

# Correlation and Regression

Page: 10/10

Date: 11/11/11

Correlation  $\rightarrow$  We know that the quantity of measure of relationship between  $X$  and  $Y$  is given by covariance between  $X$  and  $Y$ .

Here Covariance between  $X$  and  $Y$  =  $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$   
and  $-\infty < \text{Cov}(X, Y) < \infty$ ,

Correlation b/w  $X$  and  $Y$  =  $r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$

where  $\sigma_X$  = SD of  $X$  and  $\sigma_Y$  = SD of  $Y$

Limits for correlation coefficient

Let  $E(X) = \mu_X$  and  $E(Y) = \mu_Y$

$$\text{Now consider } E\left[\left(\frac{X-\mu_X}{\sigma_X} + \frac{Y-\mu_Y}{\sigma_Y}\right)^2\right] \geq 0$$

$$E\left[\left(\frac{X-\mu_X}{\sigma_X}\right)^2 + \left(\frac{Y-\mu_Y}{\sigma_Y}\right)^2 + 2\left(\frac{X-\mu_X}{\sigma_X}\right)\left(\frac{Y-\mu_Y}{\sigma_Y}\right)\right] \geq 0$$

$$E\left[\left(\frac{X-\mu_X}{\sigma_X}\right)^2\right] + E\left[\left(\frac{Y-\mu_Y}{\sigma_Y}\right)^2\right] + 2E\left[\frac{(X-\mu_X)(Y-\mu_Y)}{\sigma_X \sigma_Y}\right] \geq 0$$

$$1 + 1 + 2\text{Cov}(X, Y) \geq 0$$

$$1 + 1 + 2r(X, Y) \geq 0$$

$$-1 \leq r(X, Y) \leq 1$$

$$E\left[\left(\frac{X-\mu_X}{\sigma_X}\right)^2\right]$$

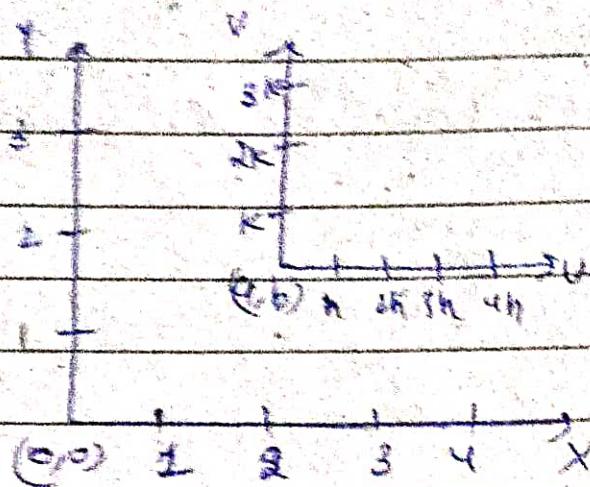
$$+ E\left[\left(\frac{Y-\mu_Y}{\sigma_Y}\right)^2\right]$$

$$\frac{\sigma_X^2}{\sigma_X^2} + \frac{\sigma_Y^2}{\sigma_Y^2} = 1$$

Note Correlation coefficient is independent of origin and scale.

Let  $U = \frac{X-a}{h}$ ,  $V = \frac{Y-b}{k}$  where  $a, b, h, k$  are constants.

Then  $\rho(X, Y) = \rho(U, V)$



Formulas to calculate correlation coefficient

(i) For ungrouped data  $\rightarrow$  If  $X$  takes the values

$x_i, i=1, 2, \dots, n$  and  $Y$  takes

the values  $y_i, i=1, 2, \dots, n$  then  $(x_i, y_i), i=1, 2, \dots, n$  are the  $n$  pairs of observations.

$$E(X) = \frac{\sum_{i=1}^n x_i}{n}$$

$$E(Y) = \frac{\sum_{i=1}^n y_i}{n}$$

$$E(X, Y) = \frac{\sum_{i=1}^n x_i y_i}{n}$$

$$\sigma_x = \sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2 - (E(X))^2}$$

$$\sigma_y = \sqrt{\frac{1}{n} \sum_{i=1}^n y_i^2 - (E(Y))^2}$$

correlation coefficient b/w X and Y

$$\Rightarrow \rho_{XY} = \frac{E(XY) - E(X)E(Y)}{\sigma_x \sigma_y}$$

(2) For grouped data (bivariate frequency distribution table)

Let  $x_i, i=1, 2, \dots, n$  be the mid points of the  $n$  class intervals defined for  $X$  with respective frequencies  $f_{xi}, i=1, 2, \dots, n$ .

Let  $y_j, j=1, 2, \dots, m$  be the mid points of  $m$  class intervals defined for  $Y$  with respective frequencies  $f_{yj}, j=1, 2, \dots, m$

Let  $f_{ij}$  be the frequency of  $i$ th class interval of  $X$  correlated with  $j$ th class interval of  $Y$

$$\text{Let } N = \sum_{i=1}^n f_{xi} = \sum_{j=1}^m f_{yj} = \sum_{i=1}^n \sum_{j=1}^m f_{ij}$$

Now we can calculate

$$E(X) = \frac{\sum_{i=1}^n f_{xi} x_i}{N} \quad E(Y) = \frac{\sum_{j=1}^m f_{yj} y_j}{N} \quad E(XY) = \frac{\sum_{i=1}^n \sum_{j=1}^m f_{ij} x_i y_j}{N}$$

$$\sigma_x = \sqrt{\frac{1}{N} \sum_{i=1}^n f_{xi} x_i^2 - [E(X)]^2} \quad \sigma_y = \sqrt{\frac{1}{N} \sum_{j=1}^m f_{yj} y_j^2 - [E(Y)]^2}$$

$$\text{correlation coeff b/w X and Y} = \rho_{XY} = \frac{E(XY) - E(X)E(Y)}{\sigma_x \sigma_y}$$

Problems

(Q1) Calculate the correlation coefficient for the following heights in inches of fathers (X) and their children (Y)

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

$$\text{So } r_{xy} = \frac{E(XY) - E(X)E(Y)}{\sigma_x \sigma_y}$$

n=8      8 pairs of information

$$E(X) = \frac{\sum_{i=1}^8 X_i}{n} = \frac{544}{8} \quad E(Y) = \frac{\sum_{i=1}^8 Y_i}{n} = \frac{552}{8}$$

$$E(XY) = \frac{\sum_{i=1}^8 x_i y_i}{n} = \frac{3756}{8} \quad \sum x_i^2 = 3728$$

$$\sum y_i^2 = 38132$$

$$E(X) = 68 \quad E(Y) = 69 \quad E(XY) = 4695$$

$$\sigma_x = \sqrt{1213} \quad \sigma_y = \sqrt{3452}$$

$$r_{xy} = 0.603$$

(Q2) To calculate the correlation coefficient bw two variables X and Y from 25 pairs of observations, the following data is obtained.

$$\sum X = 125 \quad \sum Y = 100, \quad \sum X^2 = 650, \quad \sum Y^2 = 460,$$

$$\sum XY = 508$$

Later it was discovered that the two pairs of data values (6,14), (8,6) are considered instead of

correct pair of values  $(8, 18)$ ,  $(6, 8)$  <sup>cancel</sup>  
 Hence obtain correlation coeff for the data.

Given  $\sum X = 125$  and  $\sum Y = 100$  will remain same  
 or  $6+8=14$       or  $14+6=20$   
 or  $8+6=14$       or  $12+8=20$

$$\sum X^2 = 650 \text{ remains}$$

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

$$8^2 = 64 \\ 6^2 = 36$$

$$\sum Y^2 = 436$$

$$\sum XY = 508 - 84 - 48 + 96 + 48 = 520$$

$$E(X) = \frac{\sum X}{n} = \frac{125}{25} = 5$$

$$E(Y) = \frac{\sum Y}{n} = \frac{100}{25} = 4$$

$$E(XY) = \frac{520}{n} = \frac{520}{25} = 20.8$$

$$\sigma_X = \sqrt{\frac{1}{n} \sum x_i^2 - (E(X))^2}$$

$$= \sqrt{\frac{650}{25} - 25}$$

$$\sqrt{26 - 25} = 1$$

$$\sigma_Y = \sqrt{\frac{1}{n} \sum y_i^2 - (E(Y))^2}$$

$$\sigma_Y = \sqrt{\frac{436}{25} - 16}$$

$$\sigma_Y = \sqrt{19.81 - 16} = \sqrt{3.81} \\ = 1.95$$

$$\rho(X, Y) = \frac{20.8 - 20}{1 \times 1.95} = \frac{0.8}{1.95} = 0.41$$

Q3 Calculate the correlation coefficient for the following bivariate frequency distribution table.

mid Y	20	30	40	50	60	70	Total fr <sub>2,01</sub>	
mid X	4	15-25	25-35	35-45	45-55	55-65	65-75	X
20	15-25	1	f <sub>11</sub>	f <sub>12</sub>				2 f <sub>m1</sub>
30	25-35	2	f <sub>21</sub>	f <sub>22</sub>				2 f <sub>m2</sub>
40	35-45		4	10	1			15 f <sub>m3</sub>
50	45-55			3	6	1		10 f <sub>m4</sub>
60	55-65				2	4	2	8 f <sub>m5</sub>
70	65-75					1	2	3 f <sub>m6</sub>
Total fr <sub>2,01</sub>		3	26	14	9	6	4	
of Y		f <sub>y1</sub>	f <sub>y2</sub>	f <sub>y3</sub>	f <sub>y4</sub>	f <sub>y5</sub>	f <sub>y6</sub>	

N = 62

$$E(X) = \frac{\sum_{i=1}^n f_{xi} x_i}{N}$$

$$E(Y) = \frac{\sum_{j=1}^m f_{yj} y_j}{N}$$

$$E(XY) = \sum_{i=1}^n \sum_{j=1}^m f_{ij} x_i y_j$$

n = 6 m = 6

$$\sigma_x = \sqrt{\frac{1}{N} \sum_{i=1}^n f_{xi} x_i^2 - [E(X)]^2}$$

$$\sigma_y = \sqrt{\frac{1}{N} \sum_{j=1}^m f_{yj} y_j^2 - [E(Y)]^2}$$

$$\sum_{i=1}^n f_{xi} x_i = 2280$$

$$\sum_{j=1}^m f_{yj} y_j = 2220$$

$$\sum_{i=1}^n \sum_{j=1}^m f_{ij} x_i y_j = 103800$$

$$\sum_{i=1}^n f_{xi} x_i^2 = 106800$$

$$\sum_{j=1}^m f_{yj} y_j^2 = 102600$$

(Q) Calculate the correlation coefficient for the following bivariate frequency distribution table.

60.5<sup>m1</sup> 64.5<sup>m2</sup> 68.5<sup>m3</sup> 72.5<sup>m4</sup> 76.5<sup>m5</sup>

$y$	59-62	63-66	67-70	71-74	75-78	Total freq. $f_y$
$x$						$\Sigma f_x$
100 <sup>m1</sup>	2	1				3 $f_{m1}$
120 <sup>m2</sup>	7	8	4	2		21 $f_{m2}$
140 <sup>m3</sup>	5	15	22	7	1	50 $f_{m3}$
160 <sup>m4</sup>	2	12	63	19	5	101 $f_{m4}$
180 <sup>m5</sup>		7	28	32	12	79 $f_{m5}$
200 <sup>m6</sup>		2	10	20	7	39 $f_{m6}$
220 <sup>m7</sup>			1	4	2	7 $f_{m7}$
Total freq. of $y$	16	45				
					$  h=7   m=5  $	

# Lines of Regression

Let  $(x_i, y_i)$ ,  $i=1, 2, \dots, n$  be the given set of  $n$  points

Aim → We want to get a straight line passing through the given set of  $n$  points.

If not, we want to get a best fit of straight line for the given set of  $n$  points.

Let  $y = a + bx$  be the best fit of straight line for the given set of  $n$  points if  $x$  is dependent on  $X$  and this best fit of straight line is known as regression line of  $y$  on  $x$ .

That means, we want to find values of  $a$  and  $b$  such that

$$S = \sum_{i=1}^n (y_i - (a + bx_i))^2 \text{ is minimum}$$

Here  $x_i$  and  $y_i$  are known and  $(a, b)$  are unknown.

$$\textcircled{1} f(n) \rightarrow f'(n) = 0 \Rightarrow \text{roots}$$

$$f''(n) \Big|_{\text{roots}} = \begin{cases} > 0 & \text{at that root } f(n) \text{ is min} \\ < 0 & \text{at that root } f(n) \text{ is max} \end{cases}$$

$$\textcircled{2} f(n, y) \rightarrow \boxed{\frac{\partial f}{\partial n} = 0, \frac{\partial f}{\partial y} = 0} \rightarrow \text{solving this}$$

we get stationary points

$$\boxed{\frac{\partial^2 f}{\partial n^2}, \frac{\partial^2 f}{\partial y^2}, \frac{\partial^2 f}{\partial n \partial y}} \rightarrow \text{evaluate these at stationary point}$$

(i) If  $\frac{\partial^2 f}{\partial x^2} \Big|_{st. point} > 0$ ,  $\left[ \frac{\partial^2 f}{\partial x^2} \frac{\partial^2 f}{\partial y^2} - \left( \frac{\partial^2 f}{\partial x \partial y} \right)^2 \Big|_{st. point} > 0 \right]$

at that st. point  $f(x,y)$  is minimum

(ii) If  $\frac{\partial^2 f}{\partial x^2} \Big|_{st. point} < 0$ ,  $\left[ \frac{\partial^2 f}{\partial x^2} \frac{\partial^2 f}{\partial y^2} - \left( \frac{\partial^2 f}{\partial x \partial y} \right)^2 \Big|_{st. point} < 0 \right]$

at that st. point  $f(x,y)$  is maximum

From we have

$$S = \sum_{i=1}^n (y_i - (a + b x_i))^2 = S(a, b)$$

Now  $S$  is a function of  $a$  and  $b$

for  $S$  to be minimum, the necessary condition are  $\frac{\partial S}{\partial a} = 0$ ,  $\frac{\partial S}{\partial b} = 0$ ,

$$\frac{\partial S}{\partial a} = 0$$

$$\frac{\partial S}{\partial a} = \sum_{i=1}^n (a + b x_i) - \sum_{i=1}^n y_i = 0 \quad \text{--- (1)}$$

$$\frac{\partial S}{\partial b} = \sum_{i=1}^n x_i (a + b x_i) - \sum_{i=1}^n x_i y_i = 0 \quad \text{--- (2)}$$

Solving (1) and (2) we get  $a$  and  $b$

$$b = \frac{\sum x_i y_i}{\sum x_i^2} \quad \gamma = \text{coff of correlation} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

$$a = \frac{1}{n} \left( \sum y_i - b \sum x_i \right)$$

$$a = \frac{\sum y_i}{n} - b \frac{\sum x_i}{n} = E(Y) - b E(X)$$

$$\Rightarrow E(Y) = a + b E(X) \quad \text{Here } (E(X), E(Y)) \text{ is a point of st. line}$$

$$y = a + b X \quad y = a + b X$$

Now we'll check condition for minimum by using values of  $a$  and  $b$

for these values of  $(a, b)$ ,  $S$  will be minimum if

$$\left[ \frac{\partial S}{\partial a} \Big|_{(a,b)} = 0 \text{ and } \left[ \frac{\partial^2 S}{\partial a^2} \frac{\partial^2 S}{\partial b^2} - \left( \frac{\partial S}{\partial a \partial b} \right)^2 \right]_{(a,b)} > 0 \right]$$

Hence

$$\frac{\partial S}{\partial a} = n>0 \text{ and}$$

$$\left[ \frac{\partial^2 S}{\partial a^2} \frac{\partial^2 S}{\partial b^2} - \left( \frac{\partial S}{\partial a \partial b} \right)^2 \right]_{(a,b)} = n \sigma_x^2 - (\bar{x}_i)^2 = n \left[ \sigma_x^2 - \left( \frac{\sum x_i}{n} \right)^2 \right] = n \sigma_x^2 > 0$$

Therefore at the values of  $(a, b)$ ,  $S$  is minimum.

### Note

① If  $\bar{x} = \frac{\sum x_i}{n}$  and  $\bar{y} = \frac{\sum y_i}{n}$  then regression line 'Y on X' is given by  $y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$

② If  $X$  is dependent on  $Y$  and the best fit of straight line  $x = c + d y$  is known as regression line of  $X$  on  $Y$

The values of  $c$  and  $d$  can be found by using relations

$$d = r \frac{\sigma_x}{\sigma_y} \quad \text{and} \quad c = \frac{1}{n} \left( \sum x_i - d \sum y_i \right)$$

Also we can get the regression line of  $X$  and  $Y$  by using the equation

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

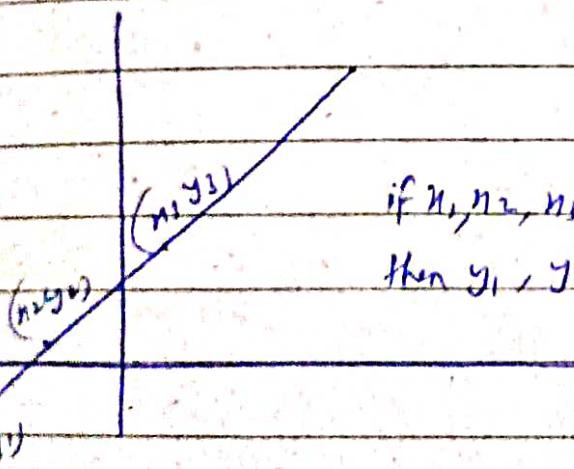
- ③ The point  $(\bar{x}, \bar{y})$  is the point of intersection of both the regression lines.
- ④ The correlation coefficient is the G.M (geometric mean) of two regression coefficients (slope of regression lines).

$$r = \pm \sqrt{\left( \frac{\partial(\bar{y})}{\partial x} \right) \left( \frac{\partial(\bar{x})}{\partial y} \right)} = \pm \sqrt{\delta^2 \pm r}$$

If both the regression coefficients are +ve, then we'll take  $r > 0$

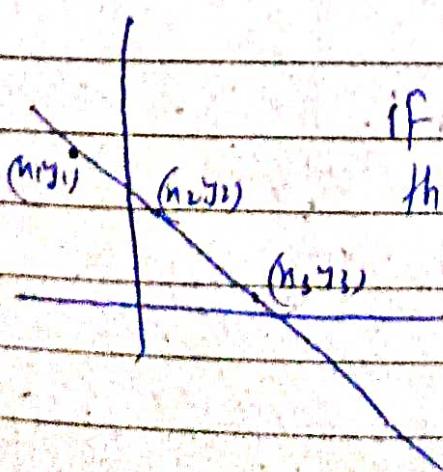
If both the regression coefficients are -ve, then we'll take  $r < 0$

$r > 0$  if x and y are directly proportional



if  $n_1, n_2, n_3, \dots = T$   
then  $y_1, y_2, y_3, \dots = T$

if  $n_1, n_2, n_3, \dots = +$   
then  $y_1, y_2, y_3, \dots = +$



if  $n_1, n_2, n_3, \dots = \downarrow$   
then  $y_1, y_2, y_3, \dots = T$

if  $n_1, n_2, n_3, \dots = T$   
then  $y_1, y_2, y_3, \dots = \downarrow$

$r < 0$  if x and y are inversely proportional

# Problems on Regression Lines

Page No. 100

- ① Find the correlation coefficient b/w X and Y for the following data. Also find regression lines

X	1	2	3	4	5	6	7	8	9	10
Y	10	12	16	28	25	36	41	49	40	50

To find

$$\text{① } r = \frac{E(XY) - E(X)E(Y)}{\sigma_x \sigma_y}$$

$$\text{② Regression line of } Y \text{ on } X \Rightarrow y - \bar{y} = \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$\text{③ Regression line of } X \text{ on } Y \Rightarrow x - \bar{x} = \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$\sum x = 55 \quad \sum y = 307 \quad \sum x^2 = 385 \quad \sum y^2 = 11387$$

$$\sum xy = 2074$$

$$\bar{x} = E(X) = \frac{\sum x}{n} = \frac{55}{10} = 5.5 \quad \bar{y} = E(Y) = \frac{\sum y}{n} = \frac{307}{10} = 30.7$$

$$E(XY) = \frac{\sum xy}{n} = 20.74$$

$$\sigma_x = \sqrt{\frac{1}{n} \sum x_i^2 - (\bar{x})^2} = 2.872$$

$$\sigma_y = \sqrt{\frac{1}{n} \sum y_i^2 - (\bar{y})^2} = 14.007$$

$$r = 0.9583$$

$$\text{Regression line of } Y \text{ on } X \Rightarrow y - \bar{y} = \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

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

Q. Develop the estimating equation that best describes the following data and predict Y for X

Page No. 10  
Date 15/11/2011

X	13	16	14	11	17	9	13	17	18	12
Y	6.2	8.6	7.2	4.5	9.0	3.5	6.5	9.3	9.5	5.7

Soln We have to find regression line of Y on X

$$Y = a + bX \quad \text{or} \quad Y - \bar{Y} = \frac{\sigma_{xy}}{\sigma_x} (X - \bar{X})$$

$$\sum X = 140 \quad \sum Y = 70 \quad \sum X^2 = 2038 \quad \sum XY = 1035$$

$$\text{reg. line} \Rightarrow Y = -2.87 + 0.705X$$

$$b = \frac{1}{n} \sum n_i x_i - \left( \frac{\sum n_i}{n} \right) \left( \frac{\sum y_i}{n} \right) = \frac{\text{cov}(X, Y)}{\sigma_x^2} = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} \cdot \frac{\sigma_y}{\sigma_x}$$

$$\frac{1}{n} \sum n_i^2 - \left( \frac{\sum n_i}{n} \right)^2 = \frac{\sigma_{xy}}{\sigma_x}$$

$$a = \frac{1}{n} \left( \sum y_i - b \sum x_i \right)$$

~~$$b = \frac{1}{10} \sum n_i x_i - \left( \frac{140}{10} \right) \left( \frac{70}{10} \right) = \frac{103.5 - 98}{203.8 - 196} = \frac{5.5}{7.8} = 0.705$$~~

$$\frac{1}{10} \sum n_i^2 - \left( \frac{140}{10} \right)^2 = 0.705$$

$$a = \frac{1}{10} \left( 70 - 0.705 \times 140 \right) = \frac{1}{10} (70 - 98.7) = -2.87$$

$$Y = -2.87 + 0.705X$$

Q3 Cost accountants often estimate overhead based on the level of production. At the Standard Knitting Co., they have collected information on overhead expenses and units produced at different plants and want to estimate a regression equation to predict future overhead.

$\downarrow Y$	Overhead expenses	191	171	272	185	280	173	234	116	153	178
$\downarrow X$	Units	40	42	53	35	56	39	48	30	37	40

X Develop the regression equation for the cost accountants

Both regression lines of  $Y$  on  $X$        $X$  is independent  
 $Y$  is dependent

If we are not able to figure out which is dependent or independent then we find both regression lines of  $Y$  on  $X$  as well as  $X$  on  $Y$ .

(i)  $n=10 \quad \sum X = 420 \quad \sum Y = 1922 \quad \sum XY = 84541$   
 $\sum X^2 = 18228$

Calculate  $a = -80.443$

and  $b = 6.49$

$$[Y = a + bX]$$

Q4 It is given that variance of  $X = 9$  and the two regression lines are  $8X - 10Y + 66 = 0$ ,  $90X - 18Y - 214 = 0$ . Calculate

- the mean values of  $X$  and  $Y$
- the correlation coeff blw  $X$  and  $Y$
- the standard deviation of  $Y$

variance  $X = 9$

$$\sigma_x^2 = 9$$

$$\sigma_x = 3$$

Page No. 10/10

Date 11/11/2023

i) the mean values of  $X$  and  $Y$

and we know mean values of  $X$  and  $Y$  are the intersection of these two lines

$$8X - 10Y + 66 = 0 \quad \text{Line 1}$$

$$40X - 18Y - 214 = 0 \quad \text{Line 2}$$

$$40X - 50Y + 330 = 0$$

$$8X - 10Y + 66 = 0$$

$$40X - 18Y - 214 = 0$$

$$8X - 17Y + 66 = 0$$

$$- 32Y + 594 = 0$$

$$8X = 10Y$$

$$Y = \frac{594}{32} = 17 \quad (1)$$

$$X = \frac{10Y}{8} = 13$$

$$Y = 17$$

$$X = 13$$

$$E(X) = 13$$

$$E(Y) = 17$$

(ii)

Let take  $8X - 10Y + 66 = 0$  as reg line  $X$  on  $Y$

$$8X = 10Y - 66$$

$$\frac{10}{8} = b = \frac{8\sigma_x}{\sigma_y}$$

$$X = \frac{10Y - 66}{8}$$

Based on this

$40X - 18Y - 214 = 0$  is reg line  $Y$  on  $X$

$$18Y = 40X - 214$$

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

$$\frac{40}{18} = b = \frac{8\sigma_y}{\sigma_x}$$

Since both slopes are  
true so  $b_1 > 0$

$$S = \sqrt{\left(\frac{\sigma_{\sigma_x}}{\sigma_y}\right) \left(\frac{\sigma_{\sigma_y}}{\sigma_x}\right)} = \sqrt{\frac{10}{8} \times \frac{40}{18}} = \sqrt{\frac{25}{4}} = \sqrt{2.7} > 0$$

But this is wrong  
as  $S$  can not be  $> 1$

So our prediction is wrong

Another case

$$8X - 10Y + 66 = 0 \quad \text{as reg line } Y \text{ on } X$$

$$8X + 66 = 10Y$$

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

On this basis

$$40X - 18Y - 214 = 0 \quad \text{as reg line } X \text{ on } Y$$

$$40X = 18Y + 214$$

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

Since both slopes are +ve. So  $\rho$  will also be +ve.

$$\rho_1 = \sqrt{\frac{8}{10} \times \frac{18}{40}} = \sqrt{\frac{9}{25}} < 1$$

$$\frac{3}{5} = 0.6 < 1$$

(iii) SD of  $Y$  ( $\sigma_Y$ )

from reg line  $Y$  on  $X$ , slope is  $\frac{8}{10}$

$$\frac{8}{10} = \rho \frac{\sigma_Y}{\sigma_X}$$

$$\frac{8}{10} = 0.6 \frac{\sigma_Y}{3} \Rightarrow \frac{8 \times 3}{10 \times 0.6} = \sigma_Y$$

$$\sigma_Y = 4 \quad \text{Any} =$$

The variable  $X$  and  $Y$  are connected by the equation  $ax + by + c = 0$ . Show that the correlation coefficient  $b/w$  them is  $-1$  if signs of  $a$  and  $b$  are alike and  $+1$  if they are different.

Soln Connected by the equation means every point lies on this straight line  $ax + by + c = 0$ .

$$ax + by + c = 0 \quad \text{--- (1)}$$

$$E[ax + by + c] = E(0)$$

$$aE(X) + bE(Y) + c = 0 \quad \text{--- (2)}$$

$$\text{--- (1)} - \text{--- (2)}$$

$$a(X - E(X)) + b(Y - E(Y)) = 0$$

$$a(X - E(X)) = -b(Y - E(Y))$$

$$X - E(X) = -\frac{b}{a} [Y - E(Y)]$$

$$\text{rl} = \text{correlation coeff} = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{E[(X - E(X))(Y - E(Y))]}{\sigma_X \sigma_Y}$$

Now  $E[(X - E(X))(Y - E(Y))]$

$$\Rightarrow E\left[-\frac{b}{a}(Y - E(Y))(Y - E(Y))\right] = -\frac{b}{a} E[(Y - E(Y))^2]$$

$$= -\frac{b}{a} \sigma_Y^2$$

$$\sigma_X^2 = E[(X - E(X))^2] = E\left[\frac{b^2}{a^2} (Y - E(Y))^2\right]$$

$$\sigma_X^2 = \frac{b^2}{a^2} E[(Y - E(Y))^2] = \frac{b^2}{a^2} \sigma_Y^2$$

$$\sigma_X = \left| \frac{b}{a} \right| \sigma_Y$$

$$\rho = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{-b \sigma_Y^2}{\sigma_X \sigma_Y}$$

# if  $b$  and  $a$  of same sign

$$\frac{b > 0}{a} \Rightarrow \left| \frac{b}{a} \right| = \frac{b}{a}$$

$$\rho = \frac{-b}{a} = \textcircled{-1} \quad \text{if } a \text{ and } b \text{ of } \cancel{\text{same}} \text{ sign}$$

# if  $a$  and  $b$  of diff sign

$$\frac{b < 0}{a} \Rightarrow \left| \frac{b}{a} \right| = -\frac{b}{a}$$

$$\rho = \frac{-b}{a} = \textcircled{+1} \quad \text{if } a \text{ and } b \text{ are of } \cancel{\text{diff}} \text{ sign.}$$

MP