

Unit-5

correlation and regression

correlation :- (Linear correlation or Karl Pearson co-efficient of correlation)

correlation co-efficient between two random variables x and y . It is denoted by $r(x, y)$ (or) r_{xy} (or) ρ_{xy}

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

$$\text{cov}(x, y) = N \sum xy - \sum x \sum y$$

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

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

$\text{var}(x) = N \sum x^2 - (\sum x)^2$
 $\text{var}(y) = N \sum y^2 - (\sum y)^2$
 $(r = -1 \text{ to } 1 \text{ (or) } -1 \leq r \leq 1)$

① problems based on correlation co-efficient for the following data.

x : 65 66 67 67 68 69 70 72

y : 67 68 65 68 72 72 69 71

Soln:

To find correlation co-efficient

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

$$N = 8$$

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

$$\begin{aligned}
 r &= \frac{N \Sigma XY - \Sigma X \Sigma Y}{\sqrt{N \Sigma X^2 - (\Sigma X)^2} \sqrt{N \Sigma Y^2 - (\Sigma Y)^2}} \\
 &= \frac{8(37560) - (544)(552)}{\sqrt{8(37028) - (544)^2} \sqrt{8(38132) - (552)^2}} \\
 &= \frac{300480 - 300288}{\sqrt{296224 - 295936} \sqrt{305056 - 304704}} \\
 &= \frac{192}{\sqrt{288} \sqrt{352}} \\
 &= \frac{192}{\sqrt{101376}} = \frac{192}{318.3959} = 0.6030
 \end{aligned}$$

$$r = 0.6030$$

$$\therefore r = 0.603$$

② Find the correlation co-efficient from the following data.

x: 10 12 14 16 18 20 22 24

y: 14 18 16 22 26 28 27 30

Soln.:-

To Find correlation co-efficient

$$r = \frac{N \sum XY - \sum X \sum Y}{\sqrt{N \sum X^2 - (\sum X)^2} \sqrt{N \sum Y^2 - (\sum Y)^2}}$$

$$N = 8$$

X	Y	X ²	Y ²	XY
10	14	100	196	140
12	18	144	324	216
14	16	196	256	224
16	22	256	484	352
18	26	324	676	468
20	28	400	784	560
22	27	484	729	594
24	30	576	900	720
$\sum X = 136$	$\sum Y = 181$	$\sum X^2 = 2480$	$\sum Y^2 = 4349$	$\sum XY = 3274$

$$r = \frac{N \sum XY - \sum X \sum Y}{\sqrt{N \sum X^2 - (\sum X)^2} \sqrt{N \sum Y^2 - (\sum Y)^2}}$$

$$= \frac{8(3274) - (136)(181)}{\sqrt{8(2480) - (136)^2} \sqrt{8(4349) - (181)^2}}$$

$$= \frac{26192 - 24616}{\sqrt{19840 - 18496} \sqrt{34792 - 32761}}$$

$$= \frac{1576}{\sqrt{1344} \sqrt{2031}}$$

$$= \frac{1576}{\sqrt{2729664}}$$

$$= \frac{1576}{1652.1694} = 0.9538$$

$$\therefore r = 0.9538$$

$$\text{var}(X) = N \sum x^2 - (\sum x)^2 = 1344$$

$$\text{var}(Y) = N \sum y^2 - (\sum y)^2 = 2031$$

$$\text{cov}(X, Y) = N \sum XY - \sum X \sum Y = 1576$$

$$\text{S.D}(X) = \sqrt{N \sum x^2 - (\sum x)^2} = 36.6606$$

$$\text{S.D}(Y) = \sqrt{N \sum y^2 - (\sum y)^2} = 45.0666$$

③ Find the correlation co-efficient for following data

X: 10 14 18 22 26 30

Y: 18 12 24 6 30 36

Soln:-

To Find: correlation co-efficient

$$r = \frac{N \sum XY - \sum X \sum Y}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

$$N = 6$$

X	Y	X ²	Y ²	XY
10	18	100	324	180
14	12	196	144	168
18	24	324	576	432
22	6	484	36	132
26	30	676	900	780
30	36	900	1296	1080
$\Sigma X = 120$	$\Sigma Y = 126$	$\Sigma X^2 = 2680$	$\Sigma Y^2 = 3276$	$\Sigma XY = 2772$

$$\begin{aligned}
 r &= \frac{N \Sigma XY - \Sigma X \Sigma Y}{\sqrt{N \Sigma X^2 - (\Sigma X)^2} \sqrt{N \Sigma Y^2 - (\Sigma Y)^2}} \\
 &= \frac{6(2772) - (120)(126)}{\sqrt{6(2680) - (120)^2} \sqrt{6(3276) - (126)^2}} \\
 &= \frac{16632 - 15120}{\sqrt{16080 - 14400} \sqrt{19656 - 15876}} \\
 &= \frac{1512}{\sqrt{1680} \sqrt{3780}} \\
 &= \frac{1512}{\sqrt{6350400}} \\
 &= \frac{1512}{2520} = 0.6
 \end{aligned}$$

$$r = 0.6$$

$$\text{var}(X) = N \Sigma X^2 - (\Sigma X)^2 = 1680$$

$$\text{var}(Y) = N \Sigma Y^2 - (\Sigma Y)^2 = 3780$$

$$\text{cov}(X, Y) = N \Sigma XY - \Sigma X \Sigma Y = 1512$$

$$\text{S.D}(X) = \sqrt{N \Sigma X^2 - (\Sigma X)^2} = 40.9878$$

$$\text{S.D}(Y) = \sqrt{N \Sigma Y^2 - (\Sigma Y)^2} = 61.4817$$

limits for co-efficient (ρ)

prove that $-1 \leq \rho_{xy} \leq 1$

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

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

$$\therefore \rho^2(x, y) = \frac{(\sum a_i b_i)^2}{(\sum a_i^2)(\sum b_i^2)} \longrightarrow \textcircled{1}$$

where, $a_i = x_i - \bar{x}$

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

we have the schwarz inequality which states that if $a_i, b_i; i=1, 2, \dots, n$ are real quantities.

then,

$$\left[\sum_{i=1}^n a_i b_i \right]^2 \leq \left[\sum_{i=1}^n a_i^2 \right] \left[\sum_{i=1}^n b_i^2 \right]$$

The sign of equality holding if and

$$\text{only if; } \frac{a_1}{b_1} = \frac{a_2}{b_2} = \dots = \frac{a_n}{b_n}$$

using schwarz inequality, we get from $\textcircled{1}$

$$\rho^2(x, y) \leq 1$$

That is $|\rho(x, y)| \leq 1 \Rightarrow -1 \leq \rho(x, y) \leq 1$.

Hence correlation co-efficient cannot exceed unity numerically.

it always lies between -1 and $+1$.

if $r = +1$, the correlation is perfect and positive and if $r = -1$, correlation is perfect and negative.

properties of correlation co-efficient :-

(i) correlation co-efficient lies between 1 and -1 ; (i.e.) $-1 \leq r \leq 1$

(ii) correlation co-efficient is independent of change of scale. for eg, if the terms of a series are $25, 50, 100, 150, 200, 225$; we can simplify these terms are $1, 2, 4, 6, 8, 9$; calculate r easily. (Dividing by 25)

(iii) r is independent of unit of measurement.

(iv) If b_{xy} and b_{yx} are two regression co-efficients; correlation co-efficient is $\sqrt{b_{xy} \times b_{yx}}$.

(v) correlation coefficient works both ways, or $r_{xy} = r_{yx}$. That is we may take any series dependent and other as independent, its value remains the same.

Regression

Definition: Regression

Regression analysis is a mathematical we assume of the average relationship between two or more variables in terms of the original units of the data.

In Regression analysis there are two types of variables. The variable whose value is influenced or is to be predicted is called dependent variable and the variable which influences the values or is used for prediction is called independent variable. In regression analysis independent variable is also known as regressor or predictor or explanatory variable while the dependent variable is also known as regressed or explained variable.

Line of Regression:-

Definition: Curve of regression

If the variables in a bivariate distribution are related, we will find that the points in the scatter diagram will cluster round some curve called the

"Curve of regression".

Definition: Linear regression

If the curve is a straight line, it is called the line of regression and there is said to be linear regression between the variables, otherwise regression is said to be curvilinear.

Derive lines of Regression formula:-

Let us suppose that in the bivariate distribution $(x_i, y_i); i=1, 2, \dots, n$; y is dependent variable and x is independent variable. Let the line of regression of y on x be $Y = a + bX$.

According to the principle of least squares, the normal equations for estimating a and b are.

$$\sum_{i=1}^n y_i = na + b \sum_{i=1}^n x_i \quad \longrightarrow \textcircled{1}$$

$$\text{and } \sum_{i=1}^n x_i y_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \quad \longrightarrow \textcircled{2}$$

From $\textcircled{1}$ on dividing by n , we get

$$\bar{y} = a + b\bar{x} \quad \longrightarrow \textcircled{3}$$

Thus the line of regression of y on x passes through the point (\bar{x}, \bar{y}) .

Now

$$\mu_{11} = \text{cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$$

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n x_i y_i = \mu_{11} + \bar{x} \bar{y} \quad \rightarrow (4)$$

Also

$$\sigma_{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \Rightarrow \frac{1}{n} \sum_{i=1}^n x_i^2 = \sigma_{x^2} + \bar{x}^2$$

$\rightarrow (5)$

Dividing (3) by n and using (4) and (5), we get

$$\mu_{11} + \bar{x} \bar{y} = a \bar{x} + b (\sigma_{x^2} + \bar{x}^2) \rightarrow (6)$$

Multiplying (3) by \bar{x} and then subtracting from (6), we get

$$\mu_{11} = b \sigma_{x^2} \Rightarrow b = \frac{\mu_{11}}{\sigma_{x^2}} \quad \rightarrow (7)$$

Since 'b' is the slope of the line of regression of Y on X and since the line of regression passes through the point (\bar{x}, \bar{y}) , its equation is

$$Y - \bar{y} = b(X - \bar{x}) = \frac{\mu_{11}}{\sigma_{x^2}} (X - \bar{x}) \rightarrow (8)$$

$$Y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (X - \bar{x}) \quad \rightarrow (9)$$

Starting with the equation $X = A + BY$ and proceeding similarly or by simply interchanging the variables X and Y in (8) and (9),

the equation of the line of regression of x on y becomes

$$x - \bar{x} = \frac{r_{11}}{\sigma_y^2} (y - \bar{y}) \longrightarrow (10)$$

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y}) \longrightarrow (11)$$

Regression Coefficients :-

'b' the slope of the line of regression of y on x is also called the coefficient of regression of y on x .

b_{yx} = Regression coefficient of y on x

$$= \frac{r_{11}}{\sigma_x^2} = r \frac{\sigma_y}{\sigma_x}$$

Similarly, the coefficient of regression of x on y .

b_{xy} = Regression coefficient of x on y

$$= \frac{r_{11}}{\sigma_y^2} = r \frac{\sigma_x}{\sigma_y}$$

Properties of Regression Coefficient

property 1 :-

Correlation coefficient is the geometric mean between the regression coefficients.

proof:-

we know that

$$b_{xy} = r \frac{\sigma_x}{\sigma_y} \longrightarrow \textcircled{1}$$

$$\text{and } b_{yx} = r \frac{\sigma_y}{\sigma_x} \longrightarrow \textcircled{2}$$

From $\textcircled{1}$ and $\textcircled{2}$

$$\begin{aligned} b_{xy} \times b_{yx} &= r \frac{\sigma_x}{\sigma_y} \times r \frac{\sigma_y}{\sigma_x} \\ &= r^2 \end{aligned}$$

$$\text{(ie), } r^2 = b_{xy} \times b_{yx}$$

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

property 2:-

If one of the regression coefficients is greater than unity, the other must be less than unity

(or)

coefficient

If the correlation between two variables is 0 or ± 1 how will the regression lines be?

proof:-

Let one of the regression coefficients (say) b_{xy} be greater than unity, then we have to show that $b_{xy} < 1$.

$$\text{Now } b_{yx} > 1 \Rightarrow \frac{1}{b_{yx}} < 1 \quad \text{---} \textcircled{1}$$

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

$$\text{Hence } b_{xy} \leq \frac{1}{b_{yx}} < 1 \quad [\text{from } \textcircled{1}]$$

property 3:-

Arithmetic mean of the regression coefficients is greater than the correlation coefficient r , provided $r > 0$.

proof:-

We have to prove that

$$\frac{1}{2}(b_{yx} + b_{xy}) \geq r$$

$$\text{(or)} \quad \frac{1}{2} \left(r \frac{\sigma_y}{\sigma_x} + r \frac{\sigma_x}{\sigma_y} \right) \geq r$$

$$\text{(or)} \quad \frac{\sigma_y}{\sigma_x} + \frac{\sigma_x}{\sigma_y} \geq 2 \quad (\because r > 0)$$

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

That is $(\sigma_y - \sigma_x)^2 \geq 0$ which is always true, since the square of real quantities is ≥ 0 .

property 4:-

Regression coefficients are independent of the change of origin but not of scale.

proof:-

$$\text{Let } U = \frac{X-a}{h}$$

$$\text{and } V = \frac{Y-b}{k}$$

$$\Rightarrow X = a + hU$$

$$Y = b + kV.$$

where $a, b, h > 0$ and $k > 0$ are constants.

Then

$$\text{cov}(X, Y) = hK \text{cov}(U, V).$$

$$\sigma_{X^2} = h^2 \sigma_{U^2} \text{ and } \sigma_{Y^2} = k^2 \sigma_{V^2}$$

$$b_{yx} = \frac{r_{11}}{\sigma_{X^2}} = \frac{hK \text{cov}(U, V)}{h^2 \sigma_{U^2}}$$

$$= \frac{k}{h} \cdot \frac{\text{cov}(U, V)}{\sigma_{U^2}} = \frac{k}{h} b_{vu}$$

Similarly, we can prove that

$$b_{xy} = \left(\frac{h}{k} \right) b_{uv}$$

Derive the formula for the angle between two lines of regression:-

Equations of the lines of regression of Y on X , and X on Y are $Y - \bar{y} = r \cdot \frac{\sigma_y}{\sigma_x} (X - \bar{x})$ and $X - \bar{x} = r \cdot \frac{\sigma_x}{\sigma_y} (Y - \bar{y})$.

Slopes of these lines are $r \cdot \frac{\sigma_y}{\sigma_x}$ and $\frac{\sigma_y}{r\sigma_x}$ respectively. If θ is the angle between the two lines of regression then

$$\begin{aligned} \tan \theta &= \frac{\frac{r\sigma_y}{\sigma_x} - \frac{\sigma_y}{r\sigma_x}}{1 + \frac{r\sigma_y}{\sigma_x} \cdot \frac{\sigma_y}{r\sigma_x}} \\ &= \frac{r^2 - 1}{r} \left(\frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \\ &= \frac{1 - r^2}{r} \left(\frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \end{aligned}$$

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

① Given $\text{var}(x) = 9$

problems based on regression

regression equation $8x - 10y + 66 = 0$ and

$40x - 18y = 214$. To Find

(i) find mean value of x and y
 (ii) the correlation coefficient between x and y .

(iii) S.D of y when variance $(X) = 9$.

Soln:-

Given: $8x - 10y + 66 = 0$, $40x - 18y = 214$

Since, both the lines of regression pass through the point (\bar{x}, \bar{y})

(i)

$$8x - 10y = -66$$

$$40x - 18y = 214$$

let $x = \bar{x}$

$y = \bar{y}$

$$8\bar{x} - 10\bar{y} = -66 \quad \text{--- (1)}$$

$$40\bar{x} - 18\bar{y} = 214 \quad \text{--- (2)}$$

$$\text{(1)} \times 5 \Rightarrow 40\bar{x} - 50\bar{y} = -330$$

$$\text{(2)} \Rightarrow 40\bar{x} - 18\bar{y} = 214$$

$$\begin{array}{r} 40\bar{x} - 50\bar{y} = -330 \\ 40\bar{x} - 18\bar{y} = 214 \\ \hline -32\bar{y} = -544 \end{array}$$

$$\bar{y} = \frac{544}{32}$$

$$\bar{y} = 17$$

$\bar{y} = 17$ sub in (1)

$$8\bar{x} - 10\bar{y} = -66$$

$$8\bar{x} - 10(17) = -66$$

$$8\bar{x} - 170 = -66$$

$$8\bar{x} = -66 + 170$$

$$8\bar{x} = 104$$

$$\bar{x} = \frac{104}{8}$$

$$\bar{x} = 13$$

$$(ii) \textcircled{2} \Rightarrow 40\bar{x} - 18\bar{y} = 214$$

$$40\bar{x} = 214 + 18\bar{y}$$

$$\bar{x} = \frac{214}{40} + \frac{18}{40}\bar{y}$$

$$\bar{x} = 5.35 + 0.45\bar{y}$$

$$b_{xy} = 0.45 \text{ (ve)}$$

$$\textcircled{1} \Rightarrow 8\bar{x} - 10\bar{y} = -66$$

$$-10\bar{y} = -66 - 8\bar{x}$$

$$10\bar{y} = 66 + 8\bar{x}$$

$$\bar{y} = \frac{66}{10} + \frac{8}{10}\bar{x}$$

$$\bar{y} = 6.6 + 0.8\bar{x}$$

$$b_{yx} = 0.8 \text{ (ve)}$$

$$\therefore r = \pm \sqrt{b_{xy} \cdot b_{yx}} = \pm \sqrt{0.45 \times 0.8}$$

$$= \pm \sqrt{0.36} = \pm 0.6$$

$$\therefore r = 0.6$$

(iii) To find S.D of y

$$\text{Given : var}(X) = 9 = \sigma_x^2$$

$$\sigma_x = \sqrt{9} = 3$$

$$\text{w.k.T } b_{xy} = r \cdot \frac{\sigma_x}{\sigma_y}$$

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

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

$$0.45 = 0.6 \times \frac{3}{\sigma_y}$$

$$0.45 \times \sigma_y = 1.8$$

$$\sigma_y = \frac{1.8}{0.45} = 4$$

$$\text{S.D of } y = 4$$

$$\sigma_y = 4$$

Formula :-

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

$$b_{yx} = \frac{N \sum xy - \sum x \sum y}{N \sum x^2 - (\sum x)^2}$$

① Obtain the equation of two lines of regression for the following data and also find estimate of x for $y = 23$

X	10	12	14	16	18	20	22	24
Y	14	18	16	22	26	28	27	30

$$\bar{x} = 17$$

$$N = 8$$

$$\bar{y} = 22.625$$

$$b_{xy} = 0.776$$

$$b_{yx} = 1.1726$$

$$x \text{ on } y \Rightarrow x = 0.776y - 0.557$$

$$y \text{ on } x \Rightarrow y = 1.1726x + 2.6908$$