



Correlation and Regression

Correlation

Introduction: → If the Distribution involving only one variable it is called Univariate distribution.

→ If the Distribution involving two variables it is called Bivariate distribution.

→ If the Distribution involving three variables it is called trivariate distribution.

If the distribution involving more than three variables are called Multivariate distribution.

In our Correlation we have to study relationship between two variables.

Q → Define Correlation? Types of Correlation? Define Karl Pearson Coefficient of Correlation.

Def: → If change in one variable effects to change in another variable then the two variables are said to be correlated then that relationship is called Correlation.



Correlation are two types.

- ① Positive Correlation (or) Direct Correlation
- ② Negative Correlation (or) Indirect Correlation

Positive Correlation:- If the two variables are deviated in same direction is called (or) Direct (or) Positive Correlation. i.e. If one variable is increases then another variable is also increases.

Negative (or) Indirect Correlation:-

If the two variables are deviated in opposite direction is called Negative (or) Indirect Correlation. i.e. If one variable is increases then another variable is decreases.

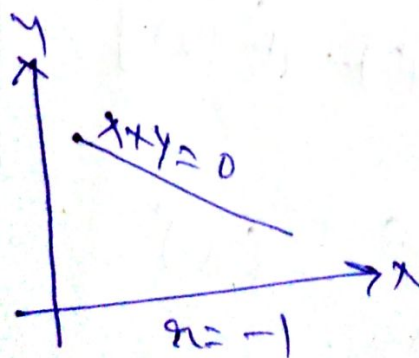
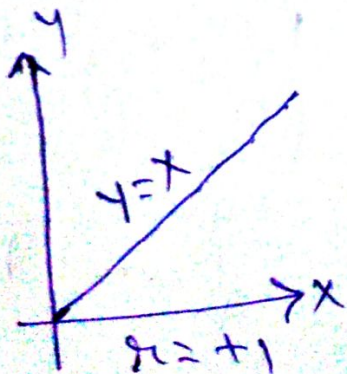
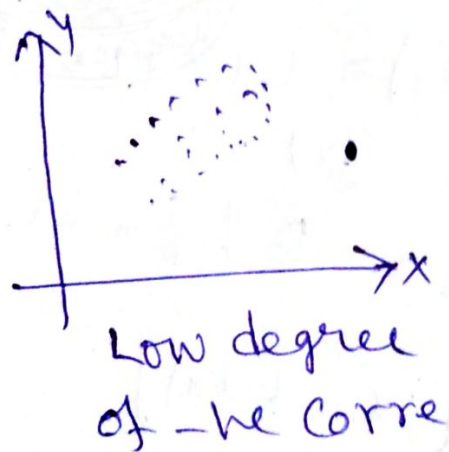
→ Scattered diagram:-

The diagrammatic representation of the Bivariate data is called

"Scattered diagram"

Using Scattered diagram we can get fairly good idea about the distribution.

If the values x and y plotted among xy plane. The diagram of dots obtained we can get fairly good idea. If the variables are correlated (or) not. If the points are scattered - very close each other. Then there is high degree of correlation. and if the points are widely scattered then there is poor correlation.





Karl Pearson Coefficient of Correlation :-

As a major intensity of association between two variables

Karl Pearson a bio-mathematician and statistician developed a formulae is called Karl Pearson Coefficient of Correlation. It is denoted with r_{xy}

defined as

$$r_{xy} = \frac{\text{Cov}(xy)}{\sigma_x \cdot \sigma_y} \quad (\text{or}) \quad r_{xy} = \frac{\text{Cov}(xy)}{\sqrt{V(x)V(y)}}$$

Where $\text{Cov}(xy) = \frac{\sum x_i y_i}{n} - (\bar{x})(\bar{y})$

(or)

$\text{Cov}(xy) = E(xy) - E(x) \cdot E(y)$

(or)

$\text{Cov}(xy) = E[\{x - E(x)\} \{y - E(y)\}]$

$V(x) = E[x - E(x)]^2$

(or)

$V(x) = E(x^2) - [E(x)]^2$

(or)

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

$V(y) = E[y - E(y)]^2$

(or)

$V(y) = E(y^2) - [E(y)]^2$

(or)

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



⇒ Properties of Karl Pearson Coefficient of Correlation

- ① Correlation Coefficient lies between -1 to $+1$ (or) $-1 \leq r_{xy} \leq +1$
- ② There is no origin and scale effect on correlation.

Theorem: — Prove limits of Correlation
(or)
Prove $-1 \leq r_{xy} \leq +1$

Proof:- Karl Pearson Coefficient of Correlation $r_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{V(x) \cdot V(y)}}$

~~Cov(x, y) = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})~~

Where

$$\text{Cov}(x, y) = E\left[\{x - E(x)\}\{y - E(y)\}\right]$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum [\{x_i - \bar{x}\}\{y_i - \bar{y}\}]$$

$$V(x) = E[x - E(x)]^2$$

$$V(y) = E[y - E(y)]^2$$

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

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

$$r_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{V(x) \cdot V(y)}}$$

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

$$r_{xy} = \frac{\frac{1}{n} \sum (x_i - \bar{x}) (y_i - \bar{y})}{\frac{1}{n} \sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}}$$

let $x_i - \bar{x} = a_i$, $y_i - \bar{y} = b_i$

$$r_{xy} = \frac{\sum a_i b_i}{\sqrt{\sum a_i^2 \cdot \sum b_i^2}}$$

Squaring on both sides.

$$r_{xy}^2 = \frac{[\sum a_i b_i]^2}{\sum a_i^2 \cdot \sum b_i^2}$$

According to Cauchy Schwartz Inequality

$$[E(xy)]^2 \leq E(x^2) \cdot E(y^2)$$

$$[\sum a_i b_i]^2 \leq \sum a_i^2 \cdot \sum b_i^2$$

$$r_{xy}^2 \leq 1$$

$$\rho_{xy} \leq +1$$

$$\boxed{-1 \leq \rho_{xy} \leq +1}$$

Hence the proof

→ Prove that there is no origin and scale effect on correlation.

Proof:- Karl Pearson Coefficient of Correlation

$$\rho_{xy} = \frac{\text{Cov}(xy)}{\sqrt{V(x) \cdot V(y)}}$$

$$\rho_{xy} = \frac{E[\{x - E(x)\}\{y - E(y)\}]}{\sqrt{E[x - E(x)]^2 \cdot E[y - E(y)]^2}}$$

To change of origin and scale in x and y variables.

$$U = \frac{x - a}{h}$$

$$V = \frac{y - b}{k}$$

$$Uh = x - a$$

$$x = a + Uh$$

$$E(x) = E[a + Uh]$$

$$= E(a) + E(Uh)$$

$$= a + h E(U)$$

$$Vk = y - b$$

$$y = Vk + b$$

Taking Expectation

$$E(y) = E(Vk + b)$$

$$E(y) = E(Vk) + E(b)$$

$$E(\bar{y}) = k E(\bar{v}) + b$$



$$\begin{aligned}
 \therefore r_{xy} &= \frac{E[\{x - E(x)\} \{y - E(y)\}]}{\sqrt{[E\{x - E(x)\}^2] [E\{y - E(y)\}^2]}} \\
 &= \frac{E[h\{u - E(u)\} k\{v - E(v)\}]}{\sqrt{E[h\{u - E(u)\}^2] \cdot E[k\{v - E(v)\}^2]}} \\
 &= \frac{hk E[\{u - E(u)\} \{v - E(v)\}]}{\sqrt{h^2 E[u - E(u)]^2 \cdot k^2 E[v - E(v)]^2}} \\
 &= \frac{hk E[\{u - E(u)\} \{v - E(v)\}]}{hk \sqrt{E[u - E(u)]^2 \cdot E[v - E(v)]^2}} \\
 &= r_{uv}
 \end{aligned}$$

$$\therefore r_{xy} = r_{uv}$$

\therefore There is no origin and scale effect on correlation.



Problem

① Compute Karl Pearson Coefficient between x and y for the given data.

x : 1 2 3 4 5

y : 10 20 30 40 50

Sol — Karl Pearson Coefficient of Correlation.

$$r_{xy} = \frac{\text{Cov}(xy)}{\sigma_x \cdot \sigma_y}$$

where $\text{Cov}(xy) = \frac{\sum x_i y_i}{n} - (\bar{x})(\bar{y})$

$$\sigma_x = \sqrt{\frac{\sum x_i^2}{n} - (\bar{x})^2}, \quad \sigma_y = \sqrt{\frac{\sum y_i^2}{n} - (\bar{y})^2}$$

$$\bar{x} = \frac{\sum x_i}{n}, \quad \bar{y} = \frac{\sum y_i}{n}$$

x_i	y_i	$\sum x_i y_i$	$\sum x_i^2$	$\sum y_i^2$
1	10	10	1	100
2	20	40	4	400
3	30	90	9	900
4	40	160	16	1600
5	50	250	25	2500
<u>15</u>	<u>150</u>	<u>550</u>	<u>55</u>	<u>5500</u>

$$\bar{x} = \frac{\sum x_i}{n} = \frac{15}{5} = 3$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{150}{5} = 30$$

$$\text{Cov}(xy) = \frac{\sum x_i y_i}{n} - (\bar{x})(\bar{y})$$

$$= \frac{550}{5} - (3)(30)$$

$$= 110 - 90$$

$$\boxed{\text{Cov}(xy) = 20}$$

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

$$= \sqrt{\frac{55}{5} - (3)^2}$$

$$= \sqrt{2}$$

$$\boxed{\sigma_x = 1.4142}$$

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

$$= \sqrt{\frac{5500}{5} - (30)^2}$$

$$= \sqrt{200}$$

$$\boxed{\sigma_y = 14.1421}$$

$$r = \frac{\text{Cov}(xy)}{\sigma_x \cdot \sigma_y} = \frac{20}{(1.4142)(14.1421)}$$

$$= \frac{20}{19.9997} = 1$$

$$\boxed{r = 1}$$

Problem: Exam Papers,

① Find correlation coefficient between industrial production and export using the following data comment on result.

Production (X): 55 56 58 59 60 60 62
 Export (Y) : 35 38 38 39 44 43 45

② Calculate coefficient of correlation from the following data

X: 50 60 70 90 100
 Y: 65 51 40 26 8

Rank Correlation Coefficient →

Def:- Arrangement of given data in to order of merit is called 'Rank Correlation Coefficient'.

Some times we come across statistical series in which the variables are not capable of qualitative measurement can be arranged in to the serial order. In that case Karl Pearson Coefficient of correlation does

not used. So we use Spearman's rank correlation formulae.

Spearman's rank correlation coefficient

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

where, $d_i = R_i - R_y$
 n = number of pairs.

Theorem obtain rank correlation coefficient formulae.

Proof:- Assuming that no two individuals get a same rank.

let, x_i is a rank of Attribute - A.

y_i is a rank of Attribute - B.

$$\bar{x} = \frac{\sum x_i}{n} = \frac{1+2+3+\dots+n}{n} =$$

$$= \frac{\frac{n(n+1)}{2}}{n} = \frac{n+1}{2}$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{1+2+3+\dots+n}{n} = \frac{\frac{n(n+1)}{2}}{n} = \frac{n+1}{2}$$

$$\therefore \bar{x} = \bar{y}$$



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

$$= \frac{1^2 + 2^2 + 3^2 + \dots + n^2}{n} - \left(\frac{n+1}{2}\right)^2$$

$$= \frac{n(n+1)(2n+1)}{6} - \left(\frac{n+1}{2}\right)^2$$

$$\sigma_x^2 = \frac{(n+1)(2n+1)}{6} - \frac{(n+1)^2}{4}$$

$$\sigma_x^2 = \frac{n+1}{2} \left[\frac{2n+1}{3} - \frac{n+1}{2} \right]$$

$$= \frac{n+1}{2} \left[\frac{2(2n+1) - 3(n+1)}{6} \right]$$

$$= \frac{n+1}{2} \left[\frac{4n+2 - 3n-3}{6} \right]$$

$$= \frac{n+1}{2} \left[\frac{n-1}{6} \right]$$

$$\sigma_x^2 = \frac{n^2-1}{6}$$

$$\therefore \sigma_y^2 = \frac{n^2-1}{6}$$

$$\therefore \sigma_x^2 = \sigma_y^2$$

~~also $d_x = R_x - R_y$~~ also $d_i = x_i - y_i$

$$d_i = x_i - \bar{x} + \bar{x} - y_i \Rightarrow d_i = x_i - \bar{x} + \bar{y} - y_i$$

$$d_i = (x_i - \bar{x}) - (y_i - \bar{y}) \quad [\because \bar{x} = \bar{y}]$$

Taking \sum divide by n and simplifying on both sides.



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

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

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

$$= \cancel{\sigma_x^2} + \sigma_y^2 - 2 \text{Cov}(xy)$$

$$= \sigma_x^2 + \sigma_y^2 - n \cdot \sigma_x \cdot \sigma_y$$

$$r = \frac{\text{Cov}(xy)}{\sigma_x \sigma_y}$$

$$\text{Cov}(xy) = r \cdot \sigma_x \sigma_y$$

$$\frac{\sum d_i^2}{n} = \sigma_x^2 + \sigma_y^2 - 2 \sigma_x \cdot \sigma_y \quad (\because \sigma_x = \sigma_y)$$

$$= 2\sigma_x^2 - 2\sigma_x^2 r$$

$$= 2\sigma_x^2 [1 - r]$$

$$\frac{\sum d_i^2}{n} = 2 \left[\frac{n^2 - 1}{12} \right] [1 - r]$$

$$\therefore \sigma_x^2 = \frac{n^2 - 1}{12}$$

$$\sum d_i^2 = \frac{n(n^2 - 1)}{6} (1 - r)$$

$$\frac{6 \sum d_i^2}{n(n^2 - 1)} = 1 - r$$

$$1 - \frac{6 \sum d_i^2}{n(n^2 - 1)} = r$$

$$\therefore P = r$$

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

→ Find Rank Correlation Coefficient between Marks of English and Computers

Marks in Eng: 48 50 47 55 43 60 68 45 52 64 58
56 72 40

Marks in Comp: 43 52 48 53 68 50 64 51 44 72
62 49 70 65

Sol: Let us denote marks in English is variable x .
marks in Computers is variable y .

x	Rank R_x	y	Rank R_y	$d_i = R_x - R_y$	d_i^2
48	10	43	14	10-14 = -4	16
50	9	52	8	9-8 = 1	1
47	11	48	12	11-12 = -1	1
55	7	53	7	0	0
43	13	68	3	10	100
60	4	50	10	-6	36
68	2	64	5	-3	9
45	12	51	9	3	9
52	8	44	13	-5	25
64	3	72	1	-5	25
58	5	62	6	-1	1
56	6	49	11	2	4
72	1	70	2	-1	1
40	14	65	4	-5	25
				10	100
					<u>328</u>



Spearman's Correlation Coefficient

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

$$r = 1 - \left[\frac{6(328)}{14(195)} \right]$$

$$r = 1 - \frac{1968}{2730}$$

$$= 1 - 0.7208$$

$$r = 0.2792$$

Problem 1 Find Rank Correlation Coefficient between Marks of Maths and Statistics.

Marks in Maths :	92	96	72	88	54	100	75	93	47	35
Marks in Stat :	100	47	74	94	75	99	72	89	50	40

Problem 2 The Coefficient of Rank Correlation of Marks obtained by 10 Students in Maths and Statistics was found to be 0.5. It was later discovered that the difference in marks two subjects obtained by one of the student was wrongly taken as 3 instead of 7. Find the correct Rank Correlation.



→ The following are the ranks assigned by two judges X and Y to 12 contestants in cooking competition. Find out what agreement the judges had in judgements.

Entry No:	A	B	C	D	E	F	G	H	I	J	K	L
Rank Judge-A:	1	9	2	10	3	11	8	4	12	7	5	6
Rank Judge-B:	2	9	1	7	4	10	8	3	12	6	5	11

Sol- Let us denote the ranks given by judge A is variable X.

Let us denote the ranks given by Judge-B is variable Y.

Spearman's rank correlation coefficient

$$r = 1 - \left[\frac{6 \sum d_i^2}{n(n^2-1)} \right]$$
 where $d_i = R_x - R_y$
 $n = \text{number of pairs}$

R_x	R_y	$d_i = R_x - R_y$	$\sum d_i^2$
1	2	1-2 = -1	1
9	9	9-9 = 0	0
2	1	2-1 = 1	1
10	7	10-7 = 3	9
3	4	3-4 = -1	1
11	10	11-10 = 1	1
8	8	8-8 = 0	0
4	3	4-3 = 1	1
12	12	12-12 = 0	0
7	6	7-6 = 1	1
5	5	5-5 = 0	0
6	11	6-11 = -5	25
			$\sum d_i^2 = 40$

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

$$r = 1 - \left[\frac{6(40)}{12(12^2-1)} \right]$$

$$= 1 - \left[\frac{240}{12(143)} \right]$$

$$= 1 - \frac{240}{1716}$$

$$= 1 - 0.1398$$

$r = 0.8602$

X	Y	Z	$d_{xy} = \frac{R_x - R_y}{5}$	$d_{yz} = \frac{R_y - R_z}{5}$	$d_{zx} = \frac{R_z - R_x}{5}$	d_{xy}^2	d_{yz}^2	d_{zx}^2
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	-2	64	81	4
7	6	5	1	1	1	1	1	1
8	9	7	-1	2	1	1	4	1
						200	214	60

$$r_{xy} = 1 - \left[\frac{6 \sum d_{xy}^2}{n(n^2-1)} \right]$$

$$= 1 - \left[\frac{6(200)}{10(10^2-1)} \right]$$

$$= 1 - \left[\frac{1200}{990} \right]$$

$$= 1 - 1.2121$$

$$r_{xy} = -0.2121$$

$$r_{yz} = 1 - \left[\frac{6 \sum d_{yz}^2}{n(n^2-1)} \right]$$

$$= 1 - \left[\frac{6(214)}{10(10^2-1)} \right]$$

$$= 1 - \left[\frac{1284}{990} \right]$$

$$= 1 - 1.2969$$

$$r_{yz} = -0.2969$$

$$r_{zx} = 1 - \left[\frac{6 \sum d_{zx}^2}{n(n^2-1)} \right]$$

$$= 1 - \left[\frac{6(60)}{10(10^2-1)} \right]$$

$$= 1 - \frac{360}{990}$$

$$r_{zx} = 1 - 0.3636$$

$$r_{zx} = 0.6364$$



Conclusion :- $r_{21} = 0.6364$ is maximum.
Judge C and Judge A are common
approach to liking music.

Repeated Ranks (or) Tie Ranks :-

If any two or more individuals are equal in any classification with respect to characteristics A and B. If there is more than one item with the same value then the Spearman's Rank Correlation Coefficient breaks down.

In this case common ranks are given to the repeated items. This common rank is average of ranks which these items would have assumed.

In this case Spearman's Rank Correlation Coefficient formulae is

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

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

where m = number of times an item is repeated.



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

$$= \frac{2(2^2-1)}{12} + \frac{3(3^2-1)}{12} + \frac{2(2^2-1)}{12}$$

$$= \frac{6}{12} + \frac{24}{12} + \frac{6}{12}$$

$$= 0.5 + 2 + 0.5$$

$$CF = 3$$

$$\sum d_i^2 + CF = 47 + 3$$

$$\sum d_i^2 + CF = 50$$

∴ Spearman's Rank Cor Coe $\rho = 1 - \frac{6[\sum d_i^2 + CF]}{n(n^2-1)}$

$$= 1 - \frac{6(50)}{10(10^2-1)}$$

$$= 1 - \frac{300}{990}$$

$$= 1 - 0.3030$$

$$\rho = 0.697$$

Problem

① Calculate Rank Correlation Coefficient from the following data

marks in statistics:	40	46	54	46	60	70	75	60	83
marks in physics:	45	45	50	43	50	75	60	75	80

② Calculate Rank Correlation Coefficient from the following data.

x:	70	65	71	62	58	69	78	64
y:	91	76	65	83	90	65	55	48

→ obtain Rank Cor Coe to the following data.

x:	115	109	112	87	98	98	120	100	98	118
y:	75	73	85	70	76	65	82	73	68	80



REGRESSION

→ Define Regression

Regression means stepping back towards the average. It is a mathematical measurement in which we have to take average relationship between two (or) more variables. It was developed by "Sir Francis Galton". There are always two lines of regression.

1. Y on X regression

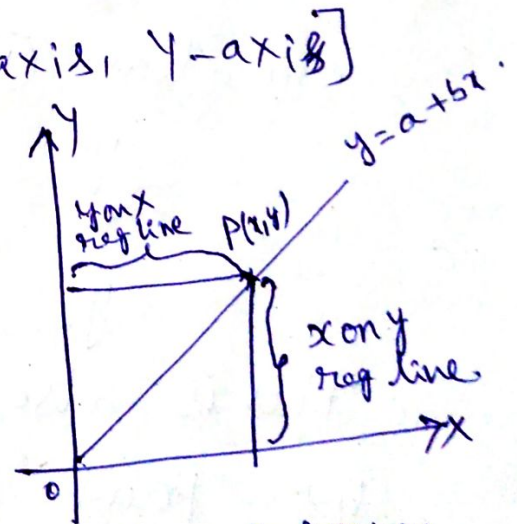
$$y - \bar{y} = r \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

2. X on Y regression

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

Let us take two axis [X-axis, Y-axis]

In XY plane we plot straight line $y = a + bx$ and $P(x, y)$ is any point on the straight line.



The parallel distance between point $P(x, y)$ and X-axis is minimum we get X on Y regression line.

The parallel distance between point $P(x, y)$ and Y-axis is minimum we get Y on X regression line.



Theorem

Derive the two regression lines

Proof:- Let (x_1, y_1) (x_2, y_2) (x_3, y_3) (x_n, y_n)

Let the set of n pairs of observations. To fit Y on X regression line we take straight line form $y = a + bx$. Here y is dependent variable and x is Independent variable. a, b are unknown constants. To estimate these unknown constants we use principles of least squares. Normal equations.

$$\sum y_i = na + b \sum x_i \quad \text{--- (1)}$$

$$\sum x_i y_i = a \sum x_i + b \sum x_i^2 \quad \text{--- (2)}$$

equation (1) is divide by n on both side.

$$\frac{\sum y_i}{n} = \frac{na}{n} + b \frac{\sum x_i}{n}$$

$$\bar{y} = a + b\bar{x} \quad \text{--- (3)}$$

This is also straight line passing through the point (\bar{x}, \bar{y})

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

$$\sigma^2 + (\bar{x})^2 = \frac{\sum x_i^2}{n} \quad \text{--- (4)}$$

equation (2) is divide by n and substitute (4)th equation.

$$\frac{\sum x_i y_i}{n} = a \frac{\sum x_i}{n} + b \frac{\sum x_i^2}{n}$$



$$\frac{\sum x_i y_i}{n} = a\bar{x} + b \frac{\sum x_i^2}{n}$$

$$\checkmark \frac{\sum x_i y_i}{n} = a\bar{x} + b[\sigma_x^2 + (\bar{x})^2] \quad \text{--- (5)}$$

and also $\text{Cov}(xy) = \frac{\sum x_i y_i}{n} - (\bar{x})(\bar{y})$

let $\text{Cov}(xy) = \mu_{11}$

$$\mu_{11} = \frac{\sum x_i y_i}{n} - (\bar{x})(\bar{y})$$

$$\mu_{11} + (\bar{x})(\bar{y}) = \frac{\sum x_i y_i}{n}$$

Substitute this value in to equation 5

$$\mu_{11} + (\bar{x})(\bar{y}) = a\bar{x} + b\sigma_x^2 + b(\bar{x})^2 \quad \text{--- (6)}$$

equation (3) is multiply by \bar{x} and subtract from equation (6)

equation (3) $\times \bar{x}$

$$\bar{x}\bar{y} = a\bar{x} + b(\bar{x})^2 \quad \text{--- (7)}$$

(6) - (7)

$$\begin{array}{r} \mu_{11} + (\bar{x})(\bar{y}) = a\bar{x} + b\sigma_x^2 + b(\bar{x})^2 \\ (\bar{x})(\bar{y}) = a\bar{x} + b(\bar{x})^2 \\ \hline \mu_{11} = b\sigma_x^2 \end{array}$$

$$b = \frac{\mu_{11}}{\sigma_x^2}$$

$$b = \frac{\text{Cov}(xy)}{\sigma_x^2}$$

$$b = \frac{r \cdot \sigma_x \cdot \sigma_y}{\sigma_x^2}$$

$$r = \frac{\text{Cov}(xy)}{\sigma_x \cdot \sigma_y}$$

$$\boxed{r \cdot \sigma_x \cdot \sigma_y = \text{Cov}(xy)}$$

$$\boxed{b = r \cdot \frac{\sigma_y}{\sigma_x}}$$



(\bar{x}, \bar{y}) is the Point $b = r \cdot \frac{\sigma_y}{\sigma_x}$ is the slope then you x regression line is

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

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

Ans:-

Q → Why two lines of Regression
 Ans:- There are always two lines of Regression.

① you x regression line

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

② x on y regression line.

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

To Estimate (or) predict the value of y to the given value of x we use you x regression line. Since in you x regression line we take straight line form $y = a + bx$. Here x is independent variable and y is dependent variable.

To Estimate (or) predict the value of x to the given value of y we use x on y regression line. Since in x on y regression line we take straight line form $x = a + by$. Here y is independent variable.



x is dependent variable.

Two regression lines are not Interchangeable or reversible. Because of simple reason that the parallel distance between the point (x, y) and x -axis is minimum we get Y on X regression line. The parallel distance between the point $P(x, y)$ and y -axis is minimum we get X on Y regression line.

In Case of Perfect Correlation $r = \pm 1$ two regression lines are coincide

Y on X regression line

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

$$\underline{r = \pm 1} \quad y - \bar{y} = \pm 1 \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \quad \text{--- (1)}$$

X on Y regression line

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

$$\underline{r = \pm 1} \quad x - \bar{x} = \pm 1 \cdot \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$(x - \bar{x}) \left(\pm 1 \cdot \frac{\sigma_y}{\sigma_x} \right) = y - \bar{y}$$

$$\therefore y - \bar{y} = \pm 1 \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \quad \text{--- (2)}$$

Equation (1) and (2) are equal.

except Perfect Correlation ($r = \pm 1$) two regression lines are coincide. In any other case two regression lines are entirely different.



Regression Coefficients :-

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

In this $r \cdot \frac{\sigma_y}{\sigma_x}$ is called regression coefficient of y on x regression line. It is denoted by b_{yx} .

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

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

In this $r \cdot \frac{\sigma_x}{\sigma_y}$ is called regression coefficient of x on y regression line. It is denoted with b_{xy}

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

Note! - (1) In Correlation $r_{xy} = r_{yx}$.

(2) In Regression $b_{xy} \neq b_{yx}$.



Q → Properties of two regression lines:-

① Two regression lines are passing through the point (\bar{x}, \bar{y})

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

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

② If (\bar{x}, \bar{y}) be the origin $(0, 0)$ the two regression lines are -

① y on x regression line $y = r \cdot \frac{\sigma_y}{\sigma_x} (x)$

② x on y regression line $x = r \cdot \frac{\sigma_x}{\sigma_y} (y)$

③ The Geometric mean of two regression coefficient gives Correlation coefficient.

$$i.e. \sqrt{b_{yx} \cdot b_{xy}}$$

$$= \sqrt{r \cdot \frac{\sigma_y}{\sigma_x} \cdot r \cdot \frac{\sigma_x}{\sigma_y}}$$

$$= \sqrt{r^2} = r$$

④ If one regression coefficient is increases then other regression coefficient decrease.

⑤ If $r=0$ then two regression lines are perpendicular.

⑥ There is no sign effect on two regression lines. But there is scale effect.



(7) two regression lines are not coincide except in case of perfect correlation i.e. $r = \pm 1$

(8) In correlation $r_{xy} = r_{yx}$. But in regression $b_{yx} \neq b_{xy}$.

Theorem Angle between two regression lines.

Proof:- y on x regression line $y - \bar{y} = r \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x})$ — (1)

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

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

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

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

in equation (1) $m_1 = r \cdot \frac{\sigma_y}{\sigma_x}$

$$m_2 = \frac{\sigma_y}{r \cdot \sigma_x}$$

Angle between two regression lines

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



Write 25 lines on regression

$$\begin{aligned} \tan \theta &= \frac{n \cdot \frac{\sum y}{\sigma} - \frac{\sum y}{n \cdot \sigma x}}{1 + \cancel{n} \cdot \frac{\sum y}{\sigma} \cdot \frac{1}{\cancel{n}} \cdot \frac{\sum y}{\sigma x}} \\ &= \frac{\frac{n^2 \sum y - \sum y}{n \sigma x}}{\frac{n \cdot \sigma_x^2 + n \cdot \sigma_y^2}{n \cdot \sigma_x^2}} = \frac{\frac{\sum y [n^2 - 1]}{n \cdot \sigma x}}{\frac{n [\sigma_x^2 + \sigma_y^2]}{n \cdot \sigma_x^2}} \\ &= \frac{\sum y [n^2 - 1]}{n [\sigma_x^2 + \sigma_y^2]} \times \frac{\cancel{n} \cdot \sigma_x^2}{\cancel{n} \cdot \sigma x} \end{aligned}$$

$$\tan \theta = \frac{\sigma x \cdot \sum y [n^2 - 1]}{n [\sigma_x^2 + \sigma_y^2]}$$

$$\theta = \tan^{-1} \left[\frac{\sigma x \cdot \sum y [n^2 - 1]}{n [\sigma_x^2 + \sigma_y^2]} \right]$$

Q → Difference between Correlation and Regression.

Correlation.
 ① If change in one variable effects to change in another variable. Then the two variables are said to be correlated. Then that relationship is called Correlation.

Regression
 ① Regression means stepping back towards the average in which we have to take average relationship between two (or) more variables.



② There is no cause and effect relationship between the variables.

③ There is no origin and scale effect on correlation.

④ In correlation $r_{xy} = r_{yx}$.

⑤ It is limited only to linear relationship between the variables.

② There is cause and effect relationship between the variables.

③ There is no origin effect but there is scale effect on regression.

④ In regression $b_{yx} \neq b_{xy}$.

⑤ It has wide scope to study linear as well as non linear relationship between the variables.



Problems

① Fit two regression equations from the following data.

x:	1	5	3	2	1	1	7	3
y:	6	1	0	0	1	2	1	5

Sol: - Two regression lines are.

① Y on X regression line $y - \bar{y} = r \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x})$

② X on Y regression line $x - \bar{x} = r \cdot \frac{\sigma_x}{\sigma_y} (y - \bar{y})$

where $r = \frac{\text{Cov}(xy)}{\sigma_x \cdot \sigma_y}$

$\text{Cov}(xy) = \frac{\sum x_i y_i}{n} - (\bar{x})(\bar{y})$

$\bar{x} = \frac{\sum x_i}{n}$

$\bar{y} = \frac{\sum y_i}{n}$

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

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

x	y	$\sum x_i y_i$	$\sum x_i^2$	$\sum y_i^2$
1	6	6	1	36
5	1	5	25	1
3	0	0	9	0
2	0	0	4	0
1	1	1	1	1
1	2	2	1	4
7	1	7	49	1
3	5	15	9	25
<u>23</u>	<u>16</u>	<u>36</u>	<u>99</u>	<u>68</u>



$$\bar{x} = \frac{\sum x_i}{n} = \frac{23}{8}$$

$$\bar{x} = 2.875$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{16}{8}$$

$$\bar{y} = 2$$

$$\begin{aligned} \text{Cov}(xy) &= \frac{\sum x_i y_i}{n} - (\bar{x})(\bar{y}) \\ &= \frac{36}{8} - (2.875)(2) \\ &= 4.5 - 5.75 \end{aligned}$$

$$\text{Cov}(xy) = -1.25$$

$$\begin{aligned} \sigma_x &= \sqrt{\frac{\sum x_i^2}{n} - (\bar{x})^2} \\ &= \sqrt{\frac{99}{8} - (2.875)^2} \\ &= \sqrt{12.375 - 8.2656} \end{aligned}$$

$$= \sqrt{4.1094}$$

$$\sigma_x = 2.0271$$

$$\begin{aligned} \sigma_y &= \sqrt{\frac{\sum y_i^2}{n} - (\bar{y})^2} \\ &= \sqrt{\frac{68}{8} - (2)^2} \\ &= \sqrt{8.5 - 4} \\ &= \sqrt{4.5} \end{aligned}$$

$$\sigma_y = 2.1213$$

$$r = \frac{\text{Cov}(xy)}{\sigma_x \cdot \sigma_y} = \frac{-1.25}{(2.0271)(2.1213)} = \frac{-1.25}{4.300}$$

$$r = -0.2906$$

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

$$y - 2 = (-0.2906) \frac{2.1213}{2.0271} (x - 2.875)$$

$$y - 2 = \frac{-0.6164}{2.0271} (x - 2.875)$$

$$= (-0.3041) (x - 2.875)$$

$$y - 2 = -0.3041x + 0.8742 \Rightarrow y = -0.3041x + 0.8742$$

$$y = -0.3041x + 0.8742$$



x on y regression line

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

$$x - 2.875 = \left(\frac{-0.2906}{2.1213} \right) (y - 2)$$

$$x - 2.875 = \frac{-0.6164}{2.1213} (y - 2)$$

$$x - 2.875 = (-0.2906) (y - 2)$$

$$x - 2.875 = -0.2906y + 0.5812$$

$$x = -0.2906y + 0.5812 + 2.875$$

$$x = -0.2906y + 3.4562$$

Problem ① Fit two regression equations

from the data given below.

x :	30	32	35	40	45
y :	20	28	30	32	35

② Fit two regression equations from the following data

x :	2	4	5	6	8	11
y :	18	12	10	8	7	5

Problem ① The regression equations are

$$8x - 10y + 66 = 0 \quad \text{and} \quad 40x - 18y - 214 = 0$$

Variance of the variable x is 9.

Then find Means of x and y and also find σ_y and σ_{xy} .

Sol:- Given two regression equations are

$$8x - 10y + 66 = 0 \quad \text{--- (1)}$$

$$40x - 18y - 214 = 0 \quad \text{--- (2)}$$

$$\textcircled{1} \times 5 - \textcircled{2}$$

$$\begin{array}{r} 40x - 50y + 330 = 0 \\ - 40x + 18y - 214 = 0 \\ \hline -32y + 544 = 0 \end{array}$$

$$y = \frac{544}{32} = 17$$

$$\boxed{\bar{y} = 17}$$

Substitute $y = 17$ in equation (1)

$$8x - 10y + 66 = 0$$

$$8x - 10(17) + 66 = 0$$

$$8x - 170 + 66 = 0$$

$$8x = 104 \Rightarrow x = \frac{104}{8} = 13$$

\therefore means of x and y is $\boxed{\bar{x} = 13, \bar{y} = 17}$

$$\textcircled{1} \Rightarrow 8x - 10y + 66 = 0$$

$$8x + 66 = 10y$$

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

$$\boxed{\therefore b_{yx} = \frac{8}{10}}$$

$$\textcircled{2} \Rightarrow 40x - 18y - 214 = 0$$

$$40x = 18y + 214$$

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

$$\boxed{b_{xy} = \frac{18}{40}}$$

We know that $r = \sqrt{b_{yx} \cdot b_{xy}}$

$$r = \sqrt{\left(\frac{8}{10}\right)\left(\frac{18}{40}\right)} = \sqrt{\frac{18}{50}}$$



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

We are given $\sigma_x^2 = 9$
 $\sigma_x = 3$

$$\frac{8}{10} = \left(\frac{0.2}{0.6}\right) \cdot \frac{\sigma_y}{3}$$

$$\frac{8}{(10)(0.2)} = \sigma_y$$

$$\therefore \sigma_y = 4$$

\therefore Variance of y is $\sigma_y^2 = 16$

Problem (2) You are given below following information about advertising and Sales.

	advertising expenditure in (Lakhs Rs)	Sales (in Crores)
Mean	20	100
S.D	5	12

Correlation Coefficient is 0.8

- Calculate two regression lines.
- Find the likely sales when advertising expenditure is Rs 25 lakhs.
- What should be advertisement if the Company wants to attain sales target of Rs 130 crores.



Sol: We are given

$$\bar{x} = 20, \quad \bar{y} = 100, \quad r = 0.8$$

$$\sigma_x = 5, \quad \sigma_y = 12$$

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

$$y - 100 = (0.8) \left(\frac{12}{5} \right) (x - 20)$$

$$y = 1.92x - 38.4 + 100$$

$$\boxed{y = 1.92x + 61.6}$$

(ii) To find likely sales (Y) when adv expenditure $x = 25$.

To find Y value we use Y on X regression

$$y = 1.92x + 61.6$$

$$y = 1.92(25) + 61.6$$

$$y = 48 + 61.6$$

$$y = 109.6 \text{ Crores}$$

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

$$x - 20 = (0.8) \left(\frac{5}{12} \right) (y - 100)$$

$$x = 0.3333y - 33 + 20$$

$$\boxed{x = 0.3333y - 13.33}$$

(iii) What should be adv (x) when Company wants to attain sales target $y = 130$ Crores

To estimate X value when $y = 130$ we use X on Y regression line



$$x = 0.33y - 13$$

$$\text{If } y = 130 \Rightarrow x = (0.33)(130) - 13$$

$$x = 42.9 - 13$$

$$x = 29.9 \text{ (Correct)}$$

Problem ② The following detail of marks in Maths and Statistics of 2nd B.Tech students in our college.

	<u>Maths</u>	<u>Statistics</u>
Average marks	39.5	49.5
S.D	10.8	16.8

The Correlation Coefficient between Maths and Statistics is 0.8 by using this find regression lines and estimate marks in Maths when mark in Statistics is 52.

Sol: let us denote marks in Maths is x and mark in Statistics is y .

We are given $\bar{x} = 39.5$, $\bar{y} = 49.5$, $r = 0.8$
 $\sigma_x = 10.8$, $\sigma_y = 16.8$

Two regression lines y on x reg line
 $y - \bar{y} = r \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x})$

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