 \neq \mu_2 \neq \mu_3$$

$$H_{02}: \mu_1' = \mu_2' = \mu_3' = \mu_4'$$

$$H_{12}: \mu_1' \neq \mu_2' \neq \mu_3' \neq \mu_4'$$

Coding method:- Subtract each value from 29

|                                        | T <sub>i</sub>                           | T <sub>i</sub> <sup>2</sup> |
|----------------------------------------|------------------------------------------|-----------------------------|
| -7 -7 8 -6                             | -12                                      | 144                         |
| 1 0 -2 -3                              | -4                                       | 16                          |
| 3 1 0 0                                | 4                                        | 16                          |
| T <sub>i</sub> -3 -6 6 -9              | G = -12                                  | $\sum T_i = 176$            |
| T <sub>j</sub> <sup>2</sup> 9 36 36 81 | <del><math>\sum T_j^2 = 162</math></del> | $\sum T_j^2 = 162$          |

$$N = 12$$

$$\bar{G} = -12$$

$$CF = \frac{\bar{G}^2}{N} = \frac{(-12)^2}{12} = 12$$

$$RSS = (-7)^2 + (-7)^2 + 8^2 + (-6)^2 + (1)^2 + 0^2 + (-2)^2 + (-3)^2 + (3)^2 + 1^2 + 0^2 + 0^2 \\ = 222$$

$$SST = \frac{1}{h} \sum T_{ij}^2 - CF$$

$$= \frac{1}{4} (176) - 12$$

$$SST = 32$$

$$SSR = \frac{1}{k} \sum T_{ij}^2 - CF$$

$$= \frac{1}{3} (162) - 12$$

$$= 42$$

$$TSS = RSS - CF$$

$$= 222 - 12$$

$$TSS = 210$$

$$SSE = TSS - SST - SSV$$

$$= 210 - 32 - 42$$

$$SSE = 136$$

ANOVA II-way table

(84)

| source of variance | d.f                          | sum of squares | mean sum of squares            | F-ratio                                    |
|--------------------|------------------------------|----------------|--------------------------------|--------------------------------------------|
| Treatments         | $3-1=2$                      | $32 = SST$     | $MST = \frac{32}{2} = 16$      | $F = \frac{MST}{MSE}$                      |
| Varieties          | $4-1=3$                      | $SSV = 42$     | $MSV = \frac{42}{3} = 14$      | $= \frac{22.667}{16}$                      |
| Error              | $(2-1)(4-1)$<br>$2*3=6$      | $SSE = 136$    | $MSE = \frac{136}{22.667} = 6$ | $F = \frac{MSV}{MSE} = 1.4167$             |
| Total              | $hk-1$<br>$= 12-1$<br>$= 11$ |                |                                | $F = \frac{MSV}{MSE} = 22.667/14 = 1.6191$ |

since the F ratios are  $< 1$  accept  $H_0$ .

$1.4167 < 19.33$  calculated value  $<$  Table value

Accept  $H_0$

$1.619 < 8.94$  calculated value  $<$  Table value

Accept  $H_0$

| df  |                                                          |
|-----|----------------------------------------------------------|
| T 2 | F at $(2, 6)$ Numerator $>$ denominator                  |
| V 3 | F at $(3, 6)$                                            |
| E 6 | F at $(6, 2)$ Numerator $<$ denominator<br>F at $(6, 3)$ |

There is no significant difference b/w the  
Seasons

There is no significant difference b/w  
the performance of  
salesmen.

Problem:- The following table represents the sales (Rs 1000) per month of 3 brands of soaps allocated among three cities

| Brands | cities |    |    |
|--------|--------|----|----|
|        | A      | B  | C  |
| I      | 12     | 48 | 30 |
| II     | 42     | 54 | 57 |
| III    | 9      | 42 | 21 |

Test whether ① mean sales of three brands are equal ② mean sales of soaps in each city are equal

| Brands               | A                                 | B   | C   | T <sub>i</sub>         | T <sub>i,2</sub>       |
|----------------------|-----------------------------------|-----|-----|------------------------|------------------------|
| I                    | 30                                | -6  | 12  | 36                     | 1296                   |
| II                   | 0                                 | -12 | -15 | -27                    | 729                    |
| III                  | 33                                | 0   | 21  | 54                     | 2916                   |
| sales T <sub>j</sub> | 63                                | -18 | 18  | G=63                   | $\sum T_{ij}^2 = 4941$ |
| total                | T <sub>j</sub> <sup>2</sup> =3969 | 324 | 324 | $\sum T_{ij}^2 = 4617$ |                        |

coding method:- Subtract with 42 (each number)

$$H_01 \Rightarrow u_1 = u_2 = u_3$$

$$H_{11} \Rightarrow u_1 \neq u_2 \neq u_3$$

$$H_{02} \Rightarrow u_1' = u_2' = u_3'$$

$$H_{12} \Rightarrow u_1' \neq u_2' \neq u_3'$$

$$N = 9$$

$$G = 63$$

$$CF = \frac{G^2}{N} = \frac{(63)^2}{9} = 441$$

(86)

$$RSS = (30)^2 + (-6)^2 + (12)^2 + (0)^2 + (-12)^2 \\ + (-15)^2 + (33)^2 + (0)^2 + (21)^2$$

$$RSS = 2979$$

$$SST = \frac{1}{n} \sum T_{ij}^2 - CF$$

rows  $\Rightarrow 3$   
columns  $\Rightarrow 3$

$$= \frac{1}{3} (4941) - 441$$

$$= 1647 - 441$$

$$SST = 1206$$

$$SSV = \frac{1}{3} \sum T_{ij}^2 - CF$$

$$= \frac{1}{3} (4617) - 441$$

$$= 1539 - 441$$

$$SSV = 1098$$

$$TSS = RSS - CF$$

$$TSS = 2538$$

$$SSE = TSS - SST - SSV$$

$$= 2538 - 1206 - 1098$$

$$SSE = 234$$

## ANOVA II way table

| Source of Variance | d.f               | Sum of Squares | mean sum of squares          | F. ratio                                    |
|--------------------|-------------------|----------------|------------------------------|---------------------------------------------|
| Treatments         | $3-1=2$           | $SST = 1206$   | $MST = \frac{1206}{2} = 603$ | $\frac{MST}{MSE} = \frac{603}{54} = 11$     |
| Varieties          | $3-1=2$           | $SSV = 1098$   | $MSV = \frac{1098}{2} = 549$ | $5.1538$                                    |
| Error              | $2*2 = 4$         | $SSE = 234$    | $MSE = \frac{234}{4} = 58.5$ | $\frac{MSV}{MSE} = \frac{549}{58.5} = 9.38$ |
| Total              | $hk-1$<br>$9-1=8$ |                |                              |                                             |

$$(2,4) \Rightarrow 19.25$$

$$(2,4) \Rightarrow 19.25$$

Tab val of F at

(2,4) df is at 5%. los is 6.94

cal val > Tab val

Reject  $H_0_1$

Tab val of F at

(2,4) df at 5%. los

is 8.94.

cal val > Tab val

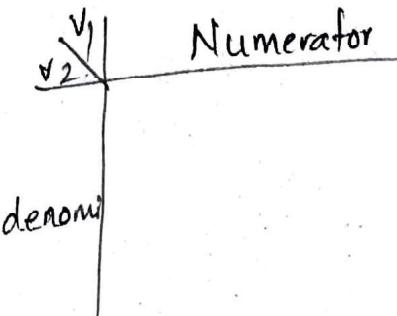
Reject  $H_0_2$

Inference:- The mean sales of three brands  
significantly differ

② The mean sales of soap in each city are not equal

$$\begin{matrix} v_1 \leftarrow D \\ v_2 \rightarrow v_2 \end{matrix}$$

(2,4)

7/12/2020

- ① A trucking company wishes to test the average life of each of the four brands of tyres the company uses all brands on randomly selected trucks. The records showing the lifes (1000 of miles) of tyres are given in the table.

Test the hypothesis that the average life of each brand of tyres is the same.  
Assume  $\alpha = 0.01$  (level of significance (1%, we have to see))

| Brand I | Brand II  | Brand III | Brand IV |
|---------|-----------|-----------|----------|
| 20      | 19        | 21        | 15       |
| 23      | 15        | 19        | 17       |
| 18      | 17        | 20        | 16       |
| 17      | 20        | 17        | 18       |
| —       | <u>16</u> | 16        | —        |

$$H_0 \Rightarrow \mu_1 = \mu_2 = \mu_3 = \mu_4$$

$$H_1 \Rightarrow \mu_1 \neq \mu_2 \neq \mu_3 \neq \mu_4$$

|           |    |    |    |    |    | T <sub>i</sub> |
|-----------|----|----|----|----|----|----------------|
|           | 20 | 23 | 18 | 17 | -  | 78             |
| Brand I   |    |    |    |    |    | 87             |
| Brand II  | 19 | 15 | 17 | 20 | 16 |                |
| Brand III | 21 | 19 | 20 | 17 | 16 | 93             |
| Brand IV  | 15 | 17 | 16 | 18 | -  |                |

coding method: - subtract each no. from 17

|  |    |    |    |    |   | T <sub>i</sub> | x <sub>ij</sub> <sup>2</sup> |
|--|----|----|----|----|---|----------------|------------------------------|
|  | -3 | -6 | -1 | 0  | - | -10            | 46                           |
|  | -2 | 2  | 0  | -3 | 1 | -2             | 18                           |
|  | -4 | -2 | -3 | 0  | 1 | -8             | 30                           |
|  | 2  | 0  | 1  | -1 | - | 2              | 6                            |
|  |    |    |    |    |   | G = 18         | RSS = 100                    |

$$N = 18$$

$$G = -18$$

$$RSS = 100$$

$$CF = \frac{G^2}{N} = \frac{(-18)^2}{18} = 18$$

$$TSS = RSS - CF$$

$$= 100 - 18$$

$$TSS = 82$$

$$SST = \left(\frac{100}{4}\right) + \frac{4}{5} + \frac{64}{5} + \frac{4}{4} \quad (90)$$

$$= \frac{198}{5} - 18$$

$$= \frac{108}{5}$$

$$SST = 21.6$$

$$SSE = TSS - SST$$

$$= 82 - 21.6$$

$$SSE = 60.4$$

| source of variation | d.f                | sum of squares | mean sum of squares               | Variance ratio         |
|---------------------|--------------------|----------------|-----------------------------------|------------------------|
| Treatments (Brands) | $k-1$<br>$4-1=3$   | $SST = 21.6$   | $MST = 21.6$<br>$= \frac{3}{7.2}$ | $F = \frac{MST}{MSE}$  |
| Error               | $N-K$<br>$18-4=14$ | $SSE = 60.4$   | $MSE = \frac{60.4}{14}$           | $= \frac{7.2}{4.3143}$ |
|                     | $18-1=17$          | $TSS = 82$     | $= 4.3143$                        | $= 1.6689$             |

$$(3,14) \Rightarrow 4.69556$$

Table value of F at 1% level of significance  
at (3,14) d.f is

Calculated value < Table value Accept H<sub>0</sub>.

Inference:- average life of each brand of  
tyres is the same.

H/10

(9)

- ② Setup an ANOVA table for the following per hectare yield in (100 kg) for the three varieties of wheat grown

| Varities of wheat | I | II | III | IV |  |
|-------------------|---|----|-----|----|--|
| A1                | 6 | 7  | 3   | 8  |  |
| A2                | 5 | 5  | 3   | 7  |  |
| A3                | 5 | 4  | 3   | 4. |  |

$$H_{01} \Rightarrow \mu_1 = \mu_2 = \mu_3$$

$$H_{11} \Rightarrow \mu_1 \neq \mu_2 \neq \mu_3$$

$$H_{02} \Rightarrow \mu_1^1 = \mu_2^1 = \mu_3^1 = \mu_4^1$$

$$H_{12} \Rightarrow \mu_1^1 \neq \mu_2^1 \neq \mu_3^1 \neq \mu_4^1$$

H/10

(91)

- ② Setup an ANOVA table for the following per hectare yield in (100 kg) for the three varieties of wheat grown

| varieties of wheat | I | II | III | IV |
|--------------------|---|----|-----|----|
| A1                 | 6 | 7  | 3   | 8  |
| A2                 | 5 | 5  | 3   | 7  |
| A3                 | 5 | 4  | 3   | 4  |

$$H_0_1 \Rightarrow \mu_1 = \mu_2 = \mu_3$$

$$H_{11} \Rightarrow \mu_1 \neq \mu_2 \neq \mu_3$$

$$H_0_2 \Rightarrow \mu_1^1 = \mu_2^1 = \mu_3^1 = \mu_4^1$$

$$H_{12} \Rightarrow \mu_1^1 + \mu_2^1 + \mu_3^1 \neq \mu_4^1$$

| varieties of wheat | I  | II | III | IV | $T_i$ | $T_{ij}^2$ |
|--------------------|----|----|-----|----|-------|------------|
| A1                 | -1 | -2 | 2   | -3 | -4    | 1816       |
| A2                 | 0  | 0  | 2   | -2 | 0     | 80         |
| A3                 | 0  | 1  | 2   | 1  | 4     | 616        |
| $T_j$              | -1 | -1 | 6   | -4 | $G=0$ | 32 = RSS   |

$$N=12$$

$$SST = \frac{16}{4} + \frac{16}{4}$$

$$RSS = 32$$

$$SSF = 8$$

$$TSS = 32$$

$$SSE = TSS - SST = SSV$$

$$SSV = \frac{1}{3} (5.6) - 0 = 18 - 32 - 8 - 8 \\ = 24.6$$

| source of variation | d.f              | sum of squares | mean sum of squares      | Variance Ratio                      |
|---------------------|------------------|----------------|--------------------------|-------------------------------------|
| treatments          | $3-1 = 2$        | $SSF = 8$      | $MST = \frac{8}{2} = 4$  | $\frac{MST}{MSE} = \frac{4}{1} = 4$ |
| varieties           | $4-1 = 3$        | $SSV = 18$     | $MSV = \frac{18}{3} = 6$ |                                     |
| error               | $2 \times 3 = 6$ | $SSE = 6$      | $MSE = \frac{6}{6} = 1$  | $\frac{MSV}{MSE} = \frac{6}{1} = 6$ |
| total               | $12-1 = 11$      |                |                          |                                     |

(2, 6) (3, 6)

Table value of F at 5%. los at (2, 6)df is

5.14 calval < tabval

Accept H<sub>0</sub>1

Table value of F at 5%. los at (3, 6)df is

4.71 calval > Tab val Reject H<sub>0</sub>2

- Inferences - ① There is no significant difference b/w the three varieties of wheat  
 ② There is significant difference b/w the four yields of land.

11/12/2020

(93)

Friday

Estimation :- The theory of statistical inference can be divided into two major areas

① Estimation of parameters

② Testing of hypothesis

Parameter :- A statistical constant derived from the population values

Ex:- population mean,  $\mu$

population standard deviation,  $\sigma$

Statistic :- A statistical constant derived from the sample values is called a 'statistic'

Estimation :- procedure it is a procedure of estimating a population parameter by using sample information (or) observation

1) Point Estimation

2) Interval Estimation

Point Estimate :- A point estimate of some population parameter  $\theta$  is a single numerical value

Point Estimator :- A point estimator is a statistic for estimating a population parameter  $\theta$  f is denoted by  $\hat{\theta}$

$\bar{x} = 25 \rightarrow$  Estimate

$\hat{\theta} \leftarrow \bar{x} = 25.2 \rightarrow$  point estimator  
sample mean

Properties of Estimators:-

An estimator is said to be a good estimator if it is

- |                                                               |                                                                             |
|---------------------------------------------------------------|-----------------------------------------------------------------------------|
| 1. unbiased<br>2. consistent<br>3. Efficient<br>4. sufficient | (on<br>1. unbiasedness<br>2. Consistency<br>3. Efficiency<br>4. sufficiency |
|---------------------------------------------------------------|-----------------------------------------------------------------------------|

Unbiased Estimator:- A statistic  $\hat{\theta}$  (say) is said to be an unbiased estimator (of) its value as an unbiased estimate if & only if the mean of the sampling distribution of the estimator equals  $\theta$  i.e;

$$\boxed{\text{Mean of } \hat{\theta} = E(\hat{\theta}) = \theta}$$

$$\boxed{E(\hat{\theta}) = \theta}$$

$$\begin{matrix} R \cdot V \\ \nearrow \\ \text{discr} \quad \text{cont} \\ E(x) = \sum x p(x) / \int f(x) dx \end{matrix}$$

Sampling distribution:- Sample size =  $n = 2$   $\Rightarrow$   $n = N$

$N=5$

$n=2$

Sample no =  $N$

1

2

3

4

5

(q8)

If the estimator is not unbiased, the difference  $E(\hat{\theta}) - \theta$  is called the bias of the estimate  $\hat{\theta}$ .  
 When the estimate estimator is unbiased,

$$\boxed{E(\hat{\theta}) - \theta = 0}$$

i.e., the bias is zero.

positive bias:- If  $E(\hat{\theta}) > \theta$ , then  $\hat{\theta}$  is called positively biased.

negative bias:- If  $E(\hat{\theta}) < \theta$ , then  $\hat{\theta}$  is called negatively biased.

①

problem:- If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two unbiased estimators of  $\theta$ , then  $w_1\hat{\theta}_1 + w_2\hat{\theta}_2$  is an unbiased estimator of  $\theta$  when  $w_1 + w_2 = 1$ .

Soln:- given that

$\hat{\theta}_1$  &  $\hat{\theta}_2$  are unbiased estimators of  $\theta$

$$\text{i.e., } E(\hat{\theta}_1) = \theta \rightarrow ①$$

$$E(\hat{\theta}_2) = \theta \rightarrow ②$$

$$\text{Now, } E(w_1\hat{\theta}_1 + w_2\hat{\theta}_2) = w_1 E(\hat{\theta}_1) + w_2 E(\hat{\theta}_2)$$

$$= w_1\theta + w_2\theta \quad (\text{from eqn } ① + ②)$$

$$= (w_1 + w_2)\theta \quad (w_1 + w_2 = 1)$$

$$E(w_1\hat{\theta}_1 + w_2\hat{\theta}_2) = \theta$$

$\therefore w_1\hat{\theta}_1 + w_2\hat{\theta}_2$  is an unbiased estimator of  $\theta$ .

(96)

Theorem:— Let  $x_1, x_2, \dots, x_n$  be a random sample drawn from a given population with the mean  $\mu$  & Variance  $\sigma^2$  show that the sample mean is an unbiased estimator of population mean  $\mu$ .  
 i.e.,  $E(\bar{x}) = \mu$ .

Proof:—  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Taking Expectations on both sides

$$E(\bar{x}) = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} E\left(\sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} E(x_1 + x_2 + \dots + x_n)$$

$$= \frac{1}{n} \{E(x_1) + E(x_2) + \dots + E(x_n)\}$$

$$= \frac{1}{n} \{ \mu + \mu + \dots + \mu \} \quad [\because E(x_i) = \mu]$$

$$= \frac{1}{n} [n\mu]$$

$E(\bar{x}) = \mu$

Theorem:— for a random sample of size  $n$ , taken from a finite population  $x_1, x_2, \dots, x_n$  taken from a finite population  $s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$  is not an unbiased estimator of  $\sigma^2$  but  $\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$  is an unbiased estimator of  $\sigma^2$ .

$$s^2 = \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2$$

MVUE (minimum variance unbiased estimator)  
(or) most efficient unbiased Estimator

If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are the two unbiased estimators of  $\theta$  if and if  $\sigma_1^2$  and  $\sigma_2^2$  are variances of their sampling distributions and  $\sigma_1^2 < \sigma_2^2$ ,  $\hat{\theta}_1$  is said to be more unbiased estimator of  $\theta$ .

Note: If  $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$  are unbiased estimators of the parameter  $\theta$

$\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2$  are variances of their sampling distributions then  $\hat{\theta}_1$  is the most efficient estimator of  $\theta$  if

$$0 < \sigma_1^2 < \sigma_2^2 < \dots < \sigma_k^2$$

e = efficiency of an estimator

$e = \frac{\text{variance of } \hat{\theta}_1}{\text{variance of } \hat{\theta}_i}$  where  $i = 2, 3, \dots, k$

$\text{variance of } \hat{\theta}_i$  ( $e < 1$ )

$$\boxed{e = \frac{\text{var of } \hat{\theta}_1}{\text{var of } \hat{\theta}_i}}$$

$$\boxed{e = \frac{\text{var } \hat{\theta}_1}{\text{var } \hat{\theta}_i}}$$

consistency: An estimator  $T_n = T(x_1, x_2, \dots, x_n)$  based on a random sample of size  $n$ , is said to be a consistent estimator of  $R(\theta)$ ,  $\theta \in \Theta$ ,  $\Theta \rightarrow$  parameter space ; if  $T_n$  converges to  $R(\theta)$

i.e., if  $T_n \xrightarrow{P} R(\theta)$  as  $n \rightarrow \infty$

(98)

In other words  $T_n$  is a consistent estimator of  $\tau(\theta)$  if for every  $\epsilon > 0$ ,  $\exists \eta > 0$ , there exists a positive integer  $n > m(\epsilon, \eta)$  such that  $P\{|T_n - \tau(\theta)| < \epsilon\} \rightarrow 1$  as  $n \rightarrow \infty \Rightarrow P\{|T_n - \tau(\theta)| < \epsilon\} \geq 1 - \eta \forall n \geq n_0$  where  $n_0$  is some very large value  $n$ .

$f, \theta \rightarrow$  parameter  
 $T_n \rightarrow$  estimator

18/12/2020

sufficiency :- An estimator is said to be sufficient for a parameter, if it contains all the information in the sample regarding the parameter.

If  $T = t(x_1, x_2, \dots, x_n)$  is an estimator of a parameter  $\theta$  based on a sample  $x_1, x_2, \dots, x_n$  of size  $n$  from the population with density  $f(x, \theta)$  such that the conditional distribution of  $x_1, x_2, \dots, x_n$  given  $T$ , is independent of  $\theta$ , then  $T$  is sufficient estimator for  $\theta$ .

## Methods of Estimation:-

- ① maximum likelihood method
- ② method of moments
- ③ method of minimum chi-square
- ④ method of minimum variance
- ⑤ method of least squares
- ⑥ method of inverse probability

Maximum likelihood Estimation:- It was formulated by C.F Gauss as a general method of estimation was introduced by professor by R.A fisher

likelihood fn:- Let  $x_1, x_2, \dots, x_n$  be a random sample of size  $n$  from a population with density fn  $f(x, \theta)$  Then the likelihood fn of the sample values  $x_1, x_2, \dots, x_n$  usually denoted by  $L = L(\theta)$  is their joint density fn

$$L = f(x_1, \theta) f(x_2, \theta) \dots f(x_n, \theta)$$

$$= \prod_{i=1}^n (f(x_i, \theta))$$

$L$  gives the relative likelihood that the random variables assume a particular set of values  $x_1, x_2, \dots, x_n$ . For a given sample  $x_1, x_2, \dots, x_n$   $L$  becomes a fn of the variable  $\theta$ , the parameter

(60)

The principle of maximum likelihood estimation consists in finding a estimator for the unknown parameter  $\theta = (\theta_1, \theta_2, \dots, \theta_k)$  which maximises the likelihood function  $L(\theta)$  for variation in parameter

$L(\theta)$  if there exists a fn  $\hat{\theta} = \hat{\theta}(x_1, x_2, \dots, x_n)$

of the sample values which maximises  $L$  for variations in  $\theta$ , then  $\hat{\theta}$  is to be taken as an estimator of  $\theta$

$\hat{\theta}$  is usually called maximum likelihood estimator

thus  $\hat{\theta}$  is the sm<sup>n</sup> of

$$\frac{\partial L}{\partial \theta} = 0 + \frac{\partial^2 L}{\partial \theta^2} < 0$$

$\hookrightarrow \textcircled{1}$      $\frac{\partial^2 L}{\partial \theta^2} < 0 \hookrightarrow \textcircled{2}$

since  $L > 0$  &  $\log L$  is nondecreasing fn of  $L$ ,  $L$  &  $\log L$  attain their extreme values (maximal minima) at the same value of  $\theta$

(1) can be rewritten as

$$\frac{1}{L} \frac{\partial L}{\partial \theta} = 0$$

This is a convenient form from practical view

If  $\theta$  is a vector valued parameter, then  $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)$  is given by the sm<sup>n</sup> of simultaneous eqns

$$\frac{\partial}{\partial \theta_i} \log L = \frac{\partial}{\partial \theta_i} \log L(\theta_1, \theta_2, \dots, \theta_k) = 0$$

$\hookrightarrow \textcircled{3}$

$$i = 1, 2, 3, \dots, k$$

(16)

eqns ①, ② & ③ are usually referred to  
as likelihood eqn for estimating the  
parameters

(91)

- H/10
- ② Setup an ANOVA table for the following per hectare yield in (100kg) for the three varieties of wheat grown

| varieties of wheat | I | II | III | IV |  |
|--------------------|---|----|-----|----|--|
| A1                 | 6 | 7  | 3   | 8  |  |
| A2                 | 5 | 5  | 3   | 7  |  |
| A3                 | 5 | 4  | 3   | 4  |  |

$$H_0_1 \Rightarrow \mu_1 = \mu_2 = \mu_3$$

$$H_{11} \Rightarrow \mu_1 \neq \mu_2 \neq \mu_3$$

$$H_{02} \Rightarrow \mu_1^1 = \mu_2^1 = \mu_3^1 = \mu_4^1$$

$$H_{12} \Rightarrow \mu_1^1 \neq \mu_2^1 \neq \mu_3^1 \neq \mu_4^1$$

| varieties of wheat | I  | II | III | IV | T <sub>i</sub> | T <sub>ij</sub> <sup>2</sup> |
|--------------------|----|----|-----|----|----------------|------------------------------|
| A1                 | -1 | -2 | 2   | -3 | -4             | 18/16                        |
| A2                 | 0  | 0  | 2   | -2 | 0              | 8/0                          |
| A3                 | 0  | 1  | 2   | 1  | 4              | 6/16                         |
| T <sub>j</sub>     | -1 | -1 | 6   | -4 | G=0            | 32 = RSS                     |

$$N = T_j = 1 + 1 + 36 = 38$$

$$SST = \frac{18}{4} + \frac{16}{4}$$

$$RSS = 32$$

$$CF = 0$$

$$SST = 8$$

$$TSS = 32$$

$$SSE = TSS - SST - SSV$$

$$SSV = \frac{1}{3} (5G) - 0 = 18 - 32 - 8 = 8$$

$$= 24/6$$

| source of variation | d.f              | sum of squares | mean sum of squares      | Variance Ratio                      |
|---------------------|------------------|----------------|--------------------------|-------------------------------------|
| Treatments          | $3-1 = 2$        | $SST = 8$      | $MST = \frac{8}{2} = 4$  | $\frac{MST}{MSE} = \frac{4}{1} = 4$ |
| varieties           | $4-1 = 3$        | $SSV = 18$     | $MSV = \frac{18}{3} = 6$ |                                     |
| Error               | $2 \times 3 = 6$ | $SSE = 6$      | $MSE = \frac{6}{6} = 1$  | $\frac{MSV}{MSE} = \frac{6}{1} = 6$ |
| total               | $12-1=11$        |                |                          |                                     |

(2,6) (3,6)

Table value of F at 5%. los at (2,6)df is  
5.14 calval < tabval

Accept  $H_0$

Table value of F at 5%. los at (3,6)df is  
4.71 calval > Tab Val Reject  $H_0$

- Inferences- ① There is no significant difference b/w the three varieties of wheat  
 ② There is significant difference b/w the four yields of land.

11/12/2020

Friday

Estimation:- The theory statistical inference  
can be divided into two major areas

① Estimation of parameters

② Testing of hypothesis

Parameter:- A statistical constant derived from  
the population values

Ex:- population mean,  $\mu$

population standard deviation,

statistic:- A statistical constant derived from  
the sample values is called a 'statistics'

Estimation:- procedure it is a procedure of  
estimating a population parameter by using  
sample information (or) observation

1) Point Estimation

2) Interval Estimation

Point Estimate:- A point estimate of some  
population parameter  $\theta$  is a single  
numerical value

Point Estimator:- A point estimator is a  
statistic for estimating a population  
parameter  $\theta$  if is denoted by  $\hat{\theta}$

$\mu = 25 \rightarrow$  Estimate

$\hat{\theta} \leftarrow \bar{x} = 25.2 \rightarrow$  point estimator  
Sample mean

Properties of Estimators:-

A statistic  $\hat{\theta}$  is said to be a good estimator if it is

1. unbiased
2. consistent
3. efficient
4. sufficient.

(on)      1. Unbiasedness  
               2. Consistency  
               3. Efficiency  
               4. Sufficiency

Unbiased Estimator:- A statistic  $\hat{\theta}$  is said to be an unbiased estimator (or) its value is an unbiased estimate if & only if the mean of the sampling distribution of the estimator equals  $\theta$  i.e;

$$\text{Mean of } \hat{\theta} = E(\hat{\theta}) = \theta$$

$$E(\hat{\theta}) = \theta$$

$$\begin{aligned} R.V \\ \text{discrete cont} \\ E(x) = \sum p(x) \\ E(n) = \int x f(x) dx \end{aligned}$$

Sampling distribution:- Sample size =  $2^m n$

Sample no =  $N$

$N=5$

1

$n=2$

2

3

4

5

If the estimator is not unbiased, the difference  $E(\hat{\theta}) - \theta$  is called the bias of the Estimate  $\hat{\theta}$   
 when the estimate estimator is unbiased,

$$\boxed{E(\hat{\theta}) - \theta = 0}$$

i.e., the bias is zero.  
 positive bias:- If  $E(\hat{\theta}) > \theta$ , then  $\hat{\theta}$  is called positively biased

negative bias:- If  $E(\hat{\theta}) < \theta$ , then  $\hat{\theta}$  is called negatively biased.

① problem:- If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are two unbiased estimators of  $\theta$ , then  $w_1\hat{\theta}_1 + w_2\hat{\theta}_2$  is an unbiased estimator of  $\theta$  when  $w_1 + w_2 = 1$

Soln:- given that  
 $\hat{\theta}_1$  &  $\hat{\theta}_2$  are unbiased estimators of  $\theta$   
 i.e.,  $E(\hat{\theta}_1) = \theta \rightarrow ①$

$$E(\hat{\theta}_2) = \theta \rightarrow ②$$

$$\text{Now, } E(w_1\hat{\theta}_1 + w_2\hat{\theta}_2) = w_1 E(\hat{\theta}_1) + w_2 E(\hat{\theta}_2)$$

$$= w_1\theta + w_2\theta \quad (\text{from eq } ① + ②)$$

$$= (w_1 + w_2)\theta \quad (w_1 + w_2 = 1)$$

$$E(w_1\hat{\theta}_1 + w_2\hat{\theta}_2) = \theta$$

$\therefore w_1\hat{\theta}_1 + w_2\hat{\theta}_2$  is an unbiased estimator of  $\theta$

Theorem:- Let  $x_1, x_2, \dots, x_n$  be a random sample drawn from a given population with the mean  $\mu$  & Variance  $\sigma^2$  show that the sample mean is an unbiased estimator of population mean  $\mu$ :

$$\text{i.e., } E(\bar{x}) = \mu.$$

Proof:-  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Taking Expectations on both sides:

$$E(\bar{x}) = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} E\left(\sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} E(x_1 + x_2 + \dots + x_n)$$

$$= \frac{1}{n} \{E(x_1) + E(x_2) + \dots + E(x_n)\}$$

$$= \frac{1}{n} \{ \mu + \mu + \dots + \mu \} \quad [ \because E(x_i) = \mu ]$$

$$= \frac{1}{n} [n\mu]$$

$E(\bar{x}) = \mu$

Theorem:- for a random sample of size  $n$ , taken from a finite population  $x_1, x_2, \dots, x_n$

$s^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$  is not an unbiased

estimator of  $\sigma^2$  but  $\frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2$  is an unbiased estimator of  $\sigma^2$

$$s^2 = \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2$$

MVUE (minimum variance unbiased estimate),  
(or) most efficient unbiased Estimator

If  $\hat{\theta}_1$  and  $\hat{\theta}_2$  are the two unbiased estimators of  $\theta$  if and if  $\sigma_1^2$  and  $\sigma_2^2$  are variances of their sampling distributions and  $\sigma_1^2 < \sigma_2^2$ ,  $\hat{\theta}_1$  is said to be more unbiased estimator of  $\theta$ .

Note: If  $\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k$  are unbiased estimators of the parameter  $\theta$ ,

$\sigma_1^2, \sigma_2^2, \dots, \sigma_k^2$  are variances of their sampling distributions then  $\hat{\theta}_1$  is the most efficient estimator of  $\theta$  if

$$\theta \sigma_1^2 < \sigma_2^2 < \dots < \sigma_k^2$$

e = efficiency of an estimator

$e = \frac{\text{Variance } \hat{\theta}_1}{\text{Variance } \hat{\theta}_i}$  where  $i = 2, 3, \dots, k$

$$e = \frac{\text{Var } \hat{\theta}_1}{\text{Var } \hat{\theta}_i} \quad (e < 1)$$

$$e = \frac{\text{Var } \hat{\theta}_1}{\text{Var } \hat{\theta}_i}$$

Consistency: An estimator  $T_n = T(x_1, x_2, \dots, x_n)$  based on a random sample of size  $n$ , is said to be a consistent estimator of  $r(\theta)$ ,  $\theta \in \Theta$ ,  $\Theta \rightarrow$  parameter space ; If  $T_n$  converges to  $r(\theta)$

i.e., if  $T_n \xrightarrow{P} r(\theta)$  as  $n \rightarrow \infty$

In other words  $T_n$  is a consistent estimator of  $\gamma(0)$  if for every  $\epsilon > 0$ ,  $\exists \eta > 0$ , there exists a positive integer  $n \geq m(\epsilon, \eta)$  such that  $P\{|T_n - \gamma(0)| < \epsilon\} \rightarrow 1$  as  $n \rightarrow \infty \Rightarrow P\{|T_n - \gamma(0)| < \epsilon\} \geq 1 - \eta \forall n \geq n_0$  where  $n_0$  is some very large value  $n$ .

$\gamma(0) \rightarrow$  parameter

$T_n \rightarrow$  estimator