

Unit-2

Correlation and Linear Regression

- Correlation:- Quantitative measure of relationship between two variables.
- Positive/Direct correlation:- If the increase/decrease in one variable results in increase/decrease in another variable.
Ex. Income & Expenditure.

- Negative/Diverse correlation:- If the increase/decrease in one variable results in decrease/increase in another variable.
Ex. Price & Demand.

- Karl Pearson's coefficient of correlation.

~~Correlation~~ Correlation coefficient between two random variables X and Y is denoted by $r(X, Y)$ and is numerical measure of relationship b/w X and Y .

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\sigma_{XY}}{\sigma_X \sigma_Y}$$

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

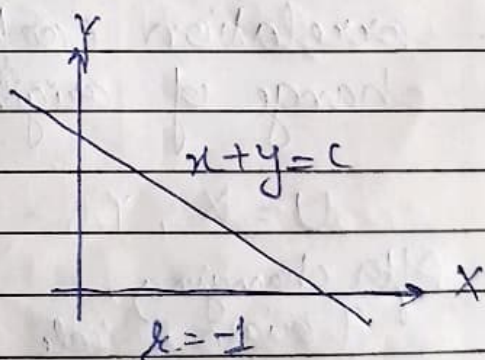
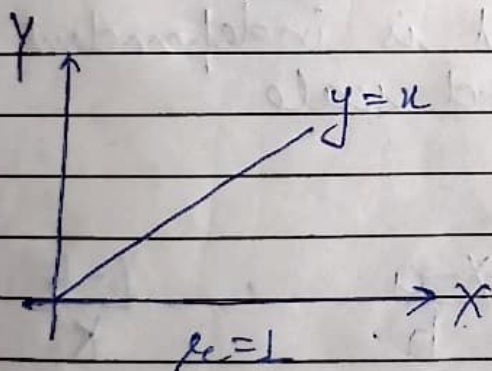
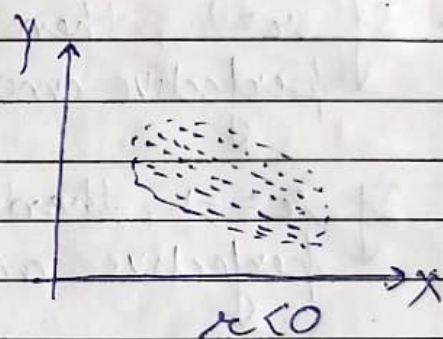
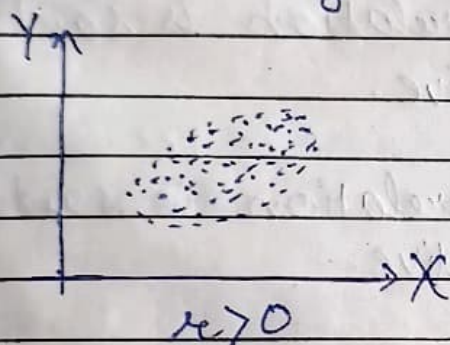
$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} \quad \bar{y} = \frac{\sum_{i=1}^n y_i}{n}$$

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

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

X	Y	(X, Y)
x_1	y_1	(x_1, y_1)
x_2	y_2	(x_2, y_2)
x_3	y_3	(x_3, y_3)
\vdots	\vdots	\vdots
x_i	y_i	(x_i, y_i)
\vdots	\vdots	\vdots
x_n	y_n	(x_n, y_n)

• Scatter Diagram:



Page: 1

* Properties of Correlation :-

1. $r(X, Y)$ is a measure of linear relationship between X and Y .
2. Karl Pearson's correlation is also called product moment correlation coefficient.
3. $-1 \leq r(X, Y) \leq 1$
4. If $r=1$, then the correlation is said to be perfective and positive.
5. If $r=-1$, then the correlation is said to be perfective and negative.
6. Correlation coefficient is independent change of origin and scale.

$$U = X, Y$$

After changing,
of origin & scale

$$U = \frac{X-a}{h}, \quad \frac{Y-b}{k}$$

~~After changing~~

7. If X and Y are random variables & a, b, c, d are any numbers provided $a \neq 0$ & $c \neq 0$ then,

$$r(aX+b, cY+d) = \frac{ac}{|ac|} r(X, Y)$$

$$r(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$$

8. If X and Y are independent then $r(X, Y) = 0$ but converse is not true.

9. Calculate the coefficient of correlation b/w X & Y .

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

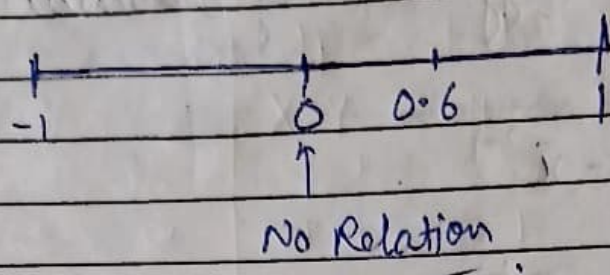
$$\sigma_y^2 = \frac{1}{n} \left(\sum y_i^2 \right) - \bar{y}^2$$

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} = \frac{\frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}}{\sqrt{\left(\frac{1}{n} \sum x_i^2 - \bar{x}^2 \right) \left(\frac{1}{n} \sum y_i^2 - \bar{y}^2 \right)}}$$

(Ht. of Father)	(Ht. of Son)				
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	4839	
72	71	5184	5041	5112	
$\sum X = 544$	$\sum Y = 552$	$\sum X_i^2$	$\sum Y_i^2$	$\sum X_i Y_i$	
$\bar{x} = 68$	$\bar{y} = 69$	\downarrow	\downarrow	\downarrow	
		37628	38132	37560	

$$\therefore r(X, Y) = \frac{\frac{1}{8} \times 37560 - 68 \times 69}{\sqrt{\left(\frac{1}{8} \times 37628 - (68)^2 \right) \left(\frac{1}{8} \times 38132 - (69)^2 \right)}} = 0.603$$

1. X Y
 height of Father height of sons



* Spearman's Rank correlation :-

Let (x_i, y_i) $i=1, 2, 3, \dots, n$ be the rank of i^{th} individual in two characteristics A and B.

then,

$$\rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2-1)}$$

not sigma

X	Y
(Characteristic A)	(Characteristic B)
x_1	y_1
x_2	y_2
x_3	y_3
\vdots	\vdots

where $d_i = x_i - y_i$

we always have $\sum_{i=1}^n d_i = 0$

Example 1: The ranks of 16 students in maths + Physics are given below. Calculate rank correlation coefficient.

X (Rank in Math)	Y (Rank in Physics)	$d_i = x_i - y_i$	d_i^2
1	1	0	0
2	10	-8	64
3	3	0	0
4	4	0	0
5	5	0	0
6	7	-1	1
7	2	5	25
8	6	2	4
9	8	1	1
10	11	-1	1
11	15	-4	16
12	9	3	9
13	14	-1	1
14	12	2	4
15	16	-1	1
16	13	3	9

$\sum d_i = 0$ $\sum d_i^2 = 136$

Should be zero always.

Here, $n = 16$

$$r = \frac{1 - 6 \cdot \sum_{i=1}^n d_i^2}{n(n^2 - 1)} = \frac{1 - 6 \times 136}{16(16^2 - 1)}$$

$r = 0.8$

This means High positive correlation b/w X & Y

Ex. Ten participants in a musical test, ~~X, Y, Z~~ were ranked by three judges X, Y, Z. discuss which pair of judges have common liking in music.

			d_1	d_2	d_3	d_1^2	d_2^2	d_3^2
X	Y	Z	X-Y	Y-Z	X-Z			
1	3	6	-2	-3	-5	4	9	25
6	5	4	1	1	2	1	1	4
5	8	9	-3	-1	-4	9	1	16
10	4	8	6	-4	2	36	16	4
3	7	1	-4	6	2	16	36	4
2	10	2	-8	8	0	64	64	0
4	2	3	2	-1	1	4	1	1
9	1	10	8	-9	-1	64	81	1
7	6	5	1	1	2	1	1	4
8	9	7	-1	2	1	1	4	1
			0	0	0	200	214	60

We have to find ρ_{xy} , ρ_{yz} , ρ_{xz} & check which of these ~~have~~ have ~~same~~ (positive) sign.

$$\rho_{xy} = 1 - 6 \cdot \frac{\sum_{i=1}^n d_1^2}{n(n^2-1)} = 1 - \frac{6 \times 200}{10(10^2-1)} = -\frac{7}{33}$$

$$\rho_{yz} = 1 - 6 \cdot \frac{\sum_{i=1}^n d_2^2}{n(n^2-1)} = 1 - \frac{6 \times 214}{10(10^2-1)} = -\frac{49}{165}$$

$$\rho_{xz} = 1 - 6 \cdot \frac{\sum_{i=1}^n d_3^2}{n(n^2-1)} = 1 - \frac{6 \times 60}{10(10^2-1)} = \frac{7}{11}$$

So, X & Z are having common liking

* Rank correlation in case of tied ranks or repeated ranks:-

$$\rho = 1 - \frac{6 \left(\sum_{i=1}^n d_i^2 + \text{Correction factor} \right)}{n(n^2 - 1)}$$

where correction factor = $\frac{m(m^2 - 1)}{12}$, where m is number of times an item repeated

Ex. Calculate the rank correlation for following data

X	Y	Rank(X)	Rank(Y)	d(X-Y)	d 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	25
80	60	1	6	-5	25
75	68	2.5	3.5	-1	1
40	48	10	9	1	1
55	50	8	8	0	0
64	70	6	2	4	16

$$\text{Correction factor} = \frac{m_1(m_1^2 - 1)}{12} + \frac{m_2(m_2^2 - 1)}{12} + \frac{m_3(m_3^2 - 1)}{12}$$

(75) $m_1 = 2 \Rightarrow 75 \Rightarrow 2 \text{ times}$

(64) $m_2 = 3 \Rightarrow 64 \Rightarrow 3 \text{ times}$

(88) $m_3 = 2 \Rightarrow 68 \Rightarrow 2 \text{ times}$

For 11 as
1 + 2 + 3
= 2.5

1
2
3
4
5
6
7
8
9
10

* Linear Regression:

* Regression: It is the study of nature of relationship between the variables so that one is able to predict the value of one variable on the basis of other. In regression analysis one variable is dependent and another is independent.

* Linear Regression: If the variables in the bivariate distribution are related, we will find that points (x_i, y_i) $i=1, 2, \dots, n$ cluster around some curve, that curve is called curve of Regression.

If the curve is straight line, then it is called the linear regression.

Suppose we have n data points (x_i, y_i) , $i=1, 2, \dots, n$

Let the line of regression of Y on X be

$$y = a + bx$$

And line of regression of X on Y be

$$x = c + dy$$

We have line of regression of Y on X

$$y = a + bx \quad \text{--- (1)}$$

$$\sum y = na + b \sum x \quad \text{--- (2)}$$

$$y_1 = a + bx_1$$

$$y_2 = a + bx_2$$

$$\vdots$$

$$y_n = a + bx_n$$

$$\sum y_i = na + b \sum x_i$$

Multiply eqⁿ (1) by x & take summation,

$$\sum xy = a \sum x + b \sum x^2 \quad \text{--- (3)}$$

Eqⁿ (2) & (3) are called normal equations of Y on X and are used to determine values of a & b .

Multiply eqⁿ (2) by $\frac{1}{n}$

$$\frac{1}{n} \sum y = a + b \frac{1}{n} \sum x$$

$$\boxed{\bar{y} = a + b \bar{x}}$$

Line of regression of Y on X passes through (\bar{x}, \bar{y}) .

Now line of regression of X on Y ,

$$x = c + dy \quad \text{--- (4)}$$

$$\sum x = nc + d \sum y \quad \text{--- (5)}$$

Multiply (4) by y & take summation,

$$\sum xy = c \sum y + d \sum y^2 \quad \text{--- (6)}$$

Equation (5) & (6) are normal equation of X on Y ,

Multiply eqⁿ (5) by $\frac{1}{n}$.

$$\frac{1}{n} \sum x = c + d \frac{1}{n} \sum y$$

$$\bar{x} = c + d\bar{y}$$

→ The line of regression of X on Y , also passes through (\bar{x}, \bar{y}) .

Q. Find the line of regression for the following data.

X	Y	X ²	Y ²	XY
1	8	1	64	8
2	7	4	49	14
3	5	9	25	15
4	9	16	81	36
5	11	25	121	55
15	40	55	340	128

Let the line of regression of Y on X be $y = a + bx$ the normal equation of Y on X are:

$$\sum y = na + b \sum x$$

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

$$\begin{cases} 40a = 5a + 15b \\ 128 = 15a + 55b \end{cases} \Rightarrow a = \frac{28}{5}, b = \frac{4}{5}$$

\therefore line of regression of Y on X is :-

$$y = \frac{28}{5} + \frac{4}{5}x$$

Let the line of regressions of X on Y be $x = c + dy$
then Normal equations are :-

$$\begin{aligned}\sum x &= nc + d \sum y \\ \sum ny &= c \sum y + d \sum y^2\end{aligned}$$

$$15 = 5c + 40d$$

$$128 = 40c + 340d$$

$$c = \frac{-1}{5}, \quad d = \frac{2}{5}$$

Equation of line of regression of X on Y is

$$x = \frac{-1}{5} + \frac{2}{5}y$$

(3,8)

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

* Another form of regression lines;

The line of regression of Y on X is given by

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

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

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

Also line of regression of X on Y is given by

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

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

where r is correlation coefficient of X & Y

σ_x is standard deviation of X

σ_y is standard deviation of Y

And b_{xy} & b_{yx} are called regression coefficients.

* Properties of Regression Coefficient

1. Correlation coefficient r of X & Y is geometric mean between regression coefficient. The sign of r , b_{xy} , b_{yx} is same.

$$r^2 = b_{xy} b_{yx}$$

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

2. If one of the regression coefficient is greater than unity, the other must be less than unity.

$$-1 \leq r \leq 1$$

3.
$$\left| \frac{b_{xy} + b_{yx}}{2} \right| \geq |r|$$

$$(A.M. \geq G.M.)$$

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

NOTE Angle b/w two lines (straight) :

$$y_1 = m_1 x + c_1$$

$$y_2 = m_2 x + c_2$$

$$\tan \theta = \left| \frac{m_1 - m_2}{1 + m_1 m_2} \right|$$

* Angle between two lines of Regression:

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

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

then,

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

Case-1: If $r=0$, then $\theta = \frac{\pi}{2}$

This means if X and Y are uncorrelated then lines of regression are perpendicular.

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

In case of perfect correlation positive or negative the lines of regression coincide.

Example: For variables X and Y ,
 $\sigma_x^2 = 9$,

Regression equations:

$$8X - 10Y + 66 = 0 \quad (Y \text{ on } X)$$

$$40X - 18Y = 214 \quad (X \text{ on } Y)$$

Find (1) mean value of X and Y

(2) Correlation coefficient b/w X & Y

(3) Standard deviation of Y .

Qⁿ. Given $\sigma_x^2 = 9$
 $\sigma_x = 3$

(1) Since (\bar{x}, \bar{y}) is point of intersection of lines of regression

$$\therefore 8\bar{x} - 10\bar{y} + 66 = 0$$

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

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

(2)

From (1) & (2),

$$Y = \frac{8}{10}x + \frac{66}{10}$$

$$X = \frac{18}{40}y + \frac{214}{40}$$

We get, $b_{yx} = \frac{8}{10}$, $b_{xy} = \frac{18}{40}$

Also, $r = \pm \sqrt{b_{xy} b_{yx}}$

$$r = \pm \sqrt{b_{xy} b_{yx}} = \pm \sqrt{\frac{8}{10} \times \frac{18}{40}}$$

$$r = \pm \frac{3}{5}$$

Since r, b_{xy}, b_{yx} have same sign -

$$\therefore r = +\frac{3}{5}$$

(3)

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

$$\sigma_x = 3 \text{ (Given)}$$

$$\sigma_y = \frac{40}{10}$$

$$\Rightarrow \sigma_y = 4$$

$$\frac{8}{10} = \frac{3}{5} \times \frac{\sigma_y}{\sigma_x}$$

Example: Find the most likely price in Mumbai corresponding to the prices of Re. 70 in Calcutta from the following data.

	X Calcutta	Y Mumbai
Avg. Price	65	67
S.D (Standard deviation)	2.5	3.5

Correlation between prices in two cities is 0.8.

Solⁿ $X = 70$ $Y = ?$

The line of regression of Y on X is given by

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

$$\bar{x} = 65$$

$$\bar{y} = 67$$

$$\sigma_x = 2.5$$

$$\sigma_y = 3.5$$

Line of regression of Y on X is:

$$Y - 67 = 0.8 \times \frac{3.5}{2.5} (X - 65)$$

At $X = 70$,

$$Y - 67 = 0.8 \times \frac{3.5}{2.5} (70 - 65)$$

$$Y = 72.6$$

$$1. \quad r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\sigma_{XY} = \text{Cov}(X, Y) = \frac{1}{n} \left(\sum_{i=1}^n x_i y_i - \bar{x} \bar{y} \right)$$

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

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

$$\sigma_Y = \frac{1}{n} \left(\sum_{i=1}^n y_i^2 - (\bar{y})^2 \right)$$

$$2. \quad \rho = 1 - 6 \cdot \frac{\sum_{i=1}^n d_i^2}{n(n^2-1)}$$

$$d_i = x_i - y_i, \quad \sum_{i=1}^n d_i = 0$$

$$3. \quad \rho = 1 - \frac{6 \left(\sum_{i=1}^n d_i^2 + \text{correction factor} \right)}{n(n^2-1)}$$

$$\text{correction factor} = \frac{m(m^2-1)}{12}$$

no. of times an item repeated

4. Line of Regression of Y on X be $y = a + bx$ & X on Y $x = c + dy$

Normal equations: Y on X: Normal eq^{ns}: X on Y

$$\sum y = na + b \sum x$$

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

$$\sum x = nc + d \sum y$$

$$\sum xy = c \sum y + d \sum y^2$$

$$5. \quad (y - \bar{y}) = r \frac{\sigma_Y}{\sigma_X} (x - \bar{x})$$

$$b_{yx} = r \frac{\sigma_Y}{\sigma_X}$$

$$(x - \bar{x}) = r \frac{\sigma_X}{\sigma_Y} (y - \bar{y})$$

$$b_{xy} = r \frac{\sigma_X}{\sigma_Y}$$

b_{yx} & b_{xy} are regression coefficients

6. Angle b/w two regression lines:-

$$\theta = \tan^{-1} \left| \frac{1-r^2}{|r|} \cdot \frac{\sigma_X \sigma_Y}{\sigma_X^2 + \sigma_Y^2} \right|$$