

Correlation

If the change in one variable affects a change in the other variable, the variables are said to be correlated.

Positive / Direct Correlation

If the increase/decrease in one results in the corresponding increase/decrease in the other, then the variables are positively correlated.

Ex: Correlation between (i) heights and weights of a group of people (ii) Income and expenditure.

Negative / Diverse Correlation

If the increase/decrease in one results in corresponding decrease/increase in other, then the variables are negatively correlated.

Ex: Correlation between demand and price.

Karl Pearson's Coefficient of Correlation

Correlation Coefficient between two random variables X and Y , denoted by $r(X, Y)$ or r_{XY} is a numerical measure of linear relationship between them and is defined as

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

$$\text{where } \text{Cov}(X, Y) = \sigma_{XY} = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{or } \sigma_{XY} = \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y},$$

$$\sigma_x^2 = \frac{1}{n} \sum (x_i - \bar{x})^2 \text{ or } \sigma_x^2 = \frac{1}{n} \sum x_i^2 - \bar{x}^2,$$

$$\sigma_y^2 = \frac{1}{n} \sum (y_i - \bar{y})^2 \text{ or } \sigma_y^2 = \frac{1}{n} \sum y_i^2 - \bar{y}^2.$$

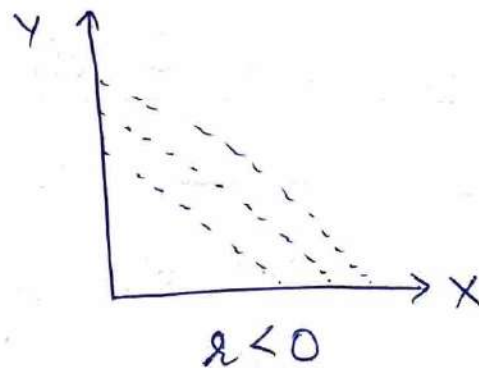
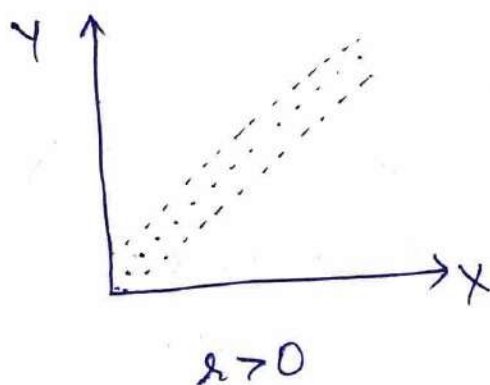
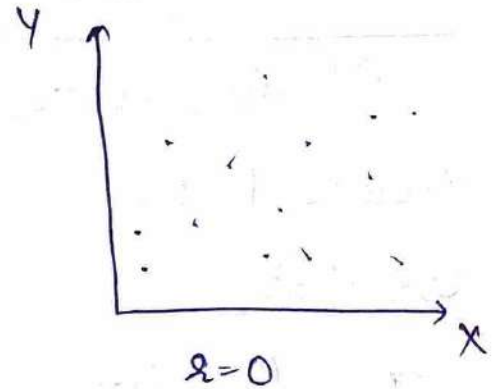
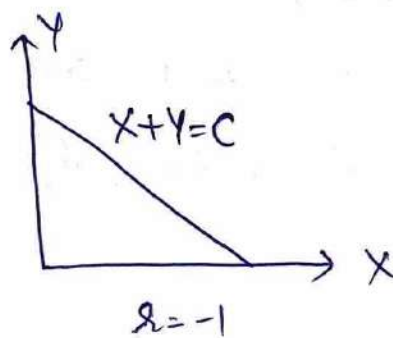
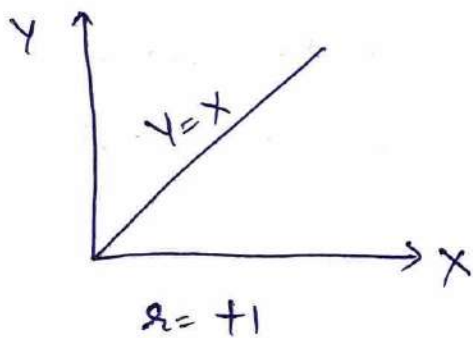
Range of Correlation Coefficient

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

If $r = -1$, then correlation is perfect and negative.

If $r = 1$, then correlation is perfect and positive.

Scatter Plots / diagrams



Results 1. Correlation Coefficient is independent of change of origin and scale.

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

$$r(ax+b, cy+d) = \frac{ac}{|ac|} r(x, y).$$

3. Two independent variables are uncorrelated.

Ex Calculate the correlation coefficient for the following heights (in inches) of fathers (X) and their sons (Y):

X:	65	66	67	67	68	69	70	72
Y:	67	68	65	68	72	72	69	71

Sol ∴ Let $U = X - 68$, $V = Y - 69$

<u>X</u>	<u>Y</u>	<u>U</u>	<u>V</u>	<u>U²</u>	<u>V²</u>	<u>UV</u>
65	67	-2 -3	-2	9	4	6
66	68	-2	-1	4	1	2
67	65	-1	-3 -4	1	16	4
67	68	-1	-1	1	1	1
68	72	0	3	0	9	0
69	72	1	3	1	9	3
70	69	2	0	4	0	0
72	71	4	2	16	4	8
		<u>0</u>	<u>0</u>	<u>36</u>	<u>44</u>	<u>24</u>

$$\bar{U} = 0, \bar{V} = 0, \sigma_{UV} = \frac{1}{n} \sum u_i v_i - \bar{U} \bar{V} = \frac{1}{8} (24) = 3$$

$$\sigma_U^2 = \frac{1}{n} \sum u_i^2 - \bar{U}^2 = \frac{1}{8} (36) = 4.5$$

$$\sigma_V^2 = \frac{1}{n} \sum v_i^2 - \bar{V}^2 = \frac{1}{8} (44) = 5.5$$

$$\therefore r(U, V) = \frac{\sigma_{UV}}{\sigma_U \sigma_V} = \frac{3}{\sqrt{4.5 \times 5.5}} = \frac{3}{\sqrt{24.75}} = \frac{3}{4.97}$$

$$\Rightarrow r(U, V) = 0.604 = r(X, Y)$$

$$\boxed{r(X, Y) = 0.604}$$

Ex A computer while calculating correlation coefficient between two variables X and Y from 25 pairs of observations obtained the following results:

$$n = 25, \sum X = 125, \sum X^2 = 650, \sum Y = 100, \sum Y^2 = 460, \sum XY = 508.$$

It was, however, later discovered at the time of checking that he had copied down two pairs as

X	Y
6	14
8	6

while the correct values were

X	Y
8	12
6	8

Obtain the correct value of correlation coefficient.

Sol. Corrected $\sum X = 125 - 6 - 8 + 8 + 6 = 125$

Corrected $\sum Y = 100 - 14 - 6 + 12 + 8 = 100$

Corrected $\sum X^2 = 650 - 36 - 64 + 64 + 36 = 650$

Corrected $\sum Y^2 = 460 - 196 - 36 + 144 + 64 = 436$

Corrected $\sum XY = 508 - 84 - 48 + \cancel{84} + 48 = 520$

$$\bar{X} = \frac{\sum X}{n} = \frac{125}{25} = 5, \quad \bar{Y} = \frac{1}{n} \sum Y = \frac{100}{25} = 4$$

$$\sigma_{XY} = \frac{1}{n} \sum XY - \bar{X} \bar{Y} = \frac{1}{25} (520) - 20 = 0.8$$

$$\sigma_x^2 = \frac{1}{n} \sum x^2 - \bar{x}^2 = \frac{1}{25} (650) - 25 = 1$$

$$\sigma_y^2 = \frac{1}{n} \sum y^2 - \bar{y}^2 = \frac{1}{25} (436) - 16 = 1.44$$

$$\therefore \text{Corrected } r(x, y) = \frac{0.8}{\sqrt{1 \times 1.44}} = \frac{0.8}{1.2} = 0.67.$$

$$\boxed{r(x, y) = 0.67}$$

Correlation coefficient for bi-variate data and probability distribution

$g(x) = \sum_y f(x,y)$ is the sum of frequencies along any column and $h(y) = \sum_x f(x,y)$ is the sum of frequencies along any row.

$$\text{Here } \sum_x \sum_y f(x,y) = \sum_y \sum_x f(x,y) = \sum_x g(x) = \sum_y h(y) = N$$

$$\bar{x} = \frac{1}{N} \sum_x \sum_y x f(x,y) = \frac{1}{N} \sum_x x g(x),$$

$$\bar{y} = \frac{1}{N} \sum_y y h(y), \quad \sigma_x^2 = \frac{1}{N} \sum_x x^2 g(x) - \bar{x}^2,$$

$$\sigma_y^2 = \frac{1}{N} \sum_y y^2 h(y) - \bar{y}^2.$$

Ex The following table gives, according to age, the frequency of marks obtained by 100 students in an intelligence

test. Ages → Marks ↓	18	19	20	21	Total
10-20	4	2	2	—	8
20-30	5	4	6	4	19
30-40	6	8	10	11	35
40-50	4	4	6	8	22
50-60	—	2	4	4	10
60-70	—	2	3	1	6
Total	19	22	31	28	100

Calculate the Correlation Coefficient.

Sol: Let $U = X - 19$, $V = \frac{Y - 35}{10}$

Correlation table

V	Mid Value	U Age X → Marks ↓	-1 18	0 19	1 20	2 21	h(v)	vh(v)	v ² h(v)	Σ uv f(u,v)
-2	15	10-20	4 ⁽⁸⁾	2 ⁽⁰⁾	2 ⁽⁴⁾	—	8	-16	32	4
-1	25	20-30	5 ⁽⁵⁾	4 ⁽⁰⁾	6 ⁽⁶⁾	4 ⁽⁸⁾	19	-19	19	-9
0	35	30-40	6 ⁽⁰⁾	8 ⁽⁰⁾	10 ⁽⁰⁾	11 ⁽⁰⁾	35	0	0	0
1	45	40-50	4 ⁽⁴⁾	4 ⁽⁰⁾	6 ⁽⁶⁾	8 ⁽¹⁶⁾	22	22	22	18
2	55	50-60	—	2 ⁽⁰⁾	4 ⁽⁸⁾	4 ⁽¹⁶⁾	10	20	40	24
3	65	60-70	—	2 ⁽⁰⁾	3 ⁽⁹⁾	1 ⁽⁶⁾	6	18	54	15
g(u)			19	22	31	28	100	25	167	52
ug(u)			-19	0	31	56	68			
u ² g(u)			19	0	31	112	162			
Σ uvf(u,v)			9	0	13	30	52			

$$\bar{U} = \frac{1}{N} \sum u g(u) = \frac{68}{100} = 0.68,$$

$$\bar{V} = \frac{1}{N} \sum v h(v) = \frac{25}{100} = 0.25$$

$$\sigma_{UV}^2 = \frac{1}{N} \sum u^2 g(u) - \bar{U}^2 = \frac{162}{100} - (0.68)^2 = 1.1576$$

$$\sigma_v^2 = \frac{1}{N} \sum_v v^2 g(v) =$$

$$\sigma_v^2 = \frac{1}{N} \sum_v v^2 h(v) - \bar{v}^2 = \frac{167}{100} - (0.25)^2 = 1.6075$$

$$\sigma_{uv} = \frac{1}{N} \sum_u \sum_v uv f(u,v) - \bar{u} \bar{v} = \frac{52}{100} - 0.68 \times 0.25 = 0.35$$

$$\therefore r(u,v) = \frac{0.35}{\sqrt{1.6075 \times 1.1576}} = \frac{0.35}{1.36} = 0.26$$

Since, correlation coefficient is independent of change of origin and scale,

$$r(X,Y) = r(u,v) = 0.26.$$

Ex The joint pdf of X and Y is given below:

Y \ X	-1	1	$h(y)$
0	$\frac{1}{8}$	$\frac{3}{8}$	$\frac{4}{8}$
1	$\frac{2}{8}$	$\frac{2}{8}$	$\frac{4}{8}$
$g(x)$	$\frac{3}{8}$	$\frac{5}{8}$	1

Find the correlation coefficient between X and Y.

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

Where $\sigma_{XY} = \frac{1}{N} \sum_x \sum_y xy f(x,y) - \bar{x} \bar{y}$

$$\sigma_X^2 = \frac{1}{N} \sum_x x^2 g(x) - \bar{x}^2, \quad \sigma_Y^2 = \frac{1}{N} \sum_y y^2 h(y) - \bar{y}^2$$

$$\bar{x} = \frac{1}{N} \sum x g(x) = \frac{1}{8} \left[-\frac{3}{8} + \frac{5}{8} \right] = \frac{2}{8} = \frac{1}{4}$$

$$\bar{y} = \frac{1}{4}, \quad \bar{y} = \frac{1}{N} \sum y h(y) = \frac{4}{8} = \frac{1}{2}$$

$$\sigma_{xy} = 0 - \frac{2}{8} + 0 + \frac{2}{8} = 0 - \frac{1}{8} = -\frac{1}{8}$$

$$\sigma_x^2 = \frac{3}{8} + \frac{5}{8} - \frac{1}{16} = \frac{15}{16},$$

$$\sigma_y^2 = 0 + \frac{4}{8} - \frac{1}{4} = \frac{2}{8} = \frac{1}{4}$$

$$\therefore \rho(x, y) = \frac{-\frac{1}{8}}{\sqrt{\frac{15}{16} \times \frac{1}{4}}} = \frac{-0.125}{0.48} = -0.26.$$

$$\rho(x, y) = -0.26$$

Rank Correlation

Spearman's Rank Correlation Coefficient

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

Ex The ranks of some 16 students in Mathematics and Physics are as follows. Two numbers within brackets denote the ranks of the students in Mathematics and Physics:

(1,1), (2,10), (3,3), (4,4), (5,5), (6,7), (7,2), (8,6), (9,8),
(10,11), (11,15), (12,9), (13,14), (14,12), (15,16), (16,13).

Calculate the rank correlation coefficient for expertise of this group in Mathematics and Physics.

<u>Sol.</u>	<u>Ranks in Maths</u> (X)	<u>Ranks in Physics</u> (Y)	<u>d = X - Y</u>	<u>d²</u>
1	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

<u>(X)</u>	<u>(Y)</u>	<u>d = X - Y</u>	<u>d²</u>
14	12	2	4
15	16	-1	1
16	13	3	9
			<u>136</u>

$$\therefore \rho = 1 - \frac{6(136)}{16(256-1)} = 1 - \frac{816}{4080} = 1 - 0.2 = 0.8$$

$$\boxed{\rho = 0.8}$$

Ex Ten competitors in a musical test were ranked by the three judges A, B and C in the following order:

<u>Ranks by A</u> :	1	6	5	10	3	2	4	9	7	8
<u>Ranks by B</u> :	3	5	8	4	7	10	2	1	6	9
<u>Ranks by C</u> :	6	4	9	8	1	2	3	10	5	7

Discuss which pair of judges has the nearest approach to Common likings in music.

<u>Solⁿ</u>	<u>X</u>	<u>Y</u>	<u>Z</u>	<u>d₁ = X - Y</u>	<u>d₂ = X - Z</u>	<u>d₃ = Y - Z</u>	<u>d₁²</u>	<u>d₂²</u>	<u>d₃²</u>
	1	3	6	-2	-5	-3	4	25	9
	6	5	4	1	2	1	1	4	1
	5	8	9	-4	-4	-1	16	16	1
	10	4	8	6	2	-4	36	4	16
	3	7	1	4	2	6	16	4	36
	2	10	2	-8	0	8	64	0	64
	4	2	3	2	1	-1	4	1	1
	9	1	10	8	-1	-9	64	1	81
	7	6	5	1	2	1	1	4	1
	8	9	7	-1	1	2	1	1	4
							<u>200</u>	<u>60</u>	<u>214</u>

$$s(X, Y) = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} = 1 - \frac{6(200)}{10(99)} = 1 - \frac{1200}{990} = 1 - 1.21 = -0.21$$

$$s(X, Y) = -0.21$$

$$s(X, Z) = 1 - \frac{6(60)}{990} = 1 - \frac{360}{990} = 1 - 0.36 = 0.64$$

$$s(Y, Z) = 1 - \frac{6(214)}{990} = 1 - 1.297 = -0.297$$

Since, $s(X, Z)$ is maximum, we conclude that the pair of judges A and C has the nearest approach to common likings in music.

Tied Ranks / Repeated Ranks

$$s = \frac{1 - \frac{6(\sum d^2 + \text{Correction factor})}{n(n^2-1)}}{n(n^2-1)}$$

To $\sum d^2$, we add $\frac{n(n^2-1)}{12}$ for each repeated value.

Ex: Obtain the rank correlation coefficient for the following data:

	68	64	75	50	64	80	75	40	55	64
X:										
Y:	62	58	68	45	81	60	68	48	50	70

Sol ∴

<u>X</u>	<u>Y</u>	^(x) <u>Rank X</u>	^(y) <u>Rank Y</u>	<u>d = x - y</u>	<u>d²</u>
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
				<u>Σd = 0</u>	<u>72</u>

$$\rho = 1 - \frac{6 \left(\Sigma d^2 + 2 + \frac{6}{12} + \frac{6}{12} \right)}{10(99)}$$

$$= 1 - \frac{6 \left(72 + 2 + \frac{1}{2} + \frac{1}{2} \right)}{990}$$

$$= 1 - \frac{450}{990} = 1 - 0.45 = 0.55$$

$$\boxed{\rho = 0.55}$$

$$(64) \quad \frac{n(n^2-1)}{12} = \frac{3(9-1)}{12} = \frac{24}{12} = 2$$

$$(75) \quad \frac{2(4-1)}{12} = \frac{6}{12}$$

$$(68) \quad \frac{6}{12}$$

Remark Limits for Rank Correlation coefficient:
 $-1 \leq \rho \leq 1$

Remarks

1. $\sum d = 0$

2. Karl Pearson's Correlation Coefficient assumes that the parent population from which sample observations are drawn is normal, whereas Spearman's Correlation Coefficient is non-parametric.

Regression

Regression is the study of the nature of relationship between the variables so that one may be able to predict the unknown value of one variable for a known value of another variable.

In regression, one variable is considered as an independent variable and another variable is taken as dependent variable.

Regression is used to predict the value of dependent variable on the basis of the value of the independent variable.

Linear 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". 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.

As the line of regression gives the best estimate of the value of one variable for any specific value of the other variable, it is called line of "best fit".

Line of Regression of Y on X is

$$Y = a + bX \quad b \rightarrow \text{slope of line, change in } y \text{ co-ordinate with change in } x \text{ co-ordinate}$$

and is given by

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

Where $b_{YX} = r \frac{\sigma_Y}{\sigma_X}$ is the regression coefficient of Y on X .

Line of Regression of X on Y is

$$X = a + bY$$

and is given by

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

where $b_{XY} = r \frac{\sigma_X}{\sigma_Y}$ is the regression coefficient of X on Y .

Ex Obtain the equations of two lines of regression for the following data. Also obtain the estimate of X

for $Y = 70$.

X :	65	66	67	67	68	69	70	72
Y :	67	68	65	68	72	72	69	71

Sol: Let $U = X - 68$, $V = Y - 69$.

$$\bar{U} = 0, \bar{V} = 0, \sigma_U^2 = 4.5, \sigma_V^2 = 5.5, \sigma_{UV} = 3,$$

$$r(U, V) = 0.604$$

Since, correlation coefficient is independent of change of origin, $r(X, Y) = r(U, V) = 0.604$

Now, $\bar{U} = \bar{X} - 68$, $\bar{V} = \bar{Y} - 69$

$\Rightarrow \bar{X} = 68, \bar{Y} = 69$

Since, Variance is independent of change of origin,

$\therefore \sigma_x^2 = \sigma_u^2 = 4.5, \sigma_y^2 = \sigma_v^2 = 4.5$

$\Rightarrow \sigma_x = \sqrt{4.5} = 2.12, \sigma_y = 2.34$

Equation of line of regression of Y on X is

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

where $b_{YX} = r \frac{\sigma_y}{\sigma_x} = (0.604) \frac{(2.34)}{2.12} = 0.67$

$\therefore Y - 69 = 0.67(X - 68)$

$\Rightarrow Y = 0.67X - 45.33 + 69$

$\Rightarrow \boxed{Y = 0.67X + 23.66}$

Equation of line of regression of X on Y is

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

where $b_{XY} = r \frac{\sigma_x}{\sigma_y} = (0.604) \frac{2.12}{2.34} = 0.55$

$\therefore X - 68 = 0.55(Y - 69)$

$\Rightarrow X = 0.55Y - 37.95 + 68$

$\Rightarrow \boxed{X = 0.55Y + 30.05}$

If Y is 70, estimated value of X is

$\hat{X} = 0.55(Y) + 30.05 = 0.55(70) + 30.05$

$\hat{X} = 68.55$, where \hat{X} is estimate of X.

Properties of Regression coefficients

- (1) Correlation coefficient is the geometric mean between the regression coefficients.

$$r^2 = b_{xy} \times b_{yx} \Rightarrow r = \pm \sqrt{b_{xy} \times b_{yx}}$$

* b_{xy} , b_{yx} , r all have same sign.

- (2) If one of the regression coefficient is greater than 1 then other must be less than one. $[\because r^2 \leq 1]$

- (3) The modulus value of the arithmetic mean of regression coefficients is not less than the modulus value of the correlation coefficient r .

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

- (4) Regression coefficients are independent of the change of origin but not of scale.

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

$$b_{xy} = \frac{h}{k} b_{uv} \quad \text{and} \quad b_{yx} = \frac{k}{h} b_{vu}$$

- (5) The point of intersection of both regression lines is (\bar{X}, \bar{Y}) .

Ex In a partially destroyed laboratory, record of an analysis of correlation data, the following results only are readable:

Variance of $X = 9$. Regression equations: $8X - 10Y + 66 = 0$, $40X - 18Y = 214$. What are (i) the mean values X and Y , (ii) the correlation coefficient between X and Y , and (iii) the standard deviation of Y ?

Sol. (i) Since, both lines pass through (\bar{X}, \bar{Y}) .

$$\therefore 8\bar{X} - 10\bar{Y} + 66 = 0 \Rightarrow 4\bar{X} - 5\bar{Y} = -33 \quad \text{--- (1)}$$

$$40\bar{X} - 18\bar{Y} = 214 \Rightarrow 20\bar{X} - 9\bar{Y} = 107 \quad \text{--- (2)}$$

$$\textcircled{1} \times 5 \Rightarrow 20\bar{X} - 25\bar{Y} = -165 \quad \text{--- (3)}$$

$$\textcircled{2} - \textcircled{3} \text{ provides, } 16\bar{Y} = 272$$

$$\Rightarrow \bar{Y} = 17$$

$$\text{and } 4\bar{X} = -33 + 85 = 52$$

$$\Rightarrow \bar{X} = 13$$

$\therefore (\bar{X}, \bar{Y})$ is $(13, 17)$.

So, $\bar{X} = 13, \bar{Y} = 17$.

(ii) Let $8X - 10Y + 66 = 0$ and $40X - 18Y = 214$ be the lines of regression of Y on X and X on Y , respectively.

$$\begin{array}{l} \text{X on Y} \\ 8X - 10Y + 66 = 0 \\ X = \frac{10Y - 66}{8} \end{array}$$

$$\therefore b_{XY} = \frac{10}{8}$$

$$\therefore Y = \frac{8}{10}X + \frac{66}{10} \text{ and } X = \frac{18}{40}Y + \frac{214}{40}$$

$$\begin{array}{l} \text{Y on X} \\ 40X - 18Y = 214 \end{array}$$

$$\Rightarrow Y = \frac{40X - 214}{18}$$

$$\therefore b_{YX} = \frac{8}{10}, b_{XY} = \frac{18}{40}$$

$$\Rightarrow b_{YX} = \frac{4}{5}, b_{XY} = \frac{9}{20}$$

$$b_{YX} = \frac{40}{18}$$

$$\therefore r^2 = \frac{4}{5} \cdot \frac{9}{20} = 0.36$$

$$\Rightarrow r = \pm 0.6$$

Since, both b_{xy} and b_{yx} are positive,

$$r = 0.6$$

$$r^2 = b_{xy} b_{yx} = \frac{10}{8} \cdot \frac{40}{18}$$

$$r^2 = 2.78$$

$$\text{But, } 0 < r^2 < 1$$

$$(iii) \quad b_{xy} = r \frac{\sigma_x}{\sigma_y} \Rightarrow \frac{9}{20} = (0.6) \frac{3}{\sigma_y} \Rightarrow \sigma_y = 1.8 \times \frac{20}{9}$$

$$\Rightarrow \boxed{\sigma_y = 4}$$

$$\text{or } b_{yx} = r \frac{\sigma_y}{\sigma_x} \Rightarrow \frac{8}{10} = (0.6) \frac{\sigma_y}{3} \Rightarrow \sigma_y = \frac{8}{10} \times \frac{3}{0.6}$$

$$\boxed{\sigma_y = 4}$$

Ex Find the most likely price in Mumbai corresponding to the price of Rs. 70 at Kolkata from the following:

	Kolkata	Mumbai
Average price	65	67
Standard deviation	2.5	3.5
Correlation coefficient between the prices of commodities in the two cities is 0.8.		

Sol: Let the prices in Kolkata and Mumbai be denoted by X and Y respectively. We have,

$$\bar{X} = 65, \bar{Y} = 67, \sigma_x = 2.5, \sigma_y = 3.5, r(X, Y) = 0.8$$

Line of regression of Y on X is

$$Y - \bar{Y} = r \frac{\sigma_y}{\sigma_x} (X - \bar{X}) \Rightarrow Y - 67 = 0.8 \left(\frac{3.5}{2.5} \right) (X - 65)$$

$$\Rightarrow Y - 67 = 1.12X - 72.8$$

$$\Rightarrow Y = 1.12X - 5.8$$

$$\text{When } X = 70, \hat{Y} = 1.12(70) - 5.8 = 72.6$$

Hence, the most likely price in Mumbai corresponding to the price of Rs. 70 at Kolkata is Rs. 72.6.

Ex Can $Y = 5 + 2.8X$ and $X = 3 - 0.5Y$ be the estimated regression equations of Y on X and X on Y respectively?

Sol.: $Y = 5 + 2.8X \Rightarrow b_{YX} = 2.8$

$$X = 3 - 0.5Y \Rightarrow b_{XY} = -0.5$$

This is not possible as the signs of b_{XY} and b_{YX} should be same.

Angle between two lines of Regression

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

Case i) If $r = 0 \Rightarrow \tan \theta = \infty \Rightarrow \theta = \frac{\pi}{2}$.

Thus, if two variables are uncorrelated, the lines of regression become perpendicular to each other.

Case ii) If $r = \pm 1$. $\tan \theta = 0 \Rightarrow \theta = 0 \text{ or } \pi$

In this case two lines of regression either coincide

or they are parallel to each other.

But, since both lines of regression pass through the point (\bar{x}, \bar{y}) , they cannot be parallel.

Hence, in the case of perfect correlation, positive or negative, two lines of regression coincide.

Fitting of a curve

(a) Fitting a straight line: $Y = a + bX$

$$\sum Y = na + b \sum X$$

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

} Normal equations

(b) Fitting a parabola: $Y = a + bX + cX^2$

$$\sum y = na + b \sum x + c \sum x^2$$

$$\sum xy = a \sum x + b \sum x^2 + c \sum x^3$$

$$\sum x^2 y = a \sum x^2 + b \sum x^3 + c \sum x^4$$

Properties of Correlation Coefficient

1. $-1 \leq r \leq 1$ or $|\text{Cov}(X, Y)| \leq \sigma_X \sigma_Y$

Proof $\therefore r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})}{\left[\frac{1}{n} \sum (x_i - \bar{x})^2 \frac{1}{n} \sum (y_i - \bar{y})^2 \right]^{1/2}}$

$$\Rightarrow r^2(X, Y) = \frac{\left(\sum a_i b_i \right)^2}{\sum a_i^2 \sum b_i^2}, \text{ where } \begin{matrix} a_i = x_i - \bar{x}, \\ b_i = y_i - \bar{y} \end{matrix}$$

By Schwartz inequality,

$$\left(\sum a_i b_i \right)^2 \leq \sum a_i^2 \sum b_i^2$$

$$\therefore r^2(X, Y) \leq 1 \Rightarrow |r(X, Y)| \leq 1 \Rightarrow -1 \leq r(X, Y) \leq 1.$$

$$\Rightarrow \left| \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \right| \leq 1$$

$$\Rightarrow |\text{Cov}(X, Y)| \leq \sigma_X \sigma_Y$$

2. Correlation Coefficient is independent of change of origin and change of scale.

Proof \therefore Let $U = \frac{X-a}{h}$, $V = \frac{Y-b}{k}$.

$$\therefore X = a + hU, Y = b + kV, \text{ where } a, b, h, k \text{ are constants;}$$

$$h > 0, k > 0.$$

T.P. $\therefore r(X, Y) = r(U, V).$

Since, $X = a + hU$, $Y = b + kV$.

Taking expectations on both sides, we get,

$$E(X) = a + hE(U), \quad E(Y) = b + kE(V).$$

$$\therefore X - E(X) = h[U - E(U)]$$

$$Y - E(Y) = k[V - E(V)]$$

$$\Rightarrow \text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))]$$

$$= hk E[(U - E(U))(V - E(V))] = hk \sigma_{UV}$$

$$\sigma_X^2 = E[(X - E(X))^2] = h^2 E[(U - E(U))^2]$$

$$\sigma_Y^2 = E[(Y - E(Y))^2] = k^2 E[(V - E(V))^2]$$

$$\Rightarrow \sigma_X^2 = h^2 E[(U - E(U))^2] \text{ and } \sigma_Y^2 = k^2 E[(V - E(V))^2]$$

$$\sigma_X = h \sigma_U \text{ and } \sigma_Y = k \sigma_V.$$

$$\rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(U - E(U))(V - E(V))]}{E[U - E(U)] E[V - E(V)]}$$

$$= \frac{\sigma_{UV}}{\sigma_U \sigma_V} = \text{Cov}(X, Y).$$

$$= \rho(U, V).$$

$$\therefore \rho(X, Y) = \rho(U, V).$$

3. Two independent variables are uncorrelated.

Proof - If X and Y are independent variables, then

$$\text{Cov}(X, Y) = 0 \Rightarrow \rho(X, Y) = 0.$$

Hence two independent variables are uncorrelated.

But, the converse may not be true; i.e., two uncorrelated variables may not be independent.

For ex-1

<u>X</u>	-3	-2	-1	1	2	3	$\Sigma X = 0$
<u>Y</u>	9	4	1	1	4	9	$\Sigma Y = 28$
<u>XY</u>	-27	-8	-1	1	8	27	$\Sigma XY = 0$

$$\bar{X} = \frac{1}{n} \Sigma X = 0, \text{Cov}(X, Y) = \frac{1}{n} \Sigma XY - \bar{X}\bar{Y} = 0.$$

$$\therefore \rho(X, Y) = 0.$$

Thus, in the above example X and Y are uncorrelated.
But, X and Y are not independent but related by the relation $Y = X^2$.

Properties of Regression Coefficients

- (1) Correlation coefficient is the geometric mean between the regression coefficients.

Proof: $b_{xy} = r \frac{\sigma_x}{\sigma_y}$ and $b_{yx} = r \frac{\sigma_y}{\sigma_x}$

$$\therefore b_{xy} \times b_{yx} = r^2 \Rightarrow r = \pm \sqrt{b_{xy} \times b_{yx}}$$

Remark: The sign of correlation coefficient is the same as that of regression coefficients.

- (2) If one of the regression coefficients is greater than unity, the other must be less than unity.

Proof: Let one of the regression coefficients, say, b_{yx} be greater than unity.

$$b_{yx} > 1 \Rightarrow \frac{1}{b_{yx}} < 1$$

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

$$\text{Hence, } b_{xy} \leq \frac{1}{b_{yx}} < 1$$

- (3) The modulus value of the ~~regression~~ arithmetic mean of the regression coefficients is not less than the modulus value of the correlation coefficient.

Proof: T.P.: $\left| \frac{b_{xy} \times b_{yx}}{2} \right| \geq |r|$

$$\Rightarrow \left| \frac{1}{2} \left(2 \frac{\sigma_x}{\sigma_y} + 2 \frac{\sigma_y}{\sigma_x} \right) \right| \geq |2|$$

$$\Rightarrow \left| \frac{2}{2} \left(\frac{\sigma_x}{\sigma_y} + \frac{\sigma_y}{\sigma_x} \right) \right| \geq |2|$$

$$\Rightarrow \frac{\sigma_x}{\sigma_y} + \frac{\sigma_y}{\sigma_x} \geq 2$$

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

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

$$\Rightarrow (\sigma_x - \sigma_y)^2 \geq 0, \text{ which is always true,}$$

since the square of a real quantity ≥ 0 .

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

Proof: Let $U = \frac{X-a}{h}$, $V = \frac{Y-b}{k}$, where $a, b, h(>0)$ and $k(>0)$ are constants.

$$\Rightarrow X = a + hU, Y = b + kV.$$

$$\text{Cov}(X, Y) = hk \text{Cov}(U, V)$$

$$\sigma_x^2 = h^2 \sigma_u^2 \text{ and } \sigma_y^2 = k^2 \sigma_v^2$$

$$b_{yx} = r \frac{\sigma_y}{\sigma_x} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \frac{\sigma_y}{\sigma_x} = \frac{\sigma_{xy}}{\sigma_x^2} = \frac{hk \sigma_{uv}}{h^2 \sigma_u^2} = \frac{k}{h} \frac{\sigma_{uv}}{\sigma_u^2}$$

$$= \frac{k}{h} \frac{\sigma_{xy}}{hk \sigma_u^2}$$

$$b_{vu} = r \frac{\sigma_v}{\sigma_u} = \frac{\sigma_{uv}}{\sigma_u \sigma_v} \frac{\sigma_v}{\sigma_u} = \frac{\sigma_{uv}}{\sigma_u^2}$$

$$\therefore b_{yx} = \frac{k}{h} b_{vu}$$

$$\text{Similarly, } b_{xy} = \frac{h}{k} b_{uv}$$

(5) Point of intersection of both regression lines is (\bar{X}, \bar{Y}) .