

## Module - III Correlation And Regression

### Correlation

Suppose two variables  $x$  and  $y$  are related in such a way that an increase in the value of one of these variables is accompanied by an increase or decrease in the value of the other. Such a relationship is called correlation.

If the values of  $x$  and  $y$  increase or decrease together, then we say that  $x$  and  $y$  are positively correlated. If the value of  $y$  decreases as the value of  $x$  increases or vice versa, then we say that  $x$  &  $y$  are negatively correlated.

The numerical measure of correlation between two variables  $x$  &  $y$  is known as Pearson's coefficient of correlation and is defined as

$$r = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} \longrightarrow \textcircled{1}$$

where  $X = x - \bar{x}$  and  $Y = y - \bar{y}$

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

$\sigma_y^2 = \frac{\sum (y - \bar{y})^2}{n}$ , employing these expressions in  $\textcircled{1}$

We get alternate formula

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

### Alternative formula for the correlation

$$r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x\sigma_y}$$

proof: Let  $z = x - y$

$$\frac{\sum z}{n} = \frac{\sum x}{n} - \frac{\sum y}{n}$$

$$\text{or } \bar{z} = \bar{x} - \bar{y}$$

$$\text{Hence } (z - \bar{z}) = (x - \bar{x}) - (y - \bar{y})$$

Squaring both side, taking summation and dividing by  $n$ , we've

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

$$\Rightarrow \sigma_z^2 = \sigma_x^2 + \sigma_y^2 - 2r\sigma_x\sigma_y$$

$$\sigma_{x-y}^2 = \sigma_x^2 + \sigma_y^2 - 2r\sigma_x\sigma_y$$

$$\Rightarrow \sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2 = 2r\sigma_x\sigma_y$$

or

$$\boxed{r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x\sigma_y}}$$

## Problems

① Find the correlation coefficient for the following data:

x	1	2	3	4	5
y	2	5	3	8	7

Soln:- Here  $n=5$ , we find that  $\bar{x} = \frac{\sum x_i}{n} = \frac{1}{5} [1+2+3+4+5]$

$$\boxed{\bar{x} = 3}$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{1}{5} [2+5+3+8+7]$$

$$\boxed{\bar{y} = 5}$$

We now prepare the following table

$x_i$	$x_i - \bar{x}$	$x_i^2$	$y_i$	$y_i - \bar{y}$	$y_i^2$	$x_i y_i$
1	-2	4	2	-3	9	6
2	-1	1	5	0	0	0
3	0	0	3	-2	4	0
4	1	1	8	3	9	3
5	2	4	7	2	4	4
		$\sum x_i^2 = 10$			$\sum y_i^2 = 26$	13

Now we find that

$$r = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \sum y_i^2}} = \frac{13}{\sqrt{10 \times 26}}$$

$$\boxed{r = 0.806}$$



- 2) The following table gives the intelligence ratio (I.R) and engineering ability (E.A) of 10 students obtained on the basis of psychological tests. calculate the coefficient of correlation.

I.R (x):	105	104	102	101	100	99	98	96	93	92
E.A (y):	101	103	100	98	95	96	104	92	97	94

Sol:- Here  $n = 10$ ,  $\bar{x} = \frac{\sum x_i}{n} = \frac{990}{10} = 99$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{980}{10} = 98$$

We now prepare the following table

$x_i$	$x_i = x_i - \bar{x}$	$x_i^2$	$y_i$	$y_i = y_i - \bar{y}$	$y_i^2$	$x_i y_i$
105	6	36	101	3	9	18
104	5	25	103	5	25	25
102	3	9	100	2	4	6
101	2	4	98	0	0	0
100	1	1	95	-3	9	-3
99	0	0	96	-2	4	0
98	-1	1	104	6	36	-6
96	-3	9	92	-6	36	18
93	-6	36	97	-1	1	6
92	-7	49	94	-4	16	28
		170			140	92

We've

$$r = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2 \sum y_i^2}} = \frac{92}{\sqrt{170 \times 140}} = 0.59$$

(3)

③ Employ the formula  $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x\sigma_y}$  do

determine the correlation coefficient 'r' for the following data:

x	92	89	87	86	83	77	71	63	53	50
y	86	83	91	77	68	85	52	82	37	57

Soln:- Here  $n=10$ ,  $\bar{x} = \frac{751}{10} = 75.1$  ;  $\bar{y} = \frac{718}{10} = 71.8$

Let  $z_i = x_i - y_i$ ,  $\bar{z} = \frac{1}{10} \sum z_i = \frac{1}{10} \sum (x_i - y_i)$

$$\bar{z} = \frac{33}{10} = 3.3$$

$$\sigma_x^2 = \frac{1}{n} \sum x_i^2 - (\bar{x})^2 = \frac{1}{10} \left[ (92)^2 + (89)^2 + (86)^2 + (83)^2 + (77)^2 + (71)^2 + (63)^2 + (53)^2 + (50)^2 \right] - (75.1)^2$$

$$= \frac{58487}{10} - (75.1)^2 = 208.69$$

$$\sigma_y^2 = \frac{1}{n} \sum y_i^2 - (\bar{y})^2 = \frac{54390}{10} - (71.8)^2 = 283.76$$

$$\sigma_{x-y}^2 = \sigma_z^2 = \frac{1}{n} \sum z_i^2 - (\bar{z})^2 = \frac{1485}{10} - (3.3)^2 = 137.61$$

$$\text{Now } r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x\sigma_y} = \frac{208.69 + 283.76 - 137.61}{2 \times \sqrt{208.69 \times 283.76}}$$

$$r = 0.729$$

## Regression

Suppose we are given  $n$  pairs of values  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  of two variables  $x$  &  $y$ . If we fit a straight line to this data by taking  $x$  as independent variable and  $y$  as dependent variable, then the straight line obtained is called the line of regression of  $y$  on  $x$ . Similarly, if we fit a straight line to the data by taking  $y$  as independent variable and  $x$  as dependent variable, the line obtained is the line of regression of  $x$  on  $y$ .

We've line of regression of  $y$  on  $x$  is

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

Line of regression of  $x$  on  $y$  is

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

Note: We've alternative formula

(i)  $r = \sqrt{(\text{coeff of } x)(\text{coeff of } y)}$  for coefficient of correlation

(ii)  $y = \frac{\sum xy}{\sum x^2} x$  &  $x = \frac{\sum xy}{\sum y^2} y$  for lines of regression



prove that if  $\theta$  is the angle between the lines of regression then  $\tan \theta = \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \left( \frac{1-r^2}{r} \right)$

Pr:- W.K.T if  $\theta$  is acute, the angle b/w the lines  $y = m_1 x + c_1$  &  $y = m_2 x + c_2$  is given by

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

We've  $m_1 = r \frac{\sigma_y}{\sigma_x} ; m_2 = r \frac{\sigma_x}{\sigma_y} = \frac{1}{r} \cdot \frac{\sigma_y}{\sigma_x} \rightarrow \textcircled{2}$

Sub  $\textcircled{2}$  in  $\textcircled{1}$

$$\tan \theta = \frac{\frac{1}{r} \frac{\sigma_y}{\sigma_x} - r \frac{\sigma_y}{\sigma_x}}{1 + r \cdot \frac{\sigma_y}{\sigma_x} \cdot \frac{1}{r} \cdot \frac{\sigma_y}{\sigma_x}}$$

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

## problems

(7)

- ① Find the coefficient of correlation and obtain the lines of regression for the following data:

$x:$	1	2	3	4	5	6	7	8	9
$y:$	9	8	10	12	11	13	14	16	15

Obtain an estimate for  $y$  which corresponds to  $x=6.2$

Sol.:- we've  $n=9$ ,  $\bar{x} = \frac{1}{n} \sum x_i = \frac{45}{9} = 5$ ;  $\bar{y} = \frac{108}{9} = 12$

we prepare the following table

$x_i$	$x_i = x_i - \bar{x}$	$x_i^2$	$y_i$	$y_i = y_i - \bar{y}$	$y_i^2$	$x_i y_i$
1	-4	16	9	-3	9	12
2	-3	9	8	-4	16	12
3	-2	4	10	-2	4	4
4	-1	1	12	0	0	0
5	0	0	11	-1	1	0
6	1	1	13	1	1	1
7	2	4	14	2	4	4
8	3	9	16	4	16	12
9	4	16	15	3	9	12
		60			60	57



Thus correlation coefficient is

$$r = \frac{\sum XY}{\sqrt{\sum X_i^2 \sum Y_i^2}} = \frac{57}{60} = 0.95$$

Also  $\sigma_x^2 = \frac{\sum X_i^2}{n} = \frac{60}{9} = 6.6667 \Rightarrow \sigma_x = \sqrt{6.6667}$   
 $= 2.5819$

$$\sigma_y^2 = \frac{\sum Y_i^2}{n} = \frac{60}{9} = 6.6667 = \sigma_y$$

Thus line of regression of  $y$  on  $x$

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

$$y - 12 = 0.95 \cdot \frac{2.5819}{2.5819} (x - 5)$$

$$y - 12 = 0.95x - 4.75$$

$$\Rightarrow y = 0.95x - 4.75 + 12$$

$$\boxed{y = 0.95x + 7.25} \Rightarrow \boxed{y_{x=6.2} = 13.14}$$

line of regression of  $x$  on  $y$  is

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

$$x - 5 = 0.95 (y - 12)$$

$$\Rightarrow \boxed{x = 0.95y - 6.4}$$

- ② Obtain the lines of regression and hence find the coefficient of correlation for the data:

x	1	3	4	2	5	8	9	10	13	15
y	8	6	10	8	12	16	16	10	32	32

Soln. Here  $n = 10$ ,  $\bar{x} = \frac{70}{10} = 7$ ;  $\bar{y} = 15$

We've  $X = x - \bar{x}$ ;  $Y = y - \bar{y}$

We now prepare the foll. table

x	y	$X = x - \bar{x}$	$Y = y - \bar{y}$	$X^2$	$Y^2$	$XY$
1	8	-6	-7	36	49	42
3	6	-4	-9	16	81	36
4	10	-3	-5	9	25	15
2	8	-5	-7	25	49	35
5	12	-2	-3	4	9	6
8	16	1	1	1	1	1
9	16	2	1	4	1	2
10	10	3	-5	9	25	-15
13	32	6	17	36	289	102
15	32	8	17	64	289	136
				204	818	360

We've line of regression of y on x

$$Y = \frac{\sum XY}{\sum X^2} \cdot X$$

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

$$y - 15 = \frac{360}{204} (x - 7)$$

$$\boxed{y = 1.76x + 2.68}$$

Also line of regression of  $x$  on  $y$  is

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

$$x - 7 = \frac{360}{818} (y - 15)$$

$$\boxed{x = 0.44y + 0.4}$$

Thus

correlation

Coefficient

$$= r = \sqrt{(\text{coeff of } x)(\text{coeff of } y)}$$

$$r = \sqrt{(1.76)(0.44)}$$

$$\boxed{r = 0.8813}$$



③ Given

	x-Series	y-Series
Mean	18	100
SD	14	20

$$r = 0.8$$

Write down the equation of the lines of regression and hence find the most probable value only when  $x = 70$

Sol: By data  $\bar{x} = 18$ ,  $\bar{y} = 100$   
 $\sigma_x = 14$   $\sigma_y = 20$

Line of regression of  $y$  on  $x$  is

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

$$y - 100 = 0.8 \cdot \frac{20}{14} (x - 18)$$

$$\Rightarrow y = 1.14x + 79.48$$

when  $x = 70$ ;  $y = 159.28$

Line of regression of  $x$  on  $y$  is

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

$$x - 18 = 0.8 \frac{14}{20} (y - 100)$$

$$x = 0.56y - 3.8$$

④ Given  $8x - 10y + 66 = 0$ ,  $40x - 18y = 214$  are two regression lines. Find the means of  $x$ 's &  $y$ 's and correlation coefficient. Find  $\sigma_y$  if  $\sigma_x = 3$ ?

Sol: - We know that lines of regression pass through

$$\bar{x} \text{ \& } \bar{y} \Rightarrow 8\bar{x} - 10\bar{y} = -66$$

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

on solving  $\boxed{\bar{x} = 13, \bar{y} = 17}$

Regression coefficients are

$$8x - 10y + 66 = 0 \quad ; \quad 40x - 18y = 214$$

$$10y = 8x + 66$$

$$\boxed{y = (0.8)x + 6.6}$$

$$40x = 18y + 214$$

$$\boxed{x = (0.45)y + 5.35}$$

Thus coefficient of correlation  $\left. \vphantom{\begin{matrix} \text{coefficient of} \\ \text{correlation} \end{matrix}} \right\} = r = \sqrt{(\text{coeff of } x)(\text{coeff of } y)}$

$$= \sqrt{(0.8)(0.45)}$$

$$\boxed{r = 0.6}$$

By data  $\sigma_x = 3$  ;  $r \frac{\sigma_y}{\sigma_x} = \text{Regression line } y \text{ on } x$

$$0.6 \cdot \frac{\sigma_y}{3} = 0.8$$

$$\Rightarrow \sigma_y = \frac{0.8 \times 3}{0.6} = \underline{\underline{4}}$$

- ⑤ If the coefficient of correlation between two variables  $x$  &  $y$  are 0.5, acute angle b/w their lines of regression is  $\tan^{-1}(3/5)$ . S.T  $\sigma_y = 2\sigma_x$  or  $\sigma_x = 2\sigma_y$

Sol: Given  $r = 0.5$ ,  $\theta = \tan^{-1}\left(\frac{3}{5}\right) \Rightarrow \tan \theta = \frac{3}{5}$

We've

$$\tan \theta = \frac{\frac{\sigma_x \sigma_y}{2}}{\sigma_x^2 + \sigma_y^2} \left[ \frac{1-r^2}{r} \right] \quad \text{Since } r = \frac{1}{2}$$

$$\frac{3}{5} = \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \left( \frac{3}{4} \cdot 2 \right)$$

$$\Rightarrow \frac{1}{5} = \frac{\sigma_x \sigma_y}{2(\sigma_x^2 + \sigma_y^2)}$$

$$\Rightarrow 2\sigma_x^2 + 2\sigma_y^2 = 5\sigma_x \sigma_y$$

$$\Rightarrow 2\sigma_x^2 - 5\sigma_x \sigma_y + 2\sigma_y^2 = 0$$

$$(2\sigma_x - \sigma_y)(\sigma_x - 2\sigma_y) = 0$$

$$\Rightarrow \boxed{2\sigma_x = \sigma_y}$$

$$\text{or } \boxed{\sigma_x = 2\sigma_y}$$



- ⑥ In a bivariate distribution  $\sigma_x = \sigma_y$  and the angle b/w the regression lines is  $\tan^{-1}(3)$ . Find the correlation coefficient.

Soln:- Given  $\theta = \tan^{-1}(3) \Rightarrow \tan \theta = 3$  &  $\sigma_x = \sigma_y$

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

$$3 = \frac{\sigma_x^2}{2\sigma_x^2} \left( \frac{1-r^2}{r} \right)$$

$$\Rightarrow 6r = 1-r^2$$

$$\text{or } r^2 + 6r - 1 = 0$$

$$r = 0.1623 \text{ or } r = -6.6123$$

But  $|r| \leq 1 \Rightarrow \boxed{r = 0.1623}$

- ⑦ If  $x$  &  $y$  are random variables with SD's,  $\sigma_x$  &  $\sigma_y$ . It was found that random variables  $x+y$ ,  $x-y$ ,  $2x+y$  respy have variance 15, 11, 29. compute SD's of  $x$  &  $y$  respy and also the coefficient of correlation.

(11)

Sol: By data  $\sigma_{x+y}^2 = 15$ ,  $\sigma_{x-y}^2 = 11$ ,  $\sigma_{2x+y}^2 = 29$

we've  $\sigma_{ax+by}^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2r ab \sigma_x \sigma_y \rightarrow (1)$

Taking  $(a, b) = (1, 1) \quad (1, -1) \quad \& \quad (2, 1)$   
 $\{x+y \quad x-y \quad 2x+y\}$

Eqn (1)  $\Rightarrow \sigma_x^2 + \sigma_y^2 + 2r \sigma_x \sigma_y = 15 \rightarrow (2)$

$$\sigma_x^2 + \sigma_y^2 - 2r \sigma_x \sigma_y = 11 \rightarrow (3)$$

$$4\sigma_x^2 + \sigma_y^2 + 4r \sigma_x \sigma_y = 29 \rightarrow (4)$$

Eqns (2) + (3) gives

$$2\sigma_x^2 + 2\sigma_y^2 = 26$$

$$\sigma_x^2 + \sigma_y^2 = 13 \rightarrow (5)$$

Eqn (2)  $\times$  (3) + (4) gives

$$6\sigma_x^2 + 3\sigma_y^2 = 51$$

$$2\sigma_x^2 + \sigma_y^2 = 17 \rightarrow (6)$$

on solving (5) & (6)

$$\sigma_x^2 = 4 \quad ; \quad \sigma_y^2 = 9$$

$$\boxed{\sigma_x = 2 \quad ; \quad \sigma_y = 3}$$

sub these values in (2)

$$4 + 9 + 2r \cdot 2 \times 3 = 15$$

$$\boxed{r = 0.17}$$