

Module - 3 : Statistical Methods

Correlation :- The changes in one variable are associated (or) followed by changes in the other is called correlation.

The degree of relation b/w 2 variables x & y is measured by a single no. 'r' is called correlation coefficient.

* The fund to find out the extent to which 2 variables are related to each other.

If 2 variables x & y are related in such a way that increase (or) decrease in one of them corresponds to increase or decrease in the other. Then variables are positively correlated.

Also, if increase or decrease in one of them corresponds to decrease or increase in the other, the variables are said to be negatively correlated.

x & y

20 40

30 30

40 22

60 15

80 10

$x + y$

50

40

30

20

$x - y$

20 10

30 15

40 22

60 30

80 40

$(x + y)$

50

40

30

20

(+vely correlated)

(-vely correlated)

If there is no relationship indicated b/w the variables, they are said to be

Ex:- 1) Income and expenditure (highly correlated)

Height & weight

2) Price & Demand (- very correlated)

① \$9.00 per kg after 4th or 5th year

3) Two winds being tested simultaneously,
(uncorrelated).

Note:- * The coefficient of correlation $r \in [-1, 1]$

* All of the situations $-1 \leq r \leq 1$

* If $r = \pm 1$, there is perfect positive or negative correlation.

* If $r = 0$, there is no linear correlation.
but a non linear correlation may exist.

Karl Pearson's Coefficient of Correlation

The correlation coefficient r is measured by the following formula

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

If $x_i = x - \bar{x}$

$y_i = y - \bar{y}$

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

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

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

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

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

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

$$\sigma_x \sqrt{n} = \sqrt{\sum x^2} \rightarrow ① \quad \sigma_y \sqrt{n} = \sqrt{\sum y^2} \rightarrow ②$$

① & ② are called generalization

Relationship between $(\bar{x})^2 \sum xy = n \sqrt{\sum x^2} \sqrt{\sum y^2}$ Correlation & -1 to +1
 i.e., $n \bar{x} \bar{y} = \sqrt{\sum x^2} \sqrt{\sum y^2} \rightarrow ③$ if \bar{x} & \bar{y} is High

∴ $\text{Coefficient of Correlation } R = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$ Lived out ④
 If we sub eq ③ in eq ④

Problems:-

i) Compute the coefficient of correlation of the following data.

x : Initial stress 2 3 4 5 6 7 in kN/mm²
 y : Tensile stress 8 10 12 11 13 14 = 57 in kN/mm²

Sln:-

	x	y	$\bar{x} = x - 5$	$\bar{y} = y - 11$	xy	x^2	y^2
1	9	8	-3	-3	6	9	4
2	8	10	-2	-1	6	4	9
3	10	12	0	1	0	0	1
4	12	11	2	0	4	4	1
5	11	13	1	2	2	1	0
6	13	14	2	3	6	4	9
7	14	12	3	1	3	9	1

$$\bar{x} = \frac{\sum x}{n} = \frac{28}{7} = 4$$

$$\bar{y} = \frac{\sum y}{n} = \frac{77}{7} = 11$$

$$\therefore r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{26}{\sqrt{28} \sqrt{28}} = \underline{\underline{0.9285}}$$

ii) Find the coefficient of correlation for the following data.

Σx	10	14	18	22	26	30
Σy	18	12	24	6	30	36

Ans:-

x	y	$x = x - \bar{x}$	$y = y - \bar{y}$	xy	x^2	y^2
10	18	-10	-3	30	100	9
14	12	-6	-9	54	36	81
18	24	-2	3	6	4	9
22	6	2	-15	30	4	225
26	30	6	9	54	36	81
30	36	10	15	150	100	225

$$\bar{x} = \frac{\sum x}{n} = \frac{120}{6} = 20$$

$$\bar{y} = \frac{\sum y}{n} = \frac{126}{6} = 21$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{252}{\sqrt{280} \sqrt{630}} = 0.6$$

Ques:- Determine r for the following data

x	50	60	70	$\bar{x} = 65$	100
y	30	40	26	$\bar{y} = 36$	81

Ans:- -0.9

Regression :- Regression is an estimation of one independent variable with respect to other.

Equation of the regression lines

► Straight line of the form $y = ax + b$ (x being the independent variable) is called the regression line of y on x .

$$y - \bar{y} = r \cdot \frac{\sum y}{\sigma_x} (x - \bar{x}) \rightarrow ①$$

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

Consider,

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

$$= \frac{\sum xy}{\sum x^2}$$

∴ eq (1) becomes

$$\boxed{y - \bar{y} = \frac{\sum xy}{\sum x^2} (x - \bar{x}) \text{ or } r = \frac{\sum xy}{\sum x^2}}$$

* Here Slope of the line of regression of y on x

* Regression Coefficient

⇒ Straight line $y = ax + b$ (y being independent variable) is called the regression line of x on y .

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

$$\Rightarrow \boxed{(x - \bar{x}) = \frac{\sum xy}{\sum y^2} (y - \bar{y}) \text{ or } x = \frac{\sum xy}{\sum y^2} y}$$

* Slope of the line of regression of x on y is

. Ratio of least square method to the geometric mean of the regression coefficients

Hence correlation coefficient is the

geometric mean of the regression coefficients

$$\sigma = \pm \sqrt{(\text{coeff of } x)(\text{coeff of } y)}$$

$$\sigma = \sqrt{\left(\frac{\partial y}{\partial x}\right) \times \left(\frac{\partial x}{\partial y}\right)}$$

Note:- The point (\bar{x}, \bar{y}) lies on regression lines.

Problems:-

- 1) Calculate (i) regression equation of x on y and y on x from the following data.
 (ii) estimate x when $y = 20$.

~~soln:-~~ $x = 10, 12, 13, 14, 18$

$y = 5, 6, 7, 8, 13$

~~soln:-~~ $\sum x = 54$, $\sum y = 43$

x	y	$x - \bar{x}$	$y - \bar{y}$	xy	x^2	y^2
10	5	-4	-3	12	16	9
12	6	-2	-2	48	4	144
13	7	-1	-1	1	1	1
14	8	3	1	3	9	1
18	13	8	5	20	64	25
				$\sum xy = 40$	$\sum x^2 = 46$	$\sum y^2 = 40$

~~soln:-~~ $\bar{x} = \frac{\sum x}{n} = \frac{54}{5} = 10.8$, $\bar{y} = \frac{\sum y}{n} = \frac{43}{5} = 8.6$

$$\text{Soln: } \bar{x} = \frac{\sum x}{n} = \frac{54}{5} = 10.8, \quad \bar{y} = \frac{\sum y}{n} = \frac{43}{5} = 8.6$$

(i)

a) Regression line of x on y

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

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

$$(x - 10.8) = \frac{40}{40} (y - 8.6)$$

$$x - 10.8 = y - 8.6$$

$$x = y - 8 + 14$$

$$\boxed{x = y + 6}$$

b) Regression line of y on x .

and we know $y = \frac{\sum xy}{\sum x^2} \cdot x$ then let's take

$$\text{this } B \text{ is } y - g = \frac{\sum xy}{\sum x^2} (x - \bar{x})$$

$$\text{total given}, \frac{40}{46} (x - 14) \text{ is no } y$$

$$(y - 8) = \frac{40}{46} (x - 14)$$

$$y - 8 = 0.8695(x - 14)$$

$$= 0.8695x - 12.1739$$

$$y = 0.8695x - 12.1739 + 8$$

$$\boxed{y = 0.8695x - 4.1739}$$

$$\text{(i) put } y = 20 \text{ in } x = y + 6$$

$$x = 20 + 6$$

$$\underline{x = 26}$$

\Rightarrow obtain the lines of regression hence,
to find the coefficient of correlation for the data

x	1	3	4	2	5	8	9	10	12	15
y	8	6	10	8	12	16	16	10	32	32

$$\bar{x} = \frac{\sum x}{n} = \frac{70}{10} = 7$$

$$\bar{y} = \frac{\sum y}{n} = \frac{150}{10} = 15$$

x	y	$x = x - \bar{x}$	$y = y - \bar{y}$	x^2	y^2
1	8	-6	-7	42	36
3	6	-4	-9	36	36
4	10	-3	-5	9	25
2	8	-5	-7	25	49
5	12	-2	-3	4	9
8	16	1	1	1	1
9	16	2	1	4	1
10	10	3	-5	9	25
13	32	6	14	102	289
15	32	8	16	64	289
				<u>= 360</u>	<u>818</u>
				<u>$\frac{360}{818}$</u>	<u>$\frac{818}{818}$</u>

Regression line x on y (i.e) x on y related

$$(x - \bar{x}) = \frac{\sum xy}{\sum y^2} (y - \bar{y}) \quad \text{relation}$$

$$(x - \bar{x}) = \frac{360}{818} (y - \bar{y})$$

$$\text{③ } \frac{360}{818} = \frac{B_2 + \epsilon_2}{B_1} = \bar{x}$$

$$x = 0.44y + 0.3985$$

Regression line y on x on y related

$$(y - \bar{y}) = \left(\frac{\sum xy}{\sum x^2} \right) (x - \bar{x}) \quad \text{relation}$$

$$\text{For } \frac{360}{818} = \frac{B_2}{B_1} = \frac{360}{204} = 1.7647$$

$$204 - B_1 = B_0 = 1.7647x - 12.3529$$

$$1.7647x + 12.3529 = y \quad \text{or} \quad y = 1.7647x + 12.3529$$

Here regression coefficient x on $y = 0.4401$

Regression y on $x = 1.7647$

$$\tau = \sqrt{(0.4401)^2 \times (1.7647)}$$

$$\tau = \underline{0.8812}$$

3) In a partially clivaged record of correlation data, the following results only are available:

Regression equations are $4x - 5y + 33 = 0$

$$20\bar{x} - 9\bar{y} = 107$$

Calculate (i) the mean values of x & y

(ii) the coefficient of correlation r_{xy}

(iii) S.D. of y .

Soln:- (i) W.K.T. (\bar{x}, \bar{y}) is point on the regression

line

$\therefore 4\bar{x} - 5\bar{y} + 33 = 0 \Rightarrow 4\bar{x} - 5\bar{y} = -33$

$$20\bar{x} - 9\bar{y} = 107 \rightarrow \textcircled{2}$$

Solve eq (1) & eq (2)

Consider,

$$4\bar{x} - 5\bar{y} = -33 \quad \text{--- (1)}$$

$$4\bar{x} = -33 + 5\bar{y}$$

$$\bar{x} = \frac{-33 + 5\bar{y}}{4} \rightarrow \textcircled{3}$$

Sub Eq (3) in Eq (2) = x

$$20\bar{x} - 9\bar{y} = 107$$

$$20\left(\frac{-33 + 5\bar{y}}{4}\right) - 9\bar{y} = 107$$

$$-165 + 25\bar{y} - 9\bar{y} = 107$$

$$16\bar{y} = 107 + 165$$

$$16\bar{y} = 272$$

$$\bar{y} = 17$$

Final Sub $\bar{y} = 17$ in eq (1)

$$4\bar{x} - 5\bar{y} = -33$$

$$4\bar{x} - 5(17) = -33$$

$$4\bar{x} - 85 = -33$$

$$4\bar{x} = -33 + 8 \Rightarrow \bar{x} = -33/4 + 2 = -5.25$$

$$4\bar{x} = 5s$$

$$\bar{x} = 13$$

Q. Is it true? Explain.

$$\therefore \bar{x} = 13, \bar{y} = 17$$

(ii) Regression line y on x

$$sy = Ax + b$$

$$y = \frac{4}{5}x + \frac{33}{5}$$

Regression line x on y

$$20x = 9y + 107$$

$$x = \left(\frac{9}{20}y + \frac{107}{20} \right)$$

Regression coefficient (of y on x) = $\frac{4}{5} = 0.8$

$$r_{xy} = \frac{9}{20} = 0.45$$

$$r = \sqrt{(0.8)(0.45)} = 0.6$$

(iii) $\sigma_x^2 = 9 \Rightarrow \sigma_x = 3$

Regression coefficient w.r.t $y = r \cdot \frac{\sigma_x}{\sigma_y} = 0.45$

$$\sigma_y = \sqrt{\frac{0.6 \times 3}{0.45}} = \sqrt{\frac{1.8}{0.45}} = \sqrt{4} = 2$$

If θ is the angle b/w the 2 regression lines, s.t. $\tan \theta = \frac{(1-r^2)}{\sqrt{1+r^2}}$.

Explain the significance when $r=0$ & $r=\pm 1$.

Soln) - Regression line y on x (Ans)

$$(y - \bar{y}) = r \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \rightarrow \text{Ans}$$

Regression line $\hat{y} = \text{const} + b_1 x$

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

Slopes of eqn ① & ②

$$m_1 = r \frac{\sigma_y}{\sigma_x} \quad m_2 = \frac{\sigma_y}{r \sigma_x}$$

$$\tan \theta = \frac{m_2 - m_1 + r m_1 m_2}{1 + m_1 m_2}$$

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

$$\tan \theta = \frac{(1 - r^2) \cdot \sigma_y}{r \sigma_x}$$

$$\tan \theta = \frac{(1 - r^2) \cdot \sigma_y}{r \sigma_x} = \frac{(1 - r^2) \cdot \sigma_y}{\sqrt{1 - r^2} \cdot \sqrt{1 - r^2} \cdot \sigma_x} = \frac{(1 - r^2)^{1/2} \cdot \sigma_y}{\sigma_x}$$

$$\tan \theta = \frac{(1 - r^2)^{1/2} \cdot \sigma_y}{\sigma_x}$$

$$\tan \theta = \frac{(1 - r^2)^{1/2} \cdot \sigma_y}{\sigma_x} = \frac{(1 - r^2)^{1/2} \cdot \sigma_y}{\sigma_x}$$

when $r = 0 \Rightarrow \tan \theta = \infty \Rightarrow \theta = \frac{\pi}{2}$

when the variables are linearly uncorrelated
then they are perpendicular to each other.

when $r = \pm 1 \Rightarrow \tan \theta = 0 \Rightarrow \theta = 0 \text{ or } \pi$

when lines of regression coincide then there

is a perfect correlation b/w 2 variables.

→ If the tangent of the angle θ b/w the lines of regression of y on x & x on y is 0.6 and σ_x is the S.D. of x , then find the coefficient of correlation b/w x & y .

∴ $\tan \theta = 0.6$

Soln)- Given α be the angle b/w the lines of regression

$$\tan \alpha = 0.6$$

w.k.t.

$$\sigma_y = 2\sigma_x$$

$$\tan \alpha = \left(\frac{1-\rho^2}{\rho} \right) \left(\frac{\sigma_x \cdot \sigma_y}{\sigma_x^2 + \sigma_y^2} \right)$$

$$0.6 = \left(\frac{1-\rho^2}{\rho} \right) \left(\frac{2\sigma_x^2}{\sigma_x^2 + \sigma_y^2} \right)$$

$\Rightarrow \frac{1-\rho^2}{\rho} = \frac{0.6}{\frac{2\sigma_x^2}{\sigma_x^2 + \sigma_y^2}}$

$$= \left(\frac{1-\rho^2}{\rho} \right) \left(\frac{2\sigma_x^2}{5\sigma_x^2} \right)$$

$$0.6 = \frac{(1-\rho^2)}{\rho} \cdot \frac{2}{5}$$

$$3\rho = 2 - 2\rho^2$$

$$2\rho^2 + 3\rho - 2 = 0$$

$$\textcircled{1} \leftarrow 1 \rightarrow \rho = \frac{1}{\sqrt{2}} = 0.5 \quad \text{or} \quad \rho = -2 \quad \text{not possible.}$$

$$\textcircled{2} \leftarrow \rho = \frac{1}{\sqrt{2}} = 0.5$$

H.W. \rightarrow find the correlation coefficient b/w x & y for the following data. Also obtain the regression line.

x : 1	2	3	4	5	6	7	8	9	10
y : 10	12	16	20	25	36	41	49	40	50

\rightarrow In a bivariate distribution, it is found that $\sigma_x = \sigma_y$ & the acute angle b/w the lines of regression is $\tan^{-1}(3)$. Find the correlation coefficient.

Soln)- w.k.t. $\tan \alpha = \left(\frac{1-\rho^2}{\rho} \right) \left(\frac{\sigma_x \cdot \sigma_y}{\sigma_x^2 + \sigma_y^2} \right)$

given $\alpha = \tan^{-1}(3)$ $\sigma_x = \sigma_y$

$$\tan(\tan^{-1}(3)) = \left(\frac{1-\rho^2}{\rho} \right) \left(\frac{\sigma_x}{\sigma_x} \right)$$

$$3 = \left(\frac{1-\rho^2}{\rho} \right) \left(\frac{1}{2} \right)$$

$$6\tau = 1 - \sigma^2 \quad \text{or} \quad \tau = \frac{1 - \sigma^2}{6}$$

$$\tau^2 + 6\tau - 1 = 0$$

$$\Rightarrow \tau = \frac{-6 \pm \sqrt{36 + 4}}{2} = \frac{-6 \pm \sqrt{40}}{2} = \frac{-6 \pm 2\sqrt{10}}{2} = -3 \pm \sqrt{10}$$

$$\Rightarrow \tau = -3 + \sqrt{10} \approx 0.1623$$

8) With usual notation, compute \bar{x}, \bar{y} & γ from the following lines of regression.

$$2x + 3y + 1 = 0$$

$$x + 6y - 4 = 0$$

Soln) - w.r.t.

(\bar{x}, \bar{y}) is point on the regression

line

$$2\bar{x} + 3\bar{y} = -1 \rightarrow ①$$

$$\bar{x} + 6\bar{y} = 4 \rightarrow ②$$

from eq ②

Sub. \bar{x} in eq ①

$$2(4 - 6\bar{y}) + 3\bar{y} = -1$$

$$8 - 12\bar{y} + 3\bar{y} = -1$$

$$8 - 9\bar{y} = -1 \rightarrow ③$$

$$-9\bar{y} = -9 \rightarrow \bar{y} = 1$$

Sub. $\bar{y} = 1$ in eq ②

$$2\bar{x} + 3(1) = 4$$

$$2\bar{x} = 4 - 3$$

$$2\bar{x} = 1 \rightarrow \bar{x} = \frac{1}{2}$$

$$\bar{x} = -2, \bar{y} = 1$$

Regression line of y on x & x on y .

$$2x + 3y + 1 = 0 \quad \text{for } x \text{ on } y$$

$$3y = -1 - 2x \quad \text{or } y = -\frac{1}{3} - \frac{2}{3}x \quad \text{for } y \text{ on } x$$

$$y = -\frac{1}{3} - \frac{2}{3}x \quad \text{(i)}$$

$$\text{Eq. } x + 6y - 4 = 0 \quad \text{for } y \text{ on } x$$

$$6y = 4 - x \quad \text{or } y = \frac{4}{6} - \frac{1}{6}x \quad \text{for } x \text{ on } y$$

$$y = \frac{2}{3} - \frac{1}{6}x \quad \text{(ii)}$$

regression coefficients are $\frac{-2}{3}$ & -6

$$\therefore r = \sqrt{\left(\frac{-2}{3}\right)(-6)} = \pm 2$$

$r \in [-1, 1]$ not possible.

Regression line of x on y & y on x .

$$2x + 3y + 1 = 0$$

$$2x = -1 - 3y$$

$$x = \frac{-3y - 1}{2}$$

$$\text{Eq. } x + 6y - 4 = 0$$

$$6y = 4 - x$$

$$y = -\frac{x}{6} + \frac{4}{6}$$

Regression coefficients are $\frac{-3}{2}$ & $\frac{1}{6}$

$$\therefore r = \sqrt{\left(-\frac{3}{2}\right)\left(\frac{1}{6}\right)} = \sqrt{\frac{1}{4}}$$

$$r = \pm \frac{1}{2} = \pm 0.5$$

Since $\frac{3}{2}$ & $\frac{1}{6}$ are negative

$$r = -0.5$$

$$\therefore \bar{x} = -2, \bar{y} = 1, r = -0.5$$

Q) From the following eqn of the lines of regression compute the mean of x (\bar{x}), mean of y (\bar{y}) & correlation coefficient r .

$$(i) y = x + 45(-y - 3 = 0 \text{ Ans!} \rightarrow \bar{x} = 1 = \bar{y}, r = 0.5)$$

$$(ii) y = 0.516x + 33.73x, \text{ Ans! } \bar{x} = 67.67 \\ x = 0.512y + 32.52 \text{ Ans! } \bar{y} = 68.65$$

$\therefore r = \frac{\bar{x}\bar{y}}{\sqrt{\bar{x}^2 - (\bar{x})^2} \sqrt{\bar{y}^2 - (\bar{y})^2}}$ ~~now~~ $r = +0.51$ m.s.p.r.

$$\text{Mean of } x = \frac{1}{n} \sum x_i = \frac{1}{10} (67.67 \times 10) = 67.67$$

$$\text{Mean of } y = \frac{1}{n} \sum y_i = \frac{1}{10} (68.65 \times 10) = 68.65$$

$$r = \frac{\bar{x}\bar{y}}{\sqrt{\bar{x}^2 - (\bar{x})^2} \sqrt{\bar{y}^2 - (\bar{y})^2}}$$

$$r = \frac{67.67 \times 68.65}{\sqrt{67.67^2 - 67.67^2} \sqrt{68.65^2 - 68.65^2}}$$

$$r = \frac{67.67 \times 68.65}{\sqrt{0} \sqrt{0}}$$

$$r = \frac{67.67 \times 68.65}{0 \times 0}$$

$$r = \frac{67.67 \times 68.65}{0 \times 0}$$

Rank Correlation Coefficient
 The coefficient of correlation in respect of the ranks of some two characteristics of an individual or an observation is called the rank correlation coefficient. It is denoted by γ (rho). γ is computed by

$$\gamma = 1 - \frac{6 \sum (x-y)^2}{n(n^2-1)} \quad (\text{or}) \quad 1 - \frac{6 \sum d^2}{n(n^2-1)}$$

where

$$d^2 = x-y$$

Note - (1) If the rankings of x, y are entirely in the same order like for example,

$$x: 1, 2, 3, 4, 5 ; y: 1, 2, 3, 4, 5 \text{ then}$$

$\sum d^2 = \sum (x-y)^2 = 0$. This will give us $\gamma = +1$ and is called perfect direct correlation.

If the ranking of x and y are entirely in the opposite order like for example,

$$x: 1, 2, 3, 4, 5 ; y: 5, 4, 3, 2, 1$$

then

$$\begin{aligned} \sum d^2 &= \sum (x-y)^2 = (-4)^2 + (-2)^2 + (0)^2 + (2)^2 + (4)^2 \\ &= 40 \end{aligned}$$

$\therefore \gamma = 1 - \frac{6(40)}{5(5^2-1)} = -1$ and is called perfect inverse correlation.

(2) γ in case of repeated ranks

If two or more magnitudes are repeated,

particular rank. In such cases we assign the average rank to all those magnitudes and use the correction factor $\frac{m(m^2-1)}{12}$ along with $\sum d^2$ in the formula for S where m denotes no. of times a magnitude repeated.

The correction factor must be added every time the 'tie' occurs for a particular rank.

$$S = \frac{6 \left[\sum d^2 + \frac{m(m^2-1)}{12} \right]}{n(n^2-1)}$$

Ex:- If there is a tie for the second rank in two magnitudes, the rank $\frac{2+3}{2} = 2.5$ is assigned to both the magnitudes, thereby welding both 2nd & 3rd ranks ($m=2$).

Again, suppose there is a tie for 5th rank among 3 magnitudes, we take the average of the 3 ranks in a row starting from 5. i.e.

$$\text{i.e. } \frac{5+6+7}{3} = 6 \quad (m=3)$$

Problem:- 0 1 2 3 4 5 6 7 8 9 10
Ten competitors in a beauty contest are ranked by 2 judges in the following order. Compute the coefficient of rank correlation:

I	1	6	5	3	10	2	4	9	7	8
II	6	4	9	8	1	2	3	10	5	7

Soln) - we have

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

For the given data, $n=10$

$$\text{we have } \sum d^2 = (1-6)^2 + (6-4)^2 + (5-9)^2 + (3-8)^2$$

$$+ (10-1)^2 + (2-2)^2 + (4-3)^2 + (9-0)^2 + \\ (7-5)^2 + (8-4)^2$$

$$\sum d^2 = 25 + 4 + 16 + 25 + 81 + 0 + 1 + 1 + 4 + 1$$

$$\therefore S = \sqrt{\frac{158}{10(10^2-1)}} = \sqrt{\frac{158}{990}} = 0.412$$

$$\text{Hence } S = \sqrt{\frac{6(158)}{10(10^2-1)}} = 0.412$$

2) Ten Competitors in music contest are ranked by 3 judges A, B, C in the following order. Use the rank correlation coefficient to decide which pair of judges have the nearest approach to common taste of music.

A 1 2 6 5 9 4 3 7 8 10

B 3 5 8 4 7 10 2 1 6 9

C 6 4 9 8 1 2 3 10 5 7

Rest & pairwise of
ranked pairs of
judges are

D 1 8 7 6 5 4 3 2 10 9

E 2 10 8 6 5 4 3 1 7 9

		Σd_{AB}^2	Σd_{BC}^2	Σd_{CA}^2
1	3	6	4	9
6	5	4	1	4
5	8	9	9	16
10	82	43	62	36
3	7	1	16	36
2	10	2	64	64
4	2	3	4	1
9	1	10	64	81
7	6	$(5-x)^2 = b$	$(4)^2 = b$	$(6)^2 = b$
8	9	$\frac{1}{\sum d_{AB}^2 = 200}$	$\frac{4}{\sum d_{BC}^2 = 214}$	$\frac{1}{\sum d_{CA}^2 = 60}$

we have,

$$S = 1 - \frac{6 \sum d^2}{n(n^2-1)} \quad n=10.$$

$$S_{AB} = 1 - \frac{6(200)}{10(10^2-1)} = -0.21$$

$$S_{BC} = 1 - \frac{6(214)}{10(10^2-1)} = -0.297$$

$$S_{CA} = 1 - \frac{6(60)}{10(10^2-1)} = +0.636$$

Here S_{AB} & S_{BC} are negative which means their tastes (A & B ; B & C) are opposite. But S_{CA} is positive and is nearer to 1 (perfect correlation).

Thus judge C is at home the nearest approach to common taste of music.

3) Ten students got the following percentage of marks in 2 subjects A & Y. Compute their rank correlation coefficient.

Marks in X: 78, 36, 98, 25, 75, 82, 90, 62, 65, 39
 Marks in Y: 84, 51, 91, 60, 68, 62, 86, 58, 53, 39

(d)
(d)

Soln:-

Marks in Rank x (x), Marks in Rank y (y) $d = (x-y)$ $d^2 = (x-y)^2$

78 4 84 3 $d = 1$ $d^2 = 1$
 36 9 51 9 $d = 0$ $d^2 = 0$

98 1 91 1 $d = 0$ $d^2 = 0$

25 10 60 6 $d = 4$ $d^2 = 16$

75 5 68 4 $d = 1$ $d^2 = 1$

82 3 62 5 $d = -2$ $d^2 = 4$

90 12 86 2 $d = 10$ $d^2 = 100$

62 7 58 0 $d = 7$ $d^2 = 49$

65 6 53 8 $d = -2$ $d^2 = 4$

39 8 39 10 $d = -2$ $d^2 = 4$

Total $\sum d^2 = 30$

we have $r_s = 1 - \frac{6 \sum d^2}{n(n^2-1)}$ $n=10$

sum $\sum d^2 = 30$ $n=10$

$r_s = 1 - \frac{6(30)}{10(10^2-1)} = 0.8181 \approx 0.82$

$$\therefore r_s = 0.82$$

4) Compute the rank correlation coefficient for the following ranks of ten candidates in 2 subjects.

	Subject 1 (marks)										Subject 2 (marks)									
Rank	1	2	3	4	5	6	7	8	9	10	1	2	3	4	5	6	7	8	9	10
Sub-1	4	8	7	6	5	9	10	3	2	1	6	7	8	1	5	10	9	2	3	4
Sub-2	6	7	8	1	5	10	9	2	3	4	5	6	7	8	1	5	10	9	2	3

Soln) - $\boxed{0.73}$ To 3 s.f.

5) Compute the rank correlation coefficient for the following data.

x	68	64	75	50	64	80	75	40	55	64
y	62	58	68	45	58	60	68	48	50	70

Soln) - $d = x - y$ d^2

x	Rank(x)	y	Rank(y)	$d = x - y$	d^2
68	4	62	25	-21	441
64	6	58	22	-2	4
75	8.5	68	3.5	-5	25

x	50	45	64	80	75	62	68	55	64
y	6	8	1	5	7	25	22	0	6
$d = x - y$	-44	-37	-63	-75	-68	-19	-10	-55	-58
d^2	1936	1369	3969	5625	4624	361	100	3025	3364

$$\frac{2+3}{2} = 2.5 \rightarrow \frac{5+6+7}{3} = \frac{18}{3} = 6 \quad \frac{3+4}{2} = \frac{7}{2} = 3.5 \quad (m=2)$$

$$(m=2) \quad 216 \quad (m=3)$$

$$S = 1 - \frac{6 \left[\sum d_{ij}^2 + \frac{m(m^2-1)}{12} \right]}{n(n^2-1)} \quad n=10$$

$$= 1 - \frac{6 \left[2(2^2-1) + 3(3^2-1) + \frac{2(2^2-1)}{12} \right]}{10(10^2-1)} = 0.545$$

$R = 0.545$ \therefore $r_s = 0.545$

Curve Fitting

The method of finding a specific relation $y = f(x)$ for the data to satisfy as accurately as possible and such an eqn is called the best fitting eqn or the curve of best fit.

fit. for 1992 London Datas. ODE into chart

+ fitting of a straight line $y = ax + b$.

Consider set of n given values (x_i, y_i) fitting the straight line $y = ax + b$ where a & b are parameters to be determined.

The residual $R = y - (ax + b)$ is the difference between the observed and estimated values of y .

By the method of least squares we find parameters a and b such that the sum of squares of the residuals is minimum.

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

$$\text{i.e. } S = \sum_{i=1}^n [y_i - (ax_i + b)]^2$$

Treating S as a function of 2 parameters a & b the necessary conditions for S to be minimum

$$\text{are } \frac{\partial S}{\partial a} = 0 \quad \& \quad \frac{\partial S}{\partial b} = 0$$

$$\text{i.e. } \sum_{i=1}^n [y_i - (ax_i + b)](-x_i) = 0$$

$$\text{and } \sum_{i=1}^n [y_i - (ax_i + b)](-1) = 0$$

Dividing both the eqns by 2, we have

$$-\sum_{i=1}^n xy_i + \sum_{i=1}^n ax_i^2 + \sum_{i=1}^n bx_i = 0$$

$$-\sum_{i=1}^n y_i + \sum_{i=1}^n ax_i + \sum_{i=1}^n b = 0$$

But $\sum_{i=1}^n b = b + b + \dots + b$ ~~times~~ n times.

So $a\sum x + nb = \sum y$

Hence $a\sum x^2 + b\sum x = \sum xy$ to ~~fitting~~

Now $a\sum x + nb = \sum y$

These eqns are called normal eqns for ~~fitting~~

fitting the straight line $y = ax + b$

Problem

> fit a curve of the form $y = ax + b$ to the following table.

x	14	13	9	5
y	14	13	9	5
x^2	196	169	81	25
xy	196	169	81	25
x^3	3024	2856	729	125

Soln - ~~by~~ ~~using~~ ~~method~~ ~~of~~ ~~minimizing~~ ~~sum~~ ~~of~~ ~~squares~~ ~~of~~ ~~residuals~~ ~~to~~ ~~fit~~ ~~straight~~ ~~line~~

$\sum a(x_i) + nb = \sum y$ ~~is~~ ~~not~~ ~~possible~~

Minimizing $\sum (y_i - ax_i - nb)^2$

$\begin{aligned} 1 & 14 & 14 & 14 & 14 \\ 2 & 13 & 26 & 4 & 9 \\ 3 & 9 & 27 & 27 & 27 \\ 4 & 5 & 40 & 16 & 16 \\ 5 & 16 & 80 & 80 & 80 \\ \hline & 65 & 143 & 97 & 55 \end{aligned}$

$n = 5$

Normal eqns are

$$a\sum x + nb = \sum y$$

$$a\sum x^2 + b\sum x = \sum xy$$

such that $a + b = 14$

$$a(15) + 5b = 43$$

$$a(55) + 15b = 97$$

$$15a + 5b = 43 \rightarrow ①$$

$$55a + 15b = 97 \rightarrow ②$$

from eq ①

$$15a = 43 - 5b$$

$$a = \frac{43}{15} - \frac{5b}{15}$$

Sub a in eq ② $5.0 + 1.0 + 1.0 + 1.0 = 8.0$

$$55a + 15b = 97$$

$$55\left(\frac{43}{15} - \frac{5b}{15}\right) + 15b = 97$$

$$\cancel{55}\left(\frac{43 - 5b}{15}\right) + 15b = 97$$

$$473 - 55b + 45b = 291$$

$$473 - 10b = 291$$

$$-10b = -182$$

$$b = 18.2$$

Now Sub $b = 18.2$ in eq ①

$$15a + 5b = 43$$

$$15a + 5(18.2) = 43$$

$$15a + 91 = 43$$

$$15a = -48$$

$$a = -3.2$$

$$d.o. = (-3.2)dF + 18.2$$

$$\therefore a = -3.2, b = 18.2$$

$$\text{so, } y = ax + b$$

$$y = (-3.2)x + 18.2 \quad \text{as a curve of best fit.}$$

2) Fit the straight line $y = ax + b$ to the following data

x	-5	-3	-1	1	2	4
y	0.4	0.1	-0.2	-0.3	0.1	0.4

$$d.o. = dF + \left(\frac{dF + 1.0}{2}\right) dF$$

$$d.o. = dF - \left(\frac{dF + 1.0}{2}\right) dF$$

Soln)-

$$\begin{array}{r}
 \text{S.L} & \text{Y} & \text{Sum} \\
 -5 & 0.4 & -2 \frac{42}{21} + \frac{25}{21} = 0 \\
 -3 & 0.1 & -0.3 \quad \text{Eq. 9 to D. side} \\
 -1 & -0.2 & +0.2 \\
 0 & -0.3 & 0 \\
 1 & -0.3 & -0.3 \\
 2 & 0.1 & 0.2 \\
 4 & 0.4 & 1.6 \\
 \hline
 -2 & & +0.2 \\
 & & \underline{-0.6} \\
 & & \underline{\underline{56}}
 \end{array}$$

$n = 7$

Normal eqns are

$$a \sum x + b \sum Y = \sum y \quad \text{Eq. 1st. row}$$

$$a \sum x^2 + b \sum x = \sum xy \quad \text{Eq. 2nd. row}$$

$$a(-2) + b(7) = -0.2 \quad \text{Eq. 3rd. row}$$

$$\Rightarrow -2a + 7b = -0.2 \rightarrow ①$$

$$a(56) + b(-2) = -0.6$$

$$56a - 2b = -0.6 \rightarrow ②$$

$$-2a + 7b = +0.2 \rightarrow ①$$

$$56a - 2b = -0.6 \rightarrow ②$$

from eq. ① $-2a = +0.2 - 7b$

$$a = -\frac{0.2}{2} + \frac{7}{2}b$$

Sub a in eq. ②

$$56a - 2b = -0.6$$

$$56\left(-\frac{0.2}{2} + \frac{7}{2}b\right) - 2b = -0.6$$

$$56\left(-\frac{0.2}{2} + \frac{7}{2}b\right) - 2b = -0.6$$

$$-5.6 + 196b - 2b^2 = -0.6$$

$$194b^2 = -0.6 + 5.6$$

$$194b^2 = 5.0$$

$$b = 0.025$$

$$-5.6 + 196 \cdot 0.025 - 2 \cdot 0.025^2 = 0$$

$$-5.6 + 4.9 - 0.01 = 0$$

$$4.9 - 0.01 = 0$$

third equat. find as $y = ax + b$ for the
fit straight line $y = ax + b$ for the
following data.

x	1	3	4	6	8	9	14
y	1	2	4	10	52	78	9

Soln -

$$\begin{aligned} & \text{Sum of } x^2 \\ & \sum x^2 = 1 + 9 = 10 \\ & 1 + 9 + 24 + 36 = 60 \end{aligned}$$

Mean of x and y is 6.5

which gives $a = 3.6$

$$64$$

$$\begin{array}{ccccccc} 5 & 40 & 20 & 10 & 02 & 02 & 0 \\ 8 & 56 & 16 & 81 & 54 & 54 & 8 \\ 9 & 72 & 16 & 12 & 12 & 12 & 9 \\ 11 & 88 & 24 & 12 & 12 & 12 & 11 \\ 14 & 126 & 36 & 196 & 126 & 126 & 14 \end{array}$$

$$\begin{array}{ccccccc} & 95 & 36 & 524 & 95 & 95 & 14 \\ \hline \Sigma x & 56 & \Sigma y & 10 & \Sigma xy & 364 & \Sigma x^2 = 524 \end{array}$$

$$\begin{array}{ccccccc} & 95 & 36 & 524 & 95 & 95 & 14 \\ \hline \Sigma x & 56 & \Sigma y & 10 & \Sigma xy & 364 & \Sigma x^2 = 524 \end{array}$$

The normal eqn are $\Sigma x + nb = cn = 8$

$$\Sigma y = a \sum x + nb$$

$$\Sigma xy = a \sum x^2 + b \sum x$$

The normal eqn becomes

$$56a + 8b = 40 \quad (1)$$

$$524a + 56b = 364 \quad (2)$$

By solving, we have

$$a = 0.636 \approx 0.64, \quad b = 0.545 \approx 0.55$$

0.5114

$$\therefore y = ax + b$$

$$y = 0.64x + 0.55$$

- 4) Find the eqn of the best fitting straight for the following data & hence estimate the value of the dependent variable to the value 30 of the independent variable.

x	5	8	10	15	20	25
y	16	19	23	26	30	

Ans:-

$$a = 0.7, \quad b = 12.3$$

$$y = (0.7)x + 12.3$$

- 5) Fit a straight line in the least square sense for the following data.

x	50	70	106	120
y	19	15	21	25

Ans:-

$$a = 0.187 \approx 0.19, \quad b = 2.27 \approx 2.28$$

$$y = (0.19)x + (2.28)$$

$$(8 = 0) \quad 19 + x \cdot 30 = 2.28 \\ 170 + x \cdot 30 = 2.28 \\ x \cdot 30 = 167.72 \\ x = 5.59$$

Fitting of a curve of the form $y = ax^b$

Consider, $y = ax^b$

Take log on both sides.

we get,

$$\log y = \log (ax^b)$$

$$\log_e y = \log_e a + \log_e x^b$$

$$\log_e y = \log_e a + b \log_e x$$

$$\text{Thus, } \log_e y = A + BX \quad \rightarrow \textcircled{1}$$

$$\text{where } \log_e y, \quad A = \log_e a, \quad B = b, \quad x = \log_e x$$

The normal eqns associated with $\textcircled{1}$ are follows

$$\sum y = nA + B\sum x \quad \rightarrow \textcircled{2}$$

$$\sum xy = A\sum x + B\sum x^2 \quad \rightarrow \textcircled{3}$$

Solving $\textcircled{2}$ & $\textcircled{3}$, we obtain A, B

But, $\log_e a = A$

$$\Rightarrow a = e^{A+Bx}$$

Substituting the values of a, b in $y = ax^b$ in the required form, we get the curve of best fit.

Remark! - We can similarly fit the curves

$$y = ae^{bx}, \quad y = ab^x \text{ etc.}$$

▷ Fit a least square geometric curve $y = ax^b$ from the following data.

x	1	2	3	4	5
y	0.5	2	4.5	8	12.5

Consider $y = ax^b$

Take log on both sides.

$$\log y = \log a + b \log x \text{ per chart}$$

$$y = (A + bX)$$

where $\gamma = \log y$, $A = \log a$, $X = \log x$

x	y	$x = \log x$	$\gamma = \log y$	XY	x^2
1	0.5	0	-0.6931	0	0
2	2	0.6931	0.4803	0.4803	1.4803
3	4.5	1.0986	1.5040	1.6522	1.2069
4	8	1.3862	2.0794	2.8824	1.9215
5	12.5	1.6094	2.5257	4.0648	2.5901
		4.7873	6.1091	9.0797	6.5198

Normal eqns

$$nA + b \sum x = \sum \gamma$$

$$A \sum x + b \sum x^2 = \sum XY$$

Here $n=5$

$$5A + 4.7873b = 6.1091$$

$$4.7873A + 6.1091b = 9.0797$$

$$A = -0.6932, b = 2$$

$$\log a = -0.6932 \Rightarrow a = e^{-0.6932} = 0.4999$$

$$\Rightarrow \log a = e^{-0.6932} = 0.4999$$

$$y = ax^b$$

$$y = (0.5) x^2 \text{ curve of best fit}$$