

$$A = \frac{0.6655}{4}$$

taking Antilog:  $b = 1.38$   $a = 0.68$

$$\therefore \text{Required Curve: } Y = (0.68)(1.38)^x$$

### Exercises

(1) Fit an equation of the form  $Y = ab^x$  to the following data

(a)	$x:$	2	3	4	5	6	$Y = (101.3)(1.196)^x$
	$Y:$	144	172.8	207.4	248.8	298.6	

(b)	$x:$	2	3	4	5	6
	$Y:$	8.3	15.4	33.1	65.2	127.4

Also estimate:  $Y$  when  $x = 4.5, 7$  and  $3.5$

(2) Fit an equation of the form  $Y = a e^{bx}$

$x:$	1	2	3	4	5	6
$Y:$	1.6	4.5	13.8	40.2	125.0	300

$$Y = (0.557) e^{1.05x}$$

### CORRELATION AND REGRESSION:

(1) Def: Correlation Co-efficient:

Correlation Co-efficient between two r.v's  $X$  and  $Y$  usually denoted by  $\rho(x, Y)$  (or  $r(x, Y)$ ) is a numerical measure of linear relationship between them and is defined by

$$\rho(x, Y) = \frac{\text{Cov}(x, Y)}{\sigma_x \sigma_Y}$$

If  $(x_i, y_i)$   $i=1, 2, \dots, n$  is  $n$  pair of values of  $x$  and  $Y$  then,

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

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

$$\sigma_x^2 = E[\{x - E(x)\}^2] = \frac{1}{n} \sum (x_i - \bar{x})^2$$

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

$$\therefore \rho(x, Y) = \frac{\frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})}{\left[ \frac{1}{n} \sum (x_i - \bar{x})^2, \frac{1}{n} \sum (y_i - \bar{y})^2 \right]^{1/2}}$$

$$\rho^2 = \frac{(\sum a_i b_i)^2}{\sum a_i^2 \sum b_i^2} \quad \text{where} \quad a_i = x_i - \bar{x} \quad b_i = y_i - \bar{y}. \quad (I)$$

Now we know a well known inequality: If  $a_i, b_i \quad i=1, 2, \dots, n$  are two sets of real quantities then

$$(\sum_{i=1}^n a_i b_i)^2 \leq (\sum_{i=1}^n a_i^2)(\sum_{i=1}^n b_i^2) \quad \text{Schwartz inequality}$$

' holds iff  $\frac{a_1}{b_1} = \frac{a_2}{b_2} = \dots = \frac{a_n}{b_n}$ .

Using Schwartz Inequality we have from (I)

$$\rho^2 \leq 1$$

$$\Rightarrow |\rho| \leq 1$$

$$\Rightarrow -1 \leq \rho \leq 1$$

Hence Correlation Coefficient Cannot exceed unity numerically. It always lies between +1 and -1

If  $r=1$  then  $\rho$  is perfect and positive  
 $=-1$  "  $\rho$  is perfect and negative.

obviously if  $x$  and  $y$  are independent random variables then  $\rho(x, y) = 0$  because in such case  $\text{Cov}(x, y) = 0$ .

#### PROBLEM

(1) Calculate the Correlation Coefficient from the following table:

$x$	$y$	$x^2$	$y^2$	$xy$
65	67	4225	4489	4355
66	68	4356	4624	4488
67	65	4489	4225	4355
67	68	4489	4624	4556
68	72	4624	5184	4896
69	72	4761	5184	4968
70	69	4900	4761	4830
72	71	5184	5041	5112
$\sum$	544	552	37028	38132
				37560

$$\bar{x} = \frac{1}{n} \sum x = 68 \quad \bar{y} = 69$$

$$\text{We have: } \text{Cov}(x, Y) = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})$$

$$\begin{aligned} &= \frac{1}{n} \sum (x_i y_i - \bar{x} y_i - \bar{y} x_i + \bar{x} \bar{y}) \\ &= \frac{1}{n} \sum x_i y_i - \bar{x} \frac{1}{n} \sum y_i - \bar{y} \frac{1}{n} \sum x_i + \bar{x} \bar{y} \\ &= \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y} - \bar{x} \bar{y} + \bar{x} \bar{y} \\ &= \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y} \end{aligned}$$

$$\begin{aligned} \sigma_x^2 &= \frac{1}{n} \sum (x_i - \bar{x})^2 \\ &= \frac{1}{n} \sum (x_i^2 - 2x_i \bar{x} + \bar{x}^2) \\ &= \frac{1}{n} \sum x_i^2 - 2\bar{x} \frac{1}{n} \sum x_i + \bar{x}^2 \\ &= \frac{1}{n} \sum x_i^2 - 2\bar{x}^2 + \bar{x}^2 \\ &= \frac{1}{n} \sum x_i^2 - \bar{x}^2 \end{aligned}$$

$$\sigma_Y^2 = \frac{1}{n} \sum y_i^2 - \bar{y}^2$$

$$\begin{aligned} \rho(x, Y) &= \frac{\frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}}{\sqrt{(\frac{1}{n} \sum x_i^2 - \bar{x}^2)(\frac{1}{n} \sum y_i^2 - \bar{y}^2)}} \\ &= \frac{\frac{1}{8} \times 37560 - 68 \times 69}{\sqrt{(\frac{37028}{8} - 68^2)(\frac{38132}{8} - 69^2)}} \\ &= 0.603 \end{aligned}$$

(2) A Computer while calculating Correlation Coefficient between two Variables X and Y from pairs of observation obtained the following results

$$n=25, \sum x = 125, \sum x^2 = 650, \sum y = 100, \sum y^2 = 640$$

$$\sum xy = 508$$

It was, however, later discovered at the time of checking that he was 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 the Correlation Coefficient

Ans

$$\text{Corrected } \sum x = 125 - 6 - 8 + 8 + 6 = 125$$

$$" \quad \sum y = 100 - 14 - 6 + 12 + 8 = 100$$

$$" \quad \sum x^2 = 650 - 6^2 - 8^2 + 8^2 + 6^2 = 650$$

$$" \quad \sum y^2 = \frac{640}{460} - 14^2 - 6^2 + 12^2 + 8^2 = 436.616$$

$$" \quad \sum xy = 508 - 6 \times 14 - 8 \times 6 + 8 \times 12 + 6 \times 8 = 520$$

$$\bar{x} = \frac{1}{n} \sum x = \frac{1}{25} \times 125 = 5, \quad \bar{y} = \frac{1}{25} \times 100 = 4$$

$$\text{Cov}(x, y) = \frac{1}{n} \sum xy - \bar{x} \bar{y} = \frac{1}{25} \times 520 - 5 \times 4 = \frac{4}{5} = 0.8$$

$$\sigma_x^2 = \frac{1}{n} \sum x^2 - \bar{x}^2 = \frac{1}{25} \times 650 - 5^2 = 1$$

$$\sigma_y^2 = \frac{1}{n} \sum y^2 - \bar{y}^2 = \frac{1}{25} \times 436.616 - 4^2 = \frac{216}{25} = 8.64$$

$$\text{Corrected } \rho(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{\frac{4}{5}}{\sqrt{1} \times \sqrt{5}} = \frac{2}{\sqrt{5}} = 0.87$$

$$= \frac{0.8}{\sqrt{8.64}} = \frac{0.8}{2.94} = 0.272$$

# ③ (a) If  $Z = ax + by$  and  $\rho \equiv \rho(x, y)$   
Show that  $\sigma_z^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab \rho \sigma_x \sigma_y$

(b) Show that the Correlation Coefficient between two r.v's  $x$  and  $y$  can be

written as:

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

Solution:

$$(a) \quad Z = ax + by$$

$$E(Z) = a E(x) + b E(y)$$

$$Z - E(Z) = a[x - E(x)] + b[y - E(y)]$$

Squaring and taking expectation on both sides

$$E[Z - E(Z)]^2 = a^2 E[x - E(x)]^2 + b^2 E[y - E(y)]^2$$

$$+ 2ab E\{[x - E(x)][y - E(y)]\}$$

$$\sigma_z^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab \text{Cov}(x, y)$$

$$\Rightarrow \sigma_z^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab \rho \sigma_x \sigma_y \quad [\because \rho = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}]$$

(b) In particular if we take  $a=1, b=-1$   
then  $Z = x - y$  and from result of (a)

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

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

## REGRESSION

1. Def: Regression Analysis is a mathematical measure of the average relationship between two or more variables in terms of the original units of the given data.

In regression analysis there are two types of variables. The variable whose value is to be predicted is called dependent variable and the variable which is used for prediction is called independent variable.

2. Lines of Regression:

The line of regression is the line which gives the best estimate to the value of one variable for any specific value of other variable. Therefore the line of regression is the line of best-fit and is obtained by Principles of least Square.

Consider  $n$  pair of points  $(x_i, y_i)$   $i=1, 2, \dots, n$ .  $Y$  is dependent and  $X$  is independent Variable.

Let the line of regression of  $Y$  on  $X$  be

$$Y = a + bX$$

By Principle of least square, Normal equation to estimate  $a, b$  are

$$\sum y_i = na + b \sum x_i \quad \dots (1)$$

$$\sum x_i y_i = a \sum x_i + b \sum x_i^2 \quad \dots (2)$$

$$\text{From (1)} \quad \bar{y} = a + b \bar{x} \quad \dots (3)$$

$\Rightarrow$  line of reg. of  $(Y$  on  $X$ ) passes through means  $(\bar{x}, \bar{y})$

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

$$\Rightarrow \mu_{11} = \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}$$

$$\Rightarrow \boxed{\frac{1}{n} \sum x_i y_i = \mu_{11} + \bar{x} \bar{y}} \quad \dots (4)$$

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

$$\Rightarrow \boxed{\frac{1}{n} \sum x_i^2 = \sigma_x^2 + \bar{x}^2} \quad \dots (5)$$

$$\text{From (2) using (4) and (5)}$$

$$\frac{1}{n} \sum x_i y_i = a \cdot \frac{1}{n} \sum x_i + b \cdot \frac{1}{n} \sum x_i^2$$

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

$$(3) \times \bar{x} \Rightarrow \bar{x} \bar{y} = a \bar{x} + b \bar{x}^2 \quad \dots (7)$$

$$(7) - (6) \Rightarrow \mu_{11} = b \sigma_x^2$$

$$\Rightarrow \boxed{b = \frac{\mu_{11}}{\sigma_x^2}} \quad \dots (8)$$

since  $b$  is the slope of the line of regression of  $y$  on  $x$  and since the line passing through  $(\bar{x}, \bar{y})$  its equation is

$$y - \bar{y} = \frac{\mu_{11}}{\sigma_x^2} (x - \bar{x})$$

$$\Rightarrow y - \bar{y} = p \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \quad [ \because p = \frac{\mu_{11}}{\sigma_x \sigma_y} ] \quad \dots \textcircled{9}$$

Starting with the equation  $x = a + b y$  and proceeding similarly the equation of the line of regression of  $x$  on  $y$  becomes:

$$x - \bar{x} = \frac{\mu_{11}}{\sigma_y^2} (y - \bar{y})$$

$$\Rightarrow x - \bar{x} = p \frac{\sigma_x}{\sigma_y} (y - \bar{y}) \quad \dots \textcircled{10}$$

Hence there are always two lines of regression one of  $y$  on  $x$  and other  $x$  on  $y$ . In fact  $y$  on  $x$  is used to estimate the value of  $y$  for any given value of  $x$ . while  $x$  on  $y$  is used to estimate the value of  $x$  for any given value of  $y$ . The estimate so obtained will be best in the sense that it will have the minimum possible error as defined by the principle of least square.

### Important Properties:

1. We know  $b$  the slope of the line of regression of  $y$  on  $x$  is also called Coefficient of regression of  $y$  on  $x$  and is denoted by  $b_{yx}$

$$\therefore b_{yx} = \frac{\mu_{11}}{\sigma_x^2} = p \frac{\sigma_y}{\sigma_x} \quad \dots \textcircled{11}$$

$b_{yx}$  represents the increment in the value of the dependent variable  $y$  corresponding to a unit change in the value of the independent variable  $x$ .

Similarly

$$b_{xy} = \text{Coefficient of regression of } x \text{ on } y$$

$$= \frac{\mu_{11}}{\sigma_y^2} = p \frac{\sigma_x}{\sigma_y} \quad \dots \textcircled{12}$$

$$\textcircled{11} \times \textcircled{12} \Rightarrow b_{yx} b_{xy} = p^2$$

$$\Rightarrow p = \pm \sqrt{b_{yx} b_{xy}} \quad \dots \textcircled{13}$$

$\Rightarrow$  Correlation Coefficient is the G.M of regression Coefficients.

Note that if the regression coefficients are positive  $p$  is +ve and if they are negative  $p$  is -ve so that the sign to be taken before the square root is that of the regression Coefficients.

If one regression coefficient is  $>1$  then other must  $<1$

$b_{yx} > 1$  then we have to show that

2. ~~It is possible that~~  $b_{xy} < 1$

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

Since  $\rho^2 \leq 1 \Rightarrow b_{yx} b_{xy} \leq 1$

$$\Rightarrow b_{xy} \leq \frac{1}{b_{yx}} < 1 \quad \text{*(PROVED)*}$$

$\Rightarrow$  If one reg. coefficient is  $>1$  then other must be  $<1$

3. A.M of regression coefficients is greater than the Correlation coefficient

We have to prove that,  $\frac{1}{2}(b_{yx} + b_{xy}) \geq \rho$

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

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

$$\Rightarrow (\sigma_x - \sigma_y)^2 \geq 0$$

which is always true. Hence the result.

Regression line of:

Y on X :  $y - \bar{y} = \rho \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \dots \textcircled{1}$

X on Y :  $x - \bar{x} = \rho \frac{\sigma_x}{\sigma_y} (y - \bar{y}) \dots \textcircled{2}$

$\therefore$  Slope of  $\textcircled{1}$   $m_1 = \rho \frac{\sigma_y}{\sigma_x}$

" "  $\textcircled{2}$   $m_2 = \frac{1}{\rho} \frac{\sigma_x}{\sigma_y}$

If  $\theta$  be the angle between  $\textcircled{1}$  &  $\textcircled{2}$

$$\tan \theta = \frac{\rho \frac{\sigma_y}{\sigma_x} - \frac{1}{\rho} \frac{\sigma_x}{\sigma_y}}{1 + \frac{\sigma_y^2}{\sigma_x^2}} = \frac{\rho^2 - 1}{\rho} \left( \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right)$$

$$= \frac{1 - \rho^2}{\rho} \left( \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \quad [\because \rho^2 \leq 1]$$

$$\therefore \theta = \tan^{-1} \left\{ \frac{1 - \rho^2}{\rho} \left( \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right) \right\} \dots \textcircled{3}$$

(i) if  $\rho = 0$ ,  $\tan \theta = \infty \Rightarrow \theta = \frac{\pi}{2} \Rightarrow$  if two variables are uncorrelated then lines of regression becomes perpendicular to each other

(ii) if  $\rho = \pm 1$  then  $\tan \theta = 0 \Rightarrow \theta = 0$  or  $\pi$  In this case the lines of regression either coincide or they are parallel to each other. But since the lines of regression passing through  $(\bar{x}, \bar{y})$  they can not be parallel. Hence in the case of perfect Correlation +ve or -ve the regression lines are coincide

(iii) Whenever two lines of intersect, there are two angles between them, one angle acute and other obtuse.

$\tan \theta > 0$  if  $0 < \theta < \pi/2$  i.e.  $\theta$  is an acute angle.  
 $\tan \theta < 0$  if  $\pi/2 < \theta < \pi$  i.e.  $\theta$  is obtuse angle

and since  $0 < \rho^2 < 1$  and acute angle  $\theta_1$  and obtuse angle  $\theta_2$  are given by

$$\theta_1 = \text{Acute angle} = \tan^{-1} \left\{ \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \frac{1 - \rho^2}{\rho} \right\}$$

$$\theta_2 = \text{Obtuse angle} = \tan^{-1} \left\{ \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \frac{\rho^2 - 1}{\rho} \right\}$$

### PROBLEMS

#① In partially destroyed Laboratory record of an analysis of Correlation data, the following results only are legible.

$$\text{Var}(x) = 9$$

Regression equations:  $8x - 10y + 66 = 0$ ,  $40x - 18y = 214$

What are (i) the mean values of  $x$  and  $y$

$$(ii) \rho(x, y)$$

$$(iii) \sigma_y$$

### Solution

(i) Since both lines pass through  $(\bar{x}, \bar{y})$

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

$$\text{Solving: } \bar{x} = 13, \bar{y} = 17$$

(ii) Let  $8x - 10y + 66 = 0$  be line of  $y$  on  $x$   
 $40x - 18y - 214 = 0$  " " "  $x$  on  $y$

$$y = \frac{8}{10}x + \frac{66}{10}, \quad x = \frac{18}{40}y + \frac{214}{40}$$

$$b_{yx} = \frac{8}{10} = \frac{4}{5}, \quad b_{xy} = \frac{18}{40} = \frac{9}{20}$$

$$\rho^2 = \frac{8}{10} \times \frac{9}{20} = \frac{9}{25} \Rightarrow \rho = \pm 0.6$$

Since both regression Co-efficient are positive

$$P = 0.6$$

$$\text{iii) We have } b_{yx} = P \frac{\sigma_y}{\sigma_x} \Rightarrow \frac{4}{5} = \frac{3}{5} \times \frac{\sigma_y}{3} \\ \Rightarrow \sigma_y = 4.$$

Note that if we assume that  $8x - 10y + 66 = 0$ :  $x$  on  $y$   
 $40x - 18y - 214 = 0$ :  $y$  on  $x$

$$x = \frac{10}{8}y - \frac{66}{8} \quad b_{xy} = \frac{10}{8}$$

$$y = \frac{40}{18}x - \frac{214}{18} \quad b_{yx} = \frac{40}{18}$$

$P^2 = 2.78$  since  $0 < P^2 \leq 1$  this assumption is not valid here.

#(2) Find the most likely price in Mumbai corresponding to the price of Rs. 70 at Calcutta from the following

	Calcutta	Mumbai
Average Price	65	67
Standard deviation	2.5	3.5

Correlation coefficient between prices of commodities in the two cities is 0.8.

Solution

Price of Mumbai =  $\bar{Y}$   
 " " Calcutta =  $\bar{X}$

$$\bar{x} = 65 \quad \bar{y} = 67 \quad \sigma_x = 2.5 \quad \sigma_y = 3.5 \\ p(x,y) = 0.8 \quad \text{We want } Y \text{ for } x = 70$$

Line of regression of  $Y$  on  $X$

$$Y - \bar{Y} = P \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$\Rightarrow Y = 67 + 0.8 \cdot \frac{3.5}{2.5} (x - 65)$$

$$\Rightarrow Y \Big|_{x=70} = 67 + 0.8 \cdot \frac{3.5}{2.5} (70 - 65) = 72.6$$

#(3)

Given  $x = 4y + 5$  regression line of  $x$  on  $y$

$y = kx + 4$  " " "  $y$  on  $x$ .

$$\text{Then } b_{xy} = 4, \quad b_{yx} = k$$

$$\therefore r^2 = b_{xy} b_{yx} = 4k$$

$$\text{Since } 0 < r^2 < 1 \Rightarrow 0 < 4k < 1$$

(i) Prove that  $0 < k < \frac{1}{4}$   
 (ii) If  $k = \frac{1}{16}$  find the value of  $P$  (Correlation Coefficient)  
 In particular if:  $k = \frac{1}{16}$  then  $r^2 = 4 \cdot \frac{1}{16} = \frac{1}{4}$   $\bar{x}$  and  $\bar{y}$ .  
 $r = \pm \frac{1}{2}$  since both  $b_{xy}$  &  $b_{yx} > 0$  ( $\because k > 0$ )