

Types of Data

Univariate Data

One variable
at a time

Bivariate Data

Two variables at
same time (Pairs)

* Bivariate Data - When data are collected on two variables simultaneously, they are known as Bivariate data. And the corresponding frequency distribution derived from it, is known as Bivariate frequency distribution.

- Ex: 1) Amount of Investment & profit
2) Advertisement Expenditure & sales
3) Demand & price

Demand (X)Price (Y)

80

20

90

18

100

16

90

17

90

18

110

15

X is changed due to change in Y.
So, X and Y are Correlated.

* Marginal Distribution

Distribution

It is the frequency of one variable (X or Y) across the other variable's full range of values.

No. of Marginal
Distribution in $\Rightarrow 2$
Bivariate Data

* Conditional Distribution

It is the Distribution of one variable (X or Y) across a particular sub-population of the other variable.

No. of Conditional $\Rightarrow m+n$
Distribution

$m =$ no. of class interval of X

$n =$ no. of class interval of Y

* Correlation - If change in one variable causes change in another variable either Directly or Inversely, then two variables are known to be Associated or Correlated.

Types of Correlation

positive
correlation



If two variables
move in the
same direction.

Ex:- price & supply



(upward sloping curve)

Negative
correlation



If two variables
move in the
opposite direction.

Ex:- price & demand



(Downward sloping curve)

* Measurement of Correlation

What to measure



Nature

Relative strength

Value

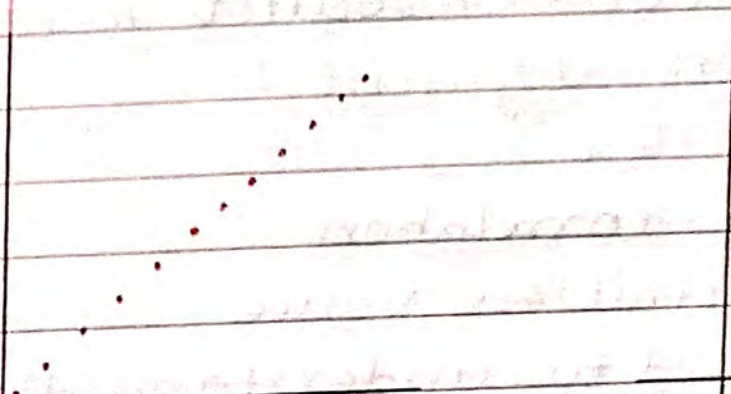
(High, low)

measures of correlation -

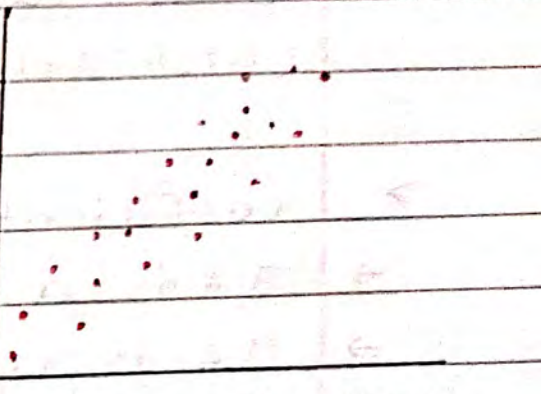
1. → Scatter Diagram
2. → Karl Pearson's product moment correlation coefficient
3. → Co-efficient of Concurrent deviation
4. → Spearman's Rank Correlation Coefficient

1. * Scatter Diagram

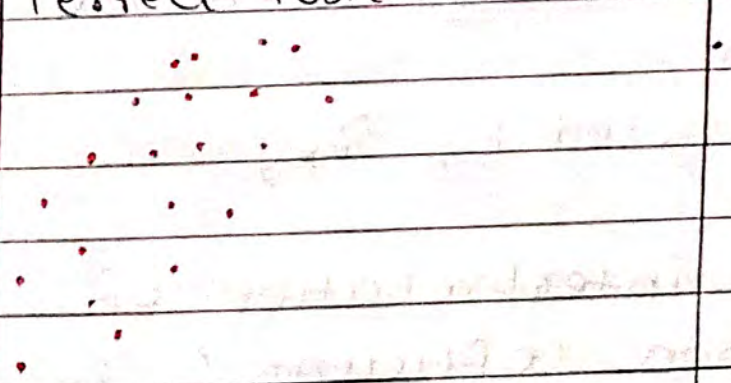
- It helps us to find Nature and Relation strength of Correlation.
- It is useful for non-linear correlation also
- It cannot be used to determine Value
- Diagrams are time taking



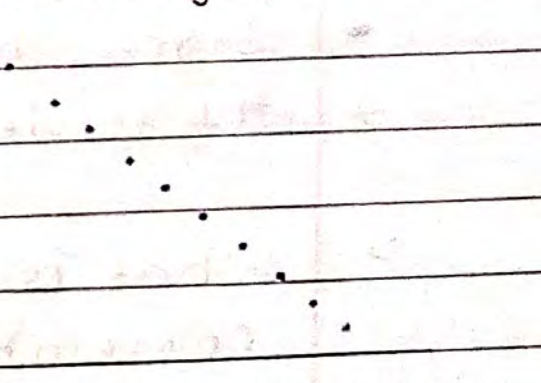
Perfect Positive



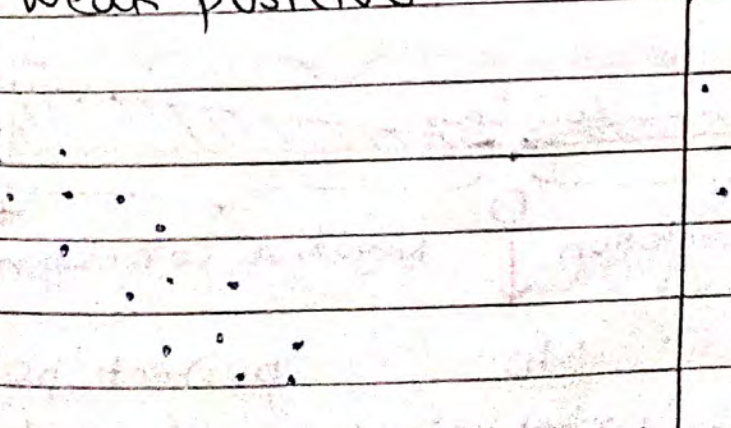
Strong positive



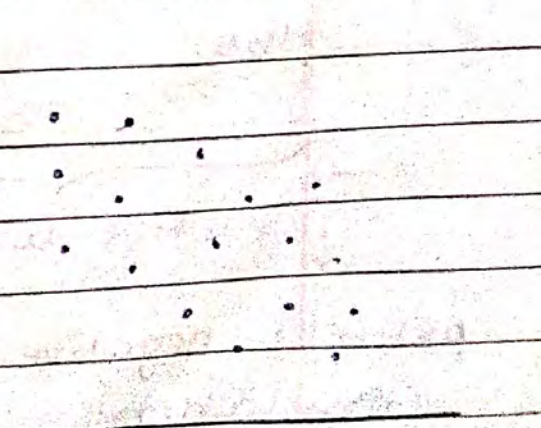
Weak positive



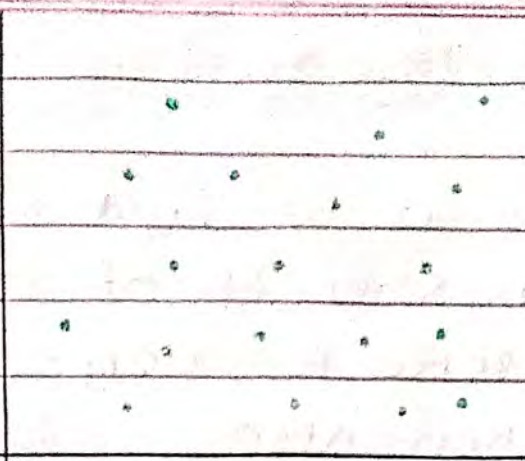
Perfect Negative



Strong Negative



Weak negative

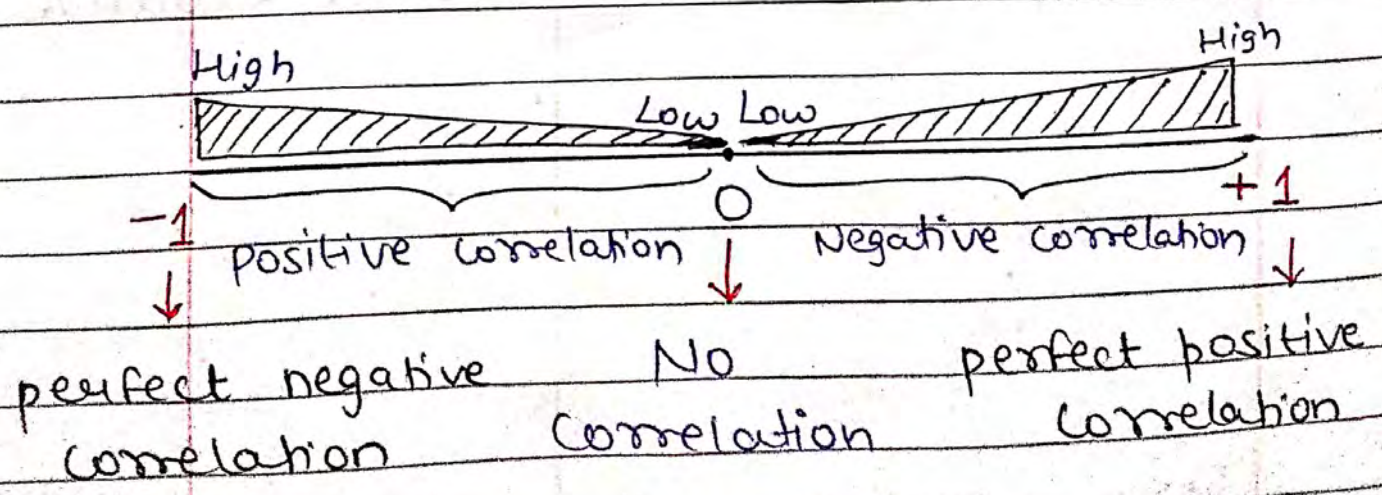


No Correlation

2. * Karl Pearson's Product Moment Correlation Coefficient :-

- > Coefficient correlation
- It is a unitless value
- It is used to understand the nature as well as accurate strength.
- It is denoted by r_{xy}

> What is interpretation of correlation coefficient? r_{xy}



> How to measure? r_{xy}

$$r_{xy} = \frac{\text{Covariance of X and Y}}{\text{SD of X} \times \text{SD of Y}}$$

Covariance

> Coefficient of x and y (bivariate)

$$\text{Variance of x} \Rightarrow \frac{\sum (x - \bar{x})^2}{n} \quad \text{or} \quad \frac{\sum x^2 - (\bar{x})^2}{n}$$

$$\Downarrow \quad \Downarrow$$

$$\frac{\sum (x - \bar{x})(x - \bar{x})}{n} \quad \text{or} \quad \frac{\sum x \cdot x - \bar{x} \cdot \bar{x}}{n}$$

$$\text{Covariance of x and y} \Rightarrow \frac{\sum (x - \bar{x})(y - \bar{y})}{n} \quad \text{or} \quad \frac{\sum x \cdot y - \bar{x} \cdot \bar{y}}{n}$$

Example - Compute Correlation Coefficient b/w x and y from the following data :-

$$n = 10 \quad \sum x^2 = 200$$

$$\sum x = 40 \quad \sum y^2 = 262$$

$$\sum y = 50 \quad \sum xy = 220$$

$$\text{Sol}^n = \quad r_{xy} = \frac{\text{COV}(x, y)}{\text{SD}_x \cdot \text{SD}_y} \quad \left| \quad \begin{aligned} \bar{x} &= \frac{\sum x}{n} = \frac{40}{10} = 4 \\ \bar{y} &= \frac{\sum y}{n} = \frac{50}{10} = 5 \end{aligned} \right.$$

$$\text{COV}(x, y) = \frac{220 - 4 \times 5}{10}$$

$$= \frac{200 - 20}{10}$$

$$= \underline{\underline{2}}$$

$$\begin{aligned} \text{SD of } x &= \sqrt{\frac{208}{10} - \left(\frac{48}{10}\right)^2} \\ &= \sqrt{20 - 16} \\ &= \sqrt{4} = \underline{\underline{2}} \end{aligned}$$

$$\begin{aligned} \text{SD of } y &= \sqrt{\frac{262}{10} - \left(\frac{50}{10}\right)^2} \\ &= \sqrt{26.2 - 25} \\ &= \sqrt{1.2} = \underline{\underline{1.0954}} \end{aligned}$$

$$\begin{aligned} r_{xy} &= \frac{\text{Cov}(x, y)}{\text{SD of } x \cdot \text{SD of } y} \\ &= \frac{2}{2 \times 1.0954} \\ &= \underline{\underline{0.9128}} \end{aligned}$$

* Properties of Correlation Coefficient

(1.)

	Change of origin	Change of scale (value)	Change of scale (sign.)
--	------------------	-------------------------	-------------------------

Central tendency	✓	✓	✓
------------------	---	---	---

Measure of Dispersion	✗	✓	✗
-----------------------	---	---	---

Correlation Coefficient	✗	✗	✓
-------------------------	---	---	---

Original x $\xrightarrow{\text{Origin/Scale}}$ modified u

y $\xrightarrow{\text{Origin/Scale}}$ v

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

* If sign of both change of scale are same, $r_{xy} = r_{uv}$

* If sign of both change of scale are different $r_{uv} = -r_{xy}$

Ex:- Given:- $r_{xy} = 0.65$

$$u = 5 - 2x \quad v = 10 + 3y$$

Change of scale of $x = -ive$

Change of scale of $y = +ive$

So,

$$r_{uv} = -r_{xy} = -0.65$$

(2.) The coefficient of correlation is a unit-free measure.

(3.) Values lies from -1 to $+1$

3. * SPEARMAN'S RANK CORRELATION COEFFICIENT

To find Correlation coefficient in Case of Qualitative data: as Ranks

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

Ex: -

No. of Candidates = 1 2 3 4 5 6 7 8 9 10

Rank by Judge A = 10 5 6 1 2 3 4 7 9 8

Rank by Judge B = 5 6 9 2 8 7 3 4 10 1

Difference = 5 1 3 1 6 4 1 3 1 7

Square of difference (d^2) = 25, 1, 9, 1, 36, 16, 1, 9, 1, 49

$$\sum d^2 = 148$$

$$r_r = 1 - \frac{6 \times 148}{10(10^2-1)} \quad [n=10]$$

$$r_r = 1 - 0.897$$

$$\underline{\underline{r_r = 0.1030}}$$

> In case of Tie in values
 - Average Rank to be given in values

$$r_r = 1 - \frac{6(\sum d^2 + \text{adj. value})}{n(n^2 - 1)}$$

$$\text{adj value} = \sum \frac{t^3 - t}{12} \quad [t = \text{tie length}]$$

Ex:-

Eco marks X	80	56	50	48	50	62	60
Stats marks Y	90	75	75	65	65	50	65
Rank R_x	1	4	5.5	7	5.5	2	3
R_y	1	2.5	2.5	5	5	7	5
d^2	0	2.25	9	4	0.25	25	4

$$\sum d^2 = 44.5$$

$$\text{Adj value} \Rightarrow \sum \frac{t^3 - t}{12}$$

	$\frac{t}{}$	$\frac{t^3 - t}{12}$
50-50	2	$\frac{2^3 - 2}{12}$
75-75	2	$\frac{2^3 - 2}{12}$
65-65-65	3	$\frac{3^3 - 3}{12}$

$$\sum \frac{t^3 - t}{12} = \underline{\underline{3}}$$

$$r_r = 1 - \frac{6[44.5 + 3]}{7(48)}$$

$$r_r = 1 - 0.849$$

$$\underline{\underline{r_r = 0.151}}$$

* 4. Co-efficient of Concurrent deviations

> Very quick, simple and casual method of finding correlation when we are not serious about the magnitude of two variables.

$$r_c = \pm \sqrt{\frac{\pm 2C - m}{m}}$$

m = no. of pairs compared ($n-1$)
 C = no. of concurrent deviations
 (no. of plus signs)

EX:-	1990	1991	1992	1993	1994	1995	1996	1997
Price	25	28	30	23	35	38	39	42
Demand	35	34	35	30	29	28	26	43
dx	na	+	+	-	+	+	+	+
dy	na	-	+	-	-	-	-	+
dx dy		-	(+)	(+)	-	-	-	(+)

$$m = n - 1 = 8 - 1 = 7$$

$$C = \text{no. of plus sign in } dx dy = 3$$

$$r_c = \pm \sqrt{\frac{\pm 2(3) - 7}{7}}$$

$$r_c = \pm \sqrt{\frac{\pm (-1)}{7}}$$

$$a_c = - \left[- \left(\frac{-1}{7} \right) \right]$$

$$a_c = - \cancel{0.317} - \underline{\underline{0.3779}}$$

* एक जगह - तो दूसरी जगह भी -
 और अगर एक जगह + तो दूसरी
 जगह भी +

$$\frac{2c-m}{m} = +ive = + \text{ both place}$$

$$\frac{2c-m}{m} = -ive = - \text{ both place}$$

* Regression

Estimation of dependent variable with the given independent variable.

- Estimation of Y when X is given
 Y: Dependent, X: Independent

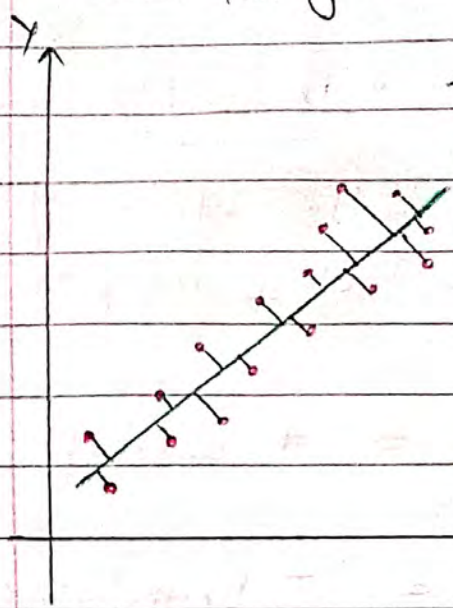
- Estimation of X, when Y is given
 X: Dependent, Y: Independent

1) Perfect Correlation

- When linear eqⁿ exist between two variables, correlation is perfect.

- perfect correlation is represented by linear equation and this eqⁿ can be used for regression purpose directly.

2) Not perfect Correlation

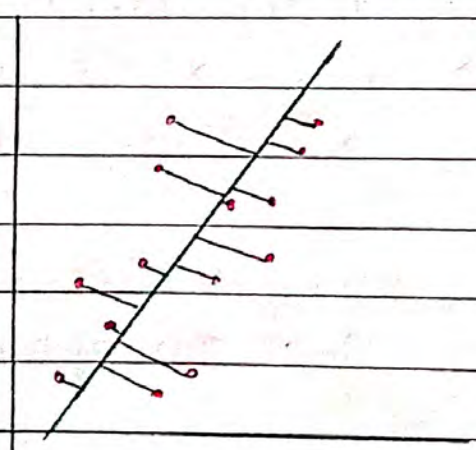


- In case of non perfect Correlation point are not in line so we don't have available line eqⁿ.

then, in that case we will
 → make a line for purpose
 X of estimation.

↙ error/residue

↗ These lines are called as regression lines.



Reg. line of x on y
 used for estimation
 of x when y is given



Reg line of y on x
 used for est. of y
 when x is given

“method of Least Square”

Objectives

Estimation of Y
 When X is given
 ↓

Estimation of X
 When Y is given
 ↓

Equation format:
 $Y - \bar{Y} = b_{yx} (X - \bar{X})$

Equation format:
 $X - \bar{X} = b_{xy} (Y - \bar{Y})$

$$Y = a + bX$$

$$X = a + bY$$

b_{yx} = reg. coeff.
 of Y on X

b_{xy} = reg. coeff. of
 X on Y

Regression = $\frac{\text{Covariance b/w two variables}}{\text{Variance of Independent Variable}}$
 Coefficient

$$b_{yx} = \frac{\text{COV}(X, Y)}{\text{Variance of } X}$$

$$b_{xy} = \frac{\text{COV}(X, Y)}{\text{Variance of } Y}$$

Ex Find two regression line/eqⁿ from:

X:	2	4	5	5	8	10
Y:	6	7	9	10	12	12

Hence estimate Y when X is 13 and
 estimate X when Y is 15.

Solⁿ → Est. of Y on X

$$Y - \bar{Y} = b_{yx}(X - \bar{X})$$

$$Y - \frac{56}{6} = b_{yx}\left(X - \frac{34}{6}\right)$$

$$Y - 9.3333 = 0.8145(X - 5.6666)$$

$$Y = 0.8145X - 4.6155 + 9.3333$$

$$Y = 0.8145X + 4.7178$$

Regression Eqⁿ of Y on X.

If X = 13

$$Y = 4.7178 + 13 \times (0.8145)$$

$$Y = 15.3063$$

$$b_{yx} = \frac{\text{Cov}(X, Y)}{\text{Var. X}}$$

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

$$b_{yx} = \frac{351}{6} - \left(\frac{34}{6} \times \frac{56}{6}\right)$$

$$\frac{234}{6} - \left(\frac{34}{6}\right)^2$$

$$b_{yx} = \frac{5.6111}{6.888} = 0.8145$$

Estimation of X on Y

$$X - \bar{X} = b_{xy}(Y - \bar{Y})$$

$$X - \frac{34}{6} = 1.0745\left(Y - \frac{56}{6}\right)$$

$$X - 5.6666 = 1.0745(Y - 9.3333)$$

$$X = 1.0745Y - 4.362$$

Reg. eqⁿ of X on Y

If Y = 15

$$X = 1.0745(15) - 4.362$$

$$X = 11.7555$$

$$b_{xy} = \frac{\text{Cov}(X, Y)}{\text{Var. of Y}}$$

$$= \frac{5.6111}{\frac{\sum y^2 - (\sum y)^2}{n}}$$

$$= \frac{5.6111}{\left(\frac{554}{6}\right) - \left(\frac{56}{6}\right)^2}$$

$$= 1.0745$$

* properties of Regression lines / Coefficient

> Change of Origin / scale
 Origin: GB (no impact)

Scale: a) value \rightarrow Yes
 b) sign \rightarrow Yes

$x \rightarrow u$
 $y \rightarrow v$

$$b_{vu} = b_{xy} \times \frac{\text{Change of scale of } x}{\text{Change of scale of } y}$$

$$b_{vu} = b_{yx} \times \frac{\text{Change of scale of } Y}{\text{Change of scale of } X}$$

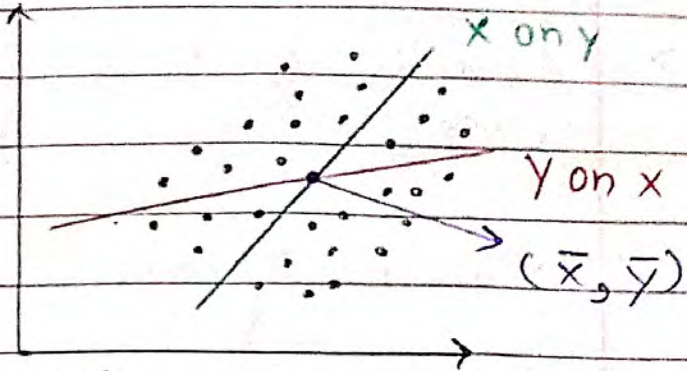
Ex. $u + 3x = 10$ $2y + 5v = 25$
 $u = 10 - 3x$ $5v = 25 - 2y$
 Change of Scale of $x = -3$ $v = 5 - \frac{2y}{5}$

$b_{yx} = 0.80$ Change of Scale of $y = -\frac{2}{5}$

$b_{vu} = ?$
 $b_{vu} = 0.80 \times \frac{-2/5}{-3} = \frac{0.80 \times 2}{5 \times 3}$

$b_{vu} = \frac{8}{75} = \underline{\underline{0.1066}}$

> Two lines of regression intersect at (\bar{x}, \bar{y}) . If lines are not identical.



Ex- $7x - 3y - 18 = 0$] x1
 $4x - y - 11 = 0$] x3

$$\begin{array}{r} 7x - 3y - 18 = 0 \\ 12x - 3y - 33 = 0 \\ - \quad + \quad + \end{array}$$

$$-5x + 15 = 0$$

$$5x = 15$$

$$\boxed{x = 3}$$

$$7x(3) - 3y - 18 = 0$$

$$21 - 3y - 18 = 0$$

$$3 = 3y$$

$$\boxed{y = 1}$$

$$\bar{x} = 3, \bar{y} = 1$$

> Relation b/w Correlation and Regression Coefficient \rightarrow

$$\boxed{r_{xy} = \pm \sqrt{b_{xy} \times b_{yx}}}$$

Ex- $b_{xy} = -1.06$ $b_{yx} = -0.73$

$$r_{xy} = \sqrt{-1.06 \times -0.73} = 0.879 \quad \times$$

$$r_{xy} = \sqrt{-(1.06) \times (-0.73)} = -0.879$$

$$r_{xy} = \frac{\text{COV}(x,y)}{SD_x SD_y} \Rightarrow \text{COV}(x,y) = r_{xy} \cdot SD_x \cdot SD_y$$

$$b_{xy} = \frac{\text{COV}(x,y)}{\text{Var. of } y}$$

$$b_{yx} = \frac{\text{COV}(x,y)}{\text{Var. of } x}$$

$$b_{xy} = \frac{r_{xy} \cdot SD_x \cdot SD_y}{(SD_y)^2}$$

$$b_{yx} = \frac{r_{xy} \cdot SD_x \cdot SD_y}{(SD_x)^2}$$

$$b_{xy} = r_{xy} \left(\frac{SD_x}{SD_y} \right)$$

$$b_{yx} = r_{xy} \left(\frac{SD_y}{SD_x} \right)$$

* you are given two regression lines/eqⁿ and you have to find out which one is Y on X and X on Y.

ex- $7x - 3y - 18 = 0$ & $4x - y - 11 = 0$

Coefficient of x
Coefficient of y
(Ignore negative sign)

Coeff. of x
Coeff. of y
(Ignore negative sign)

$$\frac{7}{3} = 2.333 <$$

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

So, ↓
y on x

↓
x on y

जो बड़ा आरगा वो x on y and दूसरा वला y on x.

$$* \text{probable error} = \frac{0.6745 \times \sqrt{1 - \mu^2}}{\sqrt{n}}$$

$$\text{EX} \rightarrow \mu = 0.7$$

$$n = 64$$

$$\text{probable error} = \frac{0.6745 \times \sqrt{1 - (0.7)^2}}{\sqrt{64}}$$

$$= \frac{0.6745 \times 0.51}{8}$$

$$= \boxed{0.04299}$$

* limits of population correlation coefficient :-

lower limit

$\mu + \text{probable error}$

Upper limit

$\mu - \text{probable error}$

$$\text{EX} - \mu = 0.7$$

$$PE = 0.0423$$

limit of population correlation coefficient :-

$$\Rightarrow 0.7 + 0.0423$$

$$\Rightarrow 0.7423$$

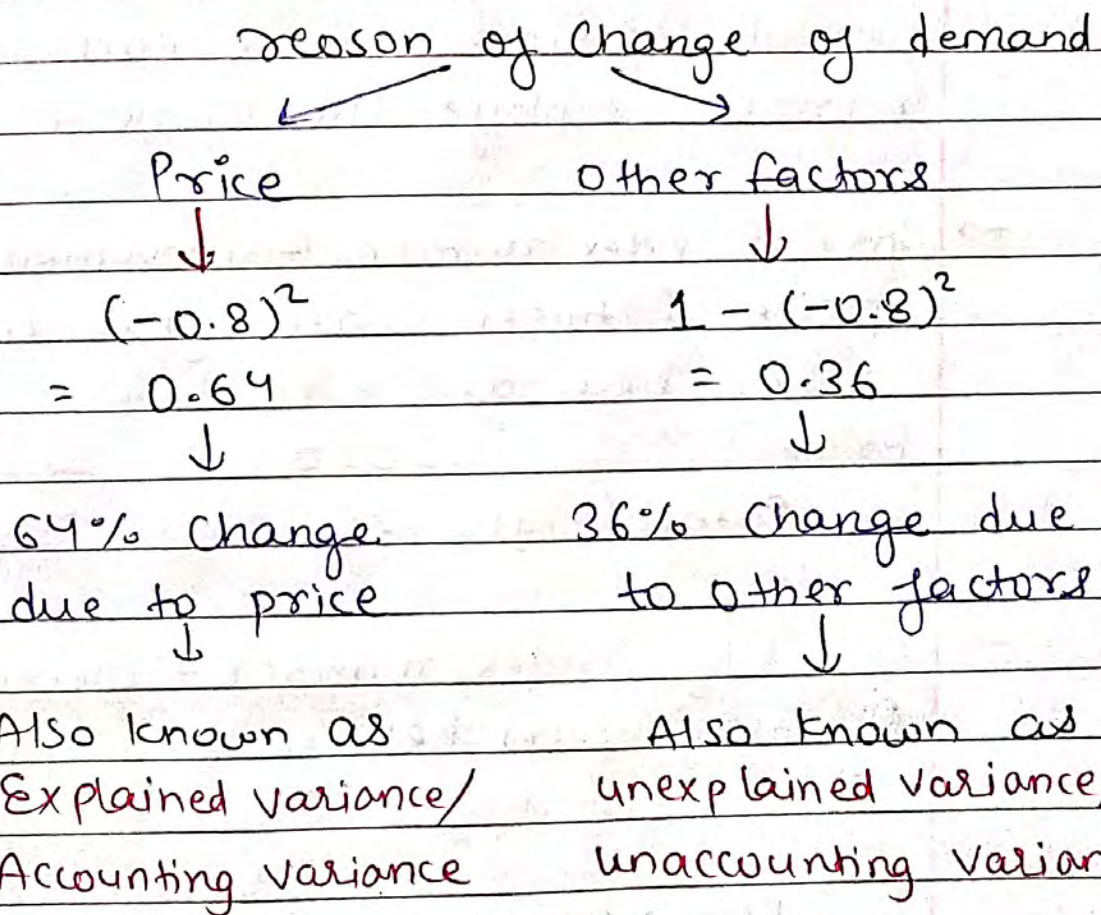
$$\Rightarrow 0.7 - 0.0423$$

$$\Rightarrow 0.6577$$

* $\underbrace{\text{Coff. of determination}}_{r^2} / \underbrace{\text{Coff. of non-determination}}_{1-r^2}$

Ex- Corr. Coff. b/w demand (X) and price (Y) is -0.8.

⇒ Demand is change due to price but there are other reasons also.



* Where no casual (reason) relationship So, it is called "Non-sense correlation" or "Spurious correlation."

* Karl Pearson correlation Coff. is best.