

$$\rho(z, x) = 1 - \frac{6 \sum d_3^2}{n(n-1)} = 1 - \frac{6 \times 60}{10 \times 99} = 0.6$$

Since $\rho(z, x)$ is maximum, the pair of judges A and C have the nearest common approach.

Problems 25.4

1. Find the correlation coefficient between x and y from the given data :

$x :$	78	89	97	69	59	79	68	57
$y :$	125	137	156	112	107	138	123	108

(J.N.T.U., 2005)

2. Find the correlation co-efficient between x and y for the given value. Find also the two regression lines.

$x :$	1	2	3	4	5	6	7	8	9	10
$y :$	10	12	16	28	25	36	41	49	40	50

(Osmania, 2003 S ; V.T.U., 2000 S)

3. Find the co-efficient of correlation between industrial production and export using the following data and comment on the result.

Production (in crore tons) :	55	56	58	59	60	60	62
Exports (in crore tons) :	35	38	38	39	44	43	45

(Madras, 2000 ; Cochin, 1999)

4. Ten people of various heights as under, were requested to read the letters on a car at 25 yards distance. The number of letters correctly read is given below :

Height (in feet) :	5.1	5.3	5.6	5.7	5.8	5.9	5.10	5.11	6.0	6.1
No. of letters :	11	17	19	14	8	15	20	6	18	12

Is there any correlation between heights and visual power ?

5. Using the formula $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x\sigma_y}$, find r from the following data :

$x :$	92	89	87	86	83	77	71	63	53	50
$y :$	86	88	91	77	68	85	52	82	37	57

6. Find the correlation between x (marks in Mathematics) and y (marks in Engineering Drawing) given in the following data :

x	10—40	40—70	70—100	Total
y				
0—30	5	20	—	25
30—60	—	28	2	30
60—90	—	32	13	45
Total	5	80	15	100

7. If the coefficient of correlation between two variables x and y is 0.5 and the acute angle between the lines of regression is $\tan^{-1}(3/8)$, show that $\sigma_x = \frac{1}{2}\sigma_y$.

(V.T.U., 2006)

8. Following table gives the data on rainfall and discharge in a certain river. Obtain the line of regression of y on x .

Rainfall x (inches) : 1.53 1.78 2.60 2.95 3.42 ,

Discharge y (1000 cc) : 33.5 36.3 40.0 45.8 53.5

9. Two random variables have the regression lines with equations $3x + 2y = 26$ and $6x + y = 31$. Find mean values and the correlation coefficient between x and y .

(Madras, 2006)

Multiplying (ii) by 0.87 and subtracting from (i), we have
 $[1 - (0.87)(0.50)] \bar{x} = 19.13 - (11.64)(0.87)$ or $0.57 \bar{x} = 9.00$ or $\bar{x} = 15.79$

$$\therefore \bar{y} = 11.64 - (0.50)(15.79) = 3.74$$

\therefore Regression co-efficient of y on x is -0.50 and that of x on y is -0.87 .
Now since the co-efficient of correlation is the geometric mean between the two regression co-efficients.

$$r = \sqrt{[(-0.50)(-0.87)]} = \sqrt{(0.43)} = -0.66.$$

Example 25.16. In the following table are recorded data showing the test scores made by salesmen on an intelligence test and their weekly sales :

Salesmen	1	2	3	4	5	6	7	8	9	10	✓
Test scores	40	70	50	60	80	50	90	40	60	60	
Sales (000)	2.5	6.0	4.5	5.0	4.5	2.0	5.5	3.0	4.5	3.0	

Calculate the regression line of sales on test scores and estimate the most probable weekly sales volume if a salesman makes a score of 70.

Sol. With the help of the table below, we have

$$\bar{x} = \text{mean of } x \text{ (test scores)} = 60 + 0/10 = 60.$$

$$\bar{y} = \text{mean of } y \text{ (sales)} = 4.5 + (-4.5)/10 = 4.05.$$

Regression line of sales (y) on scores (x) is given by

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

where

$$r \frac{\sigma_y}{\sigma_x} = \frac{\Sigma XY}{\sigma_x \sigma_y} \times \frac{\sigma_y}{\sigma_x} = \frac{\Sigma XY}{(\sigma_x)^2} = \left[\Sigma d_x d_y - \frac{\Sigma d_x \Sigma d_y}{n} \right] / \left[\Sigma d_x^2 - (\Sigma d_x)^2/n \right]$$

$$= \frac{140 - 0 \times (-4.5)}{2400 - 0^2/10} = \frac{140}{2400} = 0.06$$

\therefore The required regression line is

$$y - 4.05 = 0.06(x - 60) \quad \text{or} \quad y = 0.06x + 0.45.$$

For $x = 70$, $y = 0.06 \times 70 + 0.45 = 4.65$.

Thus the most probable weekly sales volume for a score of 70 is 4.65.

Test scores	Sales	Deviation of x from assumed mean ($= 60$)	Deviation of y from assumed average ($= 4.5$)	$d_x \times d_y$	d_x^2	d_y^2
x	y	d_x	d_y	$d_x \times d_y$	d_x^2	d_y^2
40	2.5	-20	-2	40	400	4
70	6.0	10	1.5	15	100	2.25
50	4.5	-10	0	0	100	0
60	5.0	0	0.5	0	0	0.25
80	4.5	20	0	0	400	0
50	2.0	-10	-2.5	25	100	6.25
90	5.5	30	1	30	900	1.00
40	3.0	-20	-1.5	30	400	2.25
60	4.5	0	0	0	0	0
60	3.0	0	-1.5	0	0	2.25
		$\Sigma d_x = 0$	$\Sigma d_y = -4.5$	$\Sigma d_x d_y = 140$	$\Sigma d_x^2 = 2400$	$\Sigma d_y^2 = 18.25$

$$\rho_{xy} = \frac{\sum xy - \bar{x}\bar{y}}{\sqrt{\sum x^2 - (\bar{x})^2} \sqrt{\sum y^2 - (\bar{y})^2}} = \frac{450 - 10 \times 10}{\sqrt{100 - 100}} \times \sqrt{100 - 100} = 0.0$$

Since $\rho_{xy} = 0$ is minimum, the pair of judges A and C have the lowest consistency agreement.

PROBLEMS 25-8

1. Find the correlation coefficient between x and y from the given data:

x	75	80	85	90	95	100	105	110	115	120	125
y	120	127	135	142	147	152	157	162	167	172	178

(M.T.U., 1998)

2. Find the correlation coefficient between x and y for the given value. Find also the two regression lines:

x	1	2	3	4	5	6	7	8	9	10
y	30	32	35	38	35	36	41	42	43	39

(Banaras, 2002 S., V.T.U., 2000-01)

3. Find the co-efficient of correlation between industrial production and export using the following data and comment on the result.

Production (in crore tons)	55	58	59	58	60	60	62
Exports (in crore tons)	35	38	36	38	44	43	46

(Madras, 2000, Cochran, 1998)

4. Ten people of various heights as under, were requested to read the letters on a car at 25 yards distance. The number of letters correctly read is given below:

Height (in feet)	5.1	5.3	5.6	5.7	5.8	5.9	5.10	5.11	6.0	6.1
No. of letters	11	17	19	14	8	15	20	6	18	12

Is there any correlation between heights and visual power?

5. Using the formula $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{xy}^2}{2\sigma_x\sigma_y}$, find r from the following data:

x	92	89	87	86	83	77	71	63	53	50
y	86	88	91	77	68	85	82	82	37	57

6. Find the correlation between x (marks in Mathematics) and y (marks in Engineering Drawing) given in the following data:

x	10-40	40-70	70-100	Total
y				
0-30	5	20	—	25
30-60	—	28	2	30
60-90	—	32	13	45
Total	5	80	15	100

7. If the coefficient of correlation between two variables x and y is 0.5 and the acute angle between their lines of regression is $\tan^{-1}(3/8)$, show that $\alpha_x = \frac{1}{2} \alpha_y$.

8. Following table gives the data on rainfall and discharge in a certain river. Obtain the line of regression of y on x .

Rainfall x (inches): 1.63 1.78 2.60 2.95 3.42,

Discharge y (1000 cc): 33.5 36.3 40.0 45.8 53.5

9. Two random variables have the regression lines with equations $3x + 2y = 26$ and $6x + y = 31$. Find the mean values and the correlation coefficient between x and y .

(Madras, 2002)

$$Y_i^2 - 2\sum X_i Y_i$$

$$\Sigma Y_i^2 - \Sigma d_i^2 = \frac{1}{12} (n^3 - n) - \frac{1}{2} \Sigma d_i^2.$$

It between these variables is

$$= \frac{\frac{1}{12} (n^3 - n) - \frac{1}{2} \Sigma d_i^2}{\frac{1}{12} (n^3 - n)} = 1 - \frac{6 \Sigma d_i^2}{n^3 - n}$$

ion coefficient and is denoted by ρ .

Judges in a contest are ranked by two judges as follows :

10	3	2	4	9	7	8
8	1	2	3	10	5	7

Coefficient ρ . (V.T.U., 2002)

then $d_i = -5, 2, -4, 2, 2, 0, 1, -1, 2, 1$

$$+ 16 + 4 + 4 + 0 + 1 + 1 + 4 + 1 = 60$$

$$\frac{\Sigma d_i^2}{n} = 1 - \frac{6 \times 60}{990} = 0.6 \text{ nearly.}$$

A, B, C, give the following ranks. Find which pair of judges

3	2	4	9	7	8,
7	10	2	1	6	9,
1	2	3	10	5	7

(J.N.T.U., 2003)

d_1	d_2	d_3	d_1^2	d_2^2	d_3^2
$x - y$	$y - z$	$z - x$			
-2	-3	5	4	9	25
1	1	-2	1	1	4
-3	-1	4	9	1	16
6	-4	-2	36	16	4
-4	6	-2	16	36	4
-8	8	0	64	64	0
2	-1	-1	4	1	1
8	-9	1	64	81	1
1	1	-2	1	1	4
-1	2	-1	1	4	1
0	0	0	200	214	60

REGRESSION ANALYSIS

Measuring the effect of advertising expenditure on sales.

$$y = \text{constant} + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots + b_n x_n$$

$$y = c + b_1 x_1 + b_2 x_2 + \dots + b_n x_n$$

Regression coefficient of x_1 is $b_1 = \frac{\partial y}{\partial x_1}$ and that of x_2 is $b_2 = \frac{\partial y}{\partial x_2}$

$$\therefore b_1 = \frac{\partial y}{\partial x_1} = \frac{\partial y}{\partial x_1} \cdot \frac{\partial x_1}{\partial x_1} = \frac{\partial y}{\partial x_1} = \frac{\partial y}{\partial x_1}$$

That is, b_1 is twice as large as the regression coefficient of x_1 .

Example 10.10. In the following table are recorded data showing the test scores made by students on an intelligence test and their weekly sales.

Test score	1	2	3	4	5	6	7	8	9	10
Test scores	40	70	30	90	50	80	60	40	60	40
Sales (y)	2.5	4.0	4.5	5.5	4.0	5.0	4.8	3.5	4.5	3.0

Calculate the regression line of sales on test scores and estimate the most probable weekly sales when a salesman makes a score of 70.

Sol. With the help of the table above, we have

$$\bar{x} = \text{mean of } x \text{ test scores} = \frac{40+70+30+90+50+80+60+40+60+40}{10} = 60$$

$$\bar{y} = \text{mean of } y \text{ sales} = \frac{2.5+4.0+4.5+5.5+4.0+5.0+4.8+3.5+4.5+3.0}{10} = 4.05$$

Regression line of sales (y) on scores (x) is given by

$$y = \bar{y} + b_1 d_x / \sigma_{d_x} + e$$

where

$$b_1 = \frac{\sum d_x d_y}{\sum d_x^2} = \frac{\sum d_x d_y - \left[\sum d_x \sum d_y \right] / \left(\sum d_x^2 - \sum d_x \sum d_y \right)}{\sum d_x^2 - \sum d_x^2 / 10} = \frac{140 - (20 \times 4.05)}{2400 - (20^2 / 10)} = \frac{140}{2400} = 0.06$$

The required regression line is

$$y = 4.05 + 0.06(x - 60) \quad \text{or} \quad y = 0.06x + 3.45$$

For $x = 70$, $y = 0.06 \times 70 + 3.45 = 4.45$.

Thus the most probable weekly sales volume for a score of 70 is 4.45.

Test score	Sales	Deviation of x from assumed mean (= 60)	Deviation of y from assumed average (= 4.05)	$d_x \times d_y$	d_x^2	d_y^2
40	2.5	-20	-2	40	400	4
70	4.0	10	1.5	15	100	2.25
50	4.5	-10	0	0	100	0
60	5.0	0	0.5	0	0	0.25
80	4.5	20	0	0	400	0
50	3.0	-10	-0.5	5	100	0.25
90	5.5	30	1	30	900	1.00
40	3.0	-20	-1.5	30	400	2.25
60	4.5	0	0	0	0	0
80	3.0	20	-1.5	-30	400	2.25
		$\sum d_x = 0$	$\sum d_y = -4.5$	$\sum d_x d_y = 140$	$\sum d_x^2 = 2400$	$\sum d_y^2 = 15.25$

or

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

i.e.

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

which is the required result.

(b) To find r , we have to calculate σ_x , σ_y and σ_{x-y} . We make the following table:

x	$X = x - 54$	X^2	y	$Y = y - 100$	Y^2	$y - x$	$(x - y)^2$
21	-33	1089	60	-40	1600	39	1521
23	-31	961	71	-29	841	48	2304
30	-24	576	72	-28	784	42	1764
54	0	0	83	-17	289	29	841
57	3	9	110	10	100	53	2809
58	4	16	84	-16	256	26	676
72	18	324	100	0	0	28	784
78	24	576	92	-8	64	14	196
87	33	1089	113	13	169	26	676
90	36	1296	135	35	1225	45	2025
Total	30	5936		-80	5328	350	13596

$$\therefore \sigma_x^2 = \frac{\sum X^2}{N} - \left(\frac{\sum X}{N} \right)^2 = \frac{5636}{10} - \left(\frac{30}{10} \right)^2 = 563.6 - 9 = 584.6$$

$$\sigma_y^2 = \frac{\sum Y^2}{N} - \left(\frac{\sum Y}{N} \right)^2 = \frac{5328}{10} - \left(\frac{-80}{10} \right)^2 = 532.8 - 64 = 468.8$$

$$\sigma_{x-y}^2 = \frac{\sum (x-y)^2}{N} - \left(\frac{\sum (x-y)}{N} \right)^2 = 1359.6 - 1225 = 134.6$$

From the above formula,

$$r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x\sigma_y} = \frac{584.6 + 468.8 - 134.6}{2 \times 24.18 \times 23.85} = 0.876$$

Example 25-20. While calculating correlation coefficient between two variables x and y from 25 pairs of observations, the following results were obtained : $n = 25$, $\sum x = 125$, $\sum x^2 = 650$, $\sum y = 100$, $\sum y^2 = 460$, $\sum xy = 508$. Later it was discovered at the time of checking that the pairs of values x | y were copied down as x | y . Obtain the correct value of correlation coefficient.

8	12	6	14
6	8	8	6

Sol. To get the correct results, we subtract the incorrect values and add the corresponding correct values.

\therefore The correct results would be

$$\Sigma n = 25, \Sigma x = 125 - 6 - 8 + 8 + 6 = 125, \Sigma x^2 = 650 - 6^2 - 8^2 + 8^2 + 6^2 = 650$$

$$\Sigma y = 100 - 14 - 6 + 12 + 8 = 100, \Sigma y^2 = 460 - 14^2 - 6^2 + 12^2 + 8^2 = 436$$

$$\Sigma xy = 508 - 6 \times 14 - 8 \times 6 + 8 \times 12 + 6 \times 8 = 520$$

Example 25-17. If θ is the angle between the two regression lines, show that

$$\tan \theta = \frac{1 - r^2}{r} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \quad (\text{V.T.U., 2007; Andhra, 2000; Ranchi, 1999})$$

Explain the significance when $r = 0$ and $r = \pm 1$. (V.T.U., 2001)

Sol. The equations to the line of regression of y on x and x on y are

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

\therefore Their slopes are $m_1 = r \sigma_y / \sigma_x$ and $m_2 = \sigma_x / r \sigma_y$

$$\text{Thus } \tan \theta = \frac{m_2 - m_1}{1 + m_1 m_2} = \frac{\sigma_x / r \sigma_y - r \sigma_y / \sigma_x}{1 + \sigma_y^2 / \sigma_x^2} = \frac{1 - r^2}{r} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$$

When $r = 0$, $\tan \theta \rightarrow \infty$ or $\theta = \pi/2$ i.e. when the variables are independent, the two lines of regression are perpendicular to each other.

When $r = \pm 1$, $\tan \theta = 0$ i.e. $\theta = 0$ or π . Thus the lines of regression coincide i.e. there is perfect correlation between the two variables.

Example 25-18. In a partially destroyed laboratory record, only the lines of regression of y on x and x on y are available as $4x - 5y + 33 = 0$ and $20x - 9y = 107$ respectively. Calculate \bar{x}, \bar{y} and the coefficient of correlation between x and y . (V.T.U., 2005; Anna, 2003 S; U.P.T.U., 2003)

Sol. Since the regression lines pass through (\bar{x}, \bar{y}) , therefore,

$$4\bar{x} - 5\bar{y} + 33 = 0, \quad 20\bar{x} - 9\bar{y} = 107.$$

Solving these equations, we get $\bar{x} = 13, \bar{y} = 17$.

Rewriting the line of regression of y on x as $y = \frac{4}{5}x + \frac{33}{5}$, we get

$$b_{yx} = r \frac{\sigma_y}{\sigma_x} = \frac{4}{5} \quad \dots(i)$$

Rewriting the line of regression of x on y as $x = \frac{9}{20}y + \frac{107}{9}$, we get

$$b_{xy} = r \frac{\sigma_x}{\sigma_y} = \frac{9}{20} \quad \dots(ii)$$

Multiplying (i) and (ii), we get

$$r^2 = \frac{4}{5} \times \frac{9}{20} = 0.36 \quad \therefore r = 0.6$$

Hence $r = 0.6$, the positive sign being taken as b_{yx} and b_{xy} both are positive.

Example 25-19. Establish the formula $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2 \sigma_x \sigma_y}$

Hence calculate r from the following data :

$x:$	21	23	30	54	57	58	72	78	87	90
$y:$	60	71	72	83	110	84	100	92	113	135

(U.P.T.U. 2002)

(a) Let $z = x - y$ so that $\bar{z} = \bar{x} - \bar{y}$.

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

$$(z - \bar{z})^2 = (x - \bar{x})^2 + (y - \bar{y})^2 - 2(x - \bar{x})(y - \bar{y})$$

Summing up for n terms, we have

$$\sum (z - \bar{z})^2 = \sum (x - \bar{x})^2 + \sum (y - \bar{y})^2 - 2 \sum (x - \bar{x})(y - \bar{y})$$

$$r = \frac{n \sum xy - (\bar{x}) (\bar{y})}{\sqrt{(n \sum x^2 - (\bar{x})^2)(n \sum y^2 - (\bar{y})^2)}} = \frac{25 \times 520 - 125 \times 100}{\sqrt{[(25 \times 850) - (125)^2][25 \times 436 - (100)^2]}} \\ = \frac{20}{\sqrt{25 \times 36}} = \frac{2}{3}$$

25.15. STANDARD ERROR OF ESTIMATE

The sum of the squares of the deviations of the points from the line of regression of y on x is $\sum(y - a - bx)^2 = \sum(Y - bX)^2$, where $X = x - \bar{x}$, $Y = y - \bar{y}$

$$\begin{aligned} &= \sum \left(Y - r \frac{\sigma_y}{\sigma_x} X \right)^2 = \sum Y^2 - 2r(\sigma_y/\sigma_x) \sum XY + r^2(\sigma_y^2/\sigma_x^2) \sum X^2 \\ &= n\sigma_y^2 - 2r(\sigma_y/\sigma_x) r \cdot n\sigma_x\sigma_y + r^2(\sigma_y^2/\sigma_x^2) \cdot n\sigma_x^2 = n\sigma_y^2(1 - r^2). \end{aligned}$$

Denoting this sum of squares by nS_y^2 , we have $S_y = \sigma_y \sqrt{1 - r^2}$... (1)

Since S_y is the root mean square deviation of the points from the regression line of y on x , it is called the standard error of estimate of y . Similarly the standard error of estimate of x is given by

$$S_x = \sigma_x \sqrt{1 - r^2}$$

Since the sum of the squares of deviations cannot be negative, it follows that

$$r^2 \leq 1 \quad \text{or} \quad -1 \leq r \leq 1.$$

i.e. correlation coefficient lies between -1 and 1.

(J.N.T.U., 2006)

If $r = 1$ or -1 , the sum of the squares of deviations from either line of regression is zero. Consequently each deviation is zero and all the points lie on both the lines of regression. These two lines coincide and we say that the correlation between the variables is perfect. The nearer r^2 is to unity the closer are the points to the lines of regression. Thus the departure of r^2 from unity is a measure of departure from linearity of the relationship between the variables.

25.16. RANK CORRELATION

A group of n individuals may be arranged in order of merit with respect to some characteristic. The same group would give different orders for different characteristics. Considering the orders corresponding to two characteristics A and B , the correlation between these n pairs of ranks is called the rank correlation in the characteristics A and B for that group of individuals.

Let x_i, y_i be the ranks of the i th individuals in A and B respectively. Assuming that no two individuals are bracketed equal in either case, each of the variables taking the values 1, 2, 3, ..., n , we have

$$\bar{x} = \bar{y} = \frac{1 + 2 + 3 + \dots + n}{n} = \frac{n(n+1)}{2n} = \frac{n+1}{2}$$

If X, Y be the deviations of x, y from their means, then

$$\begin{aligned} \sum X_i^2 &= \sum (x_i - \bar{x})^2 = \sum x_i^2 + n(\bar{x})^2 - 2\bar{x}\sum x_i = \sum x_i^2 + \frac{n(n+1)^2}{4} - 2 \frac{n+1}{2} \cdot \sum x_i \\ &= \frac{n(n+1)(2n+1)}{6} + \frac{n(n+1)^2}{4} - \frac{n(n+1)^2}{2} = \frac{1}{12}(n^3 - n) \end{aligned}$$

Similarly $\sum Y_i^2 = \frac{1}{12}(n^3 - n)$

Now let $d_i = x_i - y_i$ so that $d_i = (x_i - \bar{x}) - (y_i - \bar{y}) = X_i - Y_i$

With the help of the above correlation table, we have

$$r = \frac{n(\Sigma f d_x d_y) - (\Sigma f d_x)(\Sigma f d_y)}{\sqrt{[(n \Sigma f d_x^2 - (\Sigma f d_x)^2) \times (n \Sigma f d_y^2 - (\Sigma f d_y)^2)]}}$$

$$= \frac{53 \times 86 - 10 \times 16}{\sqrt{[(53 \times 98 - 100) \times (53 \times 92 - 256)]}} = \frac{4398}{\sqrt{(5094 \times 4620)}} = \frac{4398}{4850} = 0.91 \text{ (approx.)}$$

25.14. LINES OF REGRESSION

It frequently happens that the dots of the scatter diagram generally, tend to cluster along a well defined direction which suggests a linear relationship between the variables x and y . Such a line of best-fit for the given distribution of dots is called the *line of regression* (Fig. 25-6). In fact there are two such lines, one giving the best possible mean values of y for each specified value of x and the other giving the best possible mean values of x for given values of y . The former is known as the *line of regression of y on x* and the latter as the *line of regression of x on y* .

Consider first the line of regression of y on x . Let the straight line satisfying the general trend of n dots in a scatter diagram be

$$y = a + bx \quad \dots(1)$$

We have to determine the constants a and b so that (1) gives for each value of x , the best estimate for the average value of y in accordance with the *principle of least squares* (page 20), therefore, the normal equations for a and b are

$$\Sigma y = na + b \Sigma x \quad \dots(2)$$

and

$$\Sigma xy = a \Sigma x + b \Sigma x^2 \quad \dots(3)$$

(2) gives

$$\frac{1}{n} \Sigma y = a + b \cdot \frac{1}{n} \Sigma x \quad i.e. \quad \bar{y} = a + b \bar{x}.$$

This shows that (\bar{x}, \bar{y}) , i.e. the means of x and y , lie on (1).

Shifting the origin to (\bar{x}, \bar{y}) , (3) takes the form

$$\Sigma(x - \bar{x})(y - \bar{y}) = a \Sigma(x - \bar{x}) + b \Sigma(x - \bar{x})^2, \quad \text{But } a \Sigma(x - \bar{x}) = 0,$$

$$\therefore b = \frac{\Sigma(x - \bar{x})(y - \bar{y})}{\Sigma(x - \bar{x})^2} = \frac{\Sigma XY}{\Sigma X^2} = \frac{\Sigma XY}{n \sigma_x^2} = r \frac{\sigma_y}{\sigma_x} \quad \left[\therefore r = \frac{\Sigma XY}{n \sigma_x \sigma_y} \right]$$

Thus the line of best fit becomes $y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \quad \dots(4)$

which is the equation of the *line of regression of y on x* . Its slope is called the *regression co-efficient of y on x* .

Interchanging x and y , we find that the line of regression of x on y is

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y}) \quad \dots(5)$$

Thus the *regression co-efficient of y on x* = $r \sigma_y / \sigma_x$ $\dots(6)$

and the *regression co-efficient of x on y* = $r \sigma_x / \sigma_y$ $\dots(7)$

Cor. The correlation co-efficient r is the geometric mean between the two regression co-efficients.

$$\text{For } r \frac{\sigma_y}{\sigma_x} \times r \frac{\sigma_x}{\sigma_y} = r^2.$$

Example 25.15. The two regression equations of the variables x and y are $x = 19.13 - 0.87y$ and $y = 11.64 - 0.50x$. Find (i) mean of x 's, (ii) mean of y 's and (iii) the correlation co-efficient between x and y . (V.T.U., 2004; Anna, 2003; Burdwan, 2003)

Sol. Since the mean of x 's and the mean of y 's lie on the two regression lines, we have

$$\bar{x} = 19.13 - 0.87 \bar{y} \quad \dots(i)$$

$$\bar{y} = 11.64 - 0.50 \bar{x} \quad \dots(ii)$$

With the help of the above correlation table, we have

$$\begin{aligned} r &= \frac{n(\Sigma M_x M_y) - (\Sigma M_x)(\Sigma M_y)}{\sqrt{[(n\Sigma M_x^2) - (\Sigma M_x)^2] \cdot [(n\Sigma M_y^2) - (\Sigma M_y)^2]}} \\ &= \frac{53 \times 96 - 10 \times 16}{\sqrt{(53 \times 96 - 100) \times (53 \times 92 - 256)}} = \frac{4398}{\sqrt{5094 \times 4620}} = \frac{4398}{4850} = 0.91 \text{ (approx.)} \end{aligned}$$

25.14. LINES OF REGRESSION

It frequently happens that the dots of the scatter diagram generally, tend to cluster along a well defined direction which suggests a linear relationship between the variables x and y . Such a line of best-fit for the given distribution of dots is called the *line of regression* (Fig. 25-8). In fact there are two such lines, one giving the best possible mean values of y for each specified value of x and the other giving the best possible mean values of x for given values of y . The former is known as the *line of regression of y on x* and the latter as the *line of regression of x on y* .

Consider first the line of regression of y on x . Let the straight line satisfying the general trend of n dots in a scatter diagram be

$$y = a + bx \quad \dots(1)$$

We have to determine the constants a and b so that (1) gives for each value of x , the best estimate for the average value of y in accordance with the principle of least squares (page 20), therefore, the normal equations for a and b are

$$\Sigma y = na + b \Sigma x \quad \dots(2)$$

and

$$\Sigma xy = a \Sigma x + b \Sigma x^2 \quad \dots(3)$$

$$(2) \text{ gives } \frac{1}{n} \Sigma y = a + b \cdot \frac{1}{n} \Sigma x \text{ i.e. } \bar{y} = a + b \bar{x}.$$

This shows that (\bar{x}, \bar{y}) , i.e. the means of x and y , lie on (1).

Shifting the origin to (\bar{x}, \bar{y}) , (3) takes the form

$$\Sigma(x - \bar{x})(y - \bar{y}) = a \Sigma(x - \bar{x}) + b \Sigma(x - \bar{x})^2, \text{ But } a \Sigma(x - \bar{x}) = 0,$$

$$\therefore b = \frac{\Sigma(x - \bar{x})(y - \bar{y})}{\Sigma(x - \bar{x})^2} = \frac{\Sigma XY}{\Sigma X^2} = \frac{\Sigma XY}{n \sigma_x^2} = r \frac{\sigma_y}{\sigma_x} \quad \left[\because r = \frac{\Sigma XY}{n \sigma_x \sigma_y} \right]$$

$$\text{Thus the line of best fit becomes } y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \quad \dots(4)$$

which is the equation of the line of regression of y on x . Its slope is called the *regression co-efficient of y on x* .

Interchanging x and y , we find that the line of regression of x on y is

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y}) \quad \dots(5)$$

$$\text{Thus the regression co-efficient of } y \text{ on } x = r \sigma_y / \sigma_x \quad \dots(6)$$

$$\text{and the regression co-efficient of } x \text{ on } y = r \sigma_x / \sigma_y \quad \dots(7)$$

Cor. The correlation co-efficient r is the geometric mean between the two regression co-efficients.

$$\text{For } r \frac{\sigma_x}{\sigma_y} \times r \frac{\sigma_y}{\sigma_x} = r^2.$$

Example 25.15. The two regression equations of the variables x and y are $= 19.13 - 0.87y$ and $y = 11.64 - 0.50x$. Find (i) mean of x 's, (ii) mean of y 's and (iii) the correlation co-efficient between x and y . (V.T.U., 2004; Anna, 2003; Burdwan, 2003)

Sol. Since the mean of x 's and the mean of y 's lie on the two regression lines, we have

$$\bar{x} = 19.13 - 0.87 \bar{y}$$

$$\bar{y} = 11.64 - 0.50 \bar{x}$$

With the help of the above correlation ratio, we have

$$\begin{aligned} r &= \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \cdot \frac{\sigma_{xy}}{\sigma_x \sigma_y} \\ &= \frac{\sqrt{\sigma_{xy}^2}}{\sqrt{\sigma_x^2} \cdot \sqrt{\sigma_y^2}} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} = 0.91 \text{ (approx.)} \end{aligned}$$

25.14. LINES OF REGRESSION

It frequently happens that the data of the scatter diagram generally, tend to cluster along a well defined direction which suggests a linear relationship between the variables x and y . Then a line of best fit for the given distribution of data is called the *line of regression* (Fig. 25.6). In fact there are two such lines, one giving the best possible mean value of y for each specified value of x and the other giving the best possible mean value of x for given value of y . The former is known as the *line of regression of y on x* and the latter as the *line of regression of x on y* .

Consider first the line of regression of y on x . Let the straight line satisfying the general trend of n data in a scatter diagram be

$$y = a + bx \quad (1)$$

We have to determine the constants a and b so that (1) gives for each value of x , the best estimate for the average value of y in accordance with the *principle of least squares* (page 20), therefore, the normal equations for a and b are

$$\sum y = na + b\sum x \quad (2)$$

and

$$\sum xy = a\sum x + b\sum x^2 \quad (3)$$

$$(2) \text{ gives } \frac{1}{n} \sum y = a + b \cdot \frac{1}{n} \sum x, \text{ i.e., } \bar{y} = a + b\bar{x}.$$

This shows that (\bar{x}, \bar{y}) , i.e., the means of x and y , lie on (1).

Shifting the origin to (\bar{x}, \bar{y}) , (3) takes the form

$$\begin{aligned} \sum(x - \bar{x})(y - \bar{y}) &= a\sum(x - \bar{x}) + b\sum(x - \bar{x})^2, \text{ But } a\sum(x - \bar{x}) = 0, \\ b &= \frac{\sum(x - \bar{x})(y - \bar{y})}{\sum(x - \bar{x})^2} = \frac{\sum XY - \bar{X}\bar{Y}}{\sum X^2 - n\bar{x}^2} = r \frac{\sigma_y}{\sigma_x} \quad \left[\therefore r = \frac{\sum XY}{n\sigma_x \sigma_y} \right] \end{aligned}$$

$$\text{Thus the line of best fit becomes } y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x}) \quad (4)$$

which is the equation of the *line of regression of y on x* . Its slope is called the *regression co-efficient of y on x* .

Interchanging x and y , we find that the line of regression of x on y is

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y}) \quad (5)$$

Thus the *regression co-efficient of x on y* is $= r\sigma_x/\sigma_y$

$$\text{and} \quad \text{the regression co-efficient of x on y } = r\sigma_x/\sigma_y \quad (6)$$

Cor. The correlation co-efficient r is the geometric mean between the two regression co-efficients.

$$\text{For } r \frac{\sigma_x}{\sigma_y} \times r \frac{\sigma_x}{\sigma_y} = r^2.$$

Example 25.15. The two regression equations of the variables x and y are $= 19.13 - 0.87y$ and $y = 11.64 - 0.50x$. Find (i) mean of x 's, (ii) mean of y 's and (iii) the correlation co-efficient between x and y . (V.T.U., 2004; Anna, 2003; Burdwan, 2003)

Sol. Since the mean of x 's and the mean of y 's lie on the two regression lines, we have

$$\bar{x} = 19.13 - 0.87 \bar{y} \quad (i)$$

$$\bar{y} = 11.64 - 0.50 \bar{x} \quad (ii)$$

