

**CHAPTER TEN**

# **Correlation and Regression**

---

**10-1. Bivariate Distribution, Correlation.** So far we have confined ourselves to univariate distributions, i.e., the distributions involving only one variable. We may, however, come across certain series where each term of the series may assume the values of two or more variables. For example, if we measure the heights and weights of a certain group of persons, we shall get what is known as *Bivariate distribution*—one variable relating to height and other variable relating to weight.

In a bivariate distribution we may be interested to find out if there is any correlation or covariation between the two variables under study. If the change in one variable affects a change in the other variable, the variables are said to be correlated. If the two variables deviate in the same direction, i.e., if the increase (or decrease) in one results in a corresponding increase (or decrease) in the other, correlation is said to be *direct* or *positive*. But if they constantly deviate in the opposite directions, i.e., if increase (or decrease) in one results in corresponding decrease (or increase) in the other, correlation is said to be *diverse* or *negative*. For example, the correlation between (i) the heights and weights of a group of persons, (ii) the income and expenditure is positive and the correlation between (i) price and demand of a commodity, (ii) the volume and pressure of a perfect gas, is negative. Correlation is said to be *perfect* if the deviation in one variable is followed by a corresponding and proportional deviation in the other.

**10-2. Scatter Diagram.** It is the simplest way of the diagrammatic representation of bivariate data. Thus for the bivariate distribution  $(x_i, y_i); i = 1, 2, \dots, n$ , if the values of the variables  $X$  and  $Y$  be plotted along the  $x$ -axis and  $y$ -axis respectively in the  $xy$  plane, the diagram of dots so obtained is known as *scatter diagram*. From the scatter diagram, we can form a fairly good, though vague, idea whether the variables are correlated or not, e.g., if the points are very dense, i.e., very close to each other, we should expect a fairly good amount of correlation between the variables and if the points are widely scattered, a poor correlation is expected. This method, however, is not suitable if the number of observations is fairly large.

**10-3. Karl Pearson Coefficient of Correlation.** As a measure of intensity or degree of linear relationship between two variables, Karl Pearson (1867—1936), a British Biometrician, developed a formula called *Correlation Coefficient*.

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

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \quad (10-1)$$

If  $(x_i, y_i) ; i = 1, 2, \dots, n$  is the bivariate distribution, then

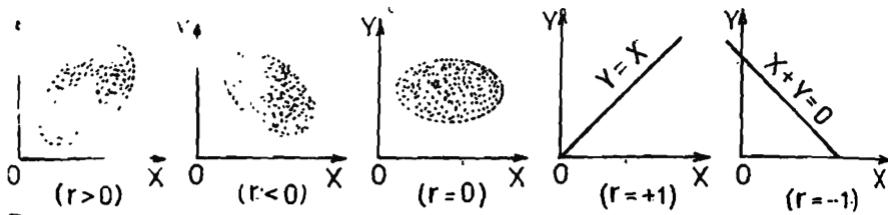
$$\left. \begin{aligned} \text{Cov}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) = \mu_{11} \\ \sigma_X^2 &= E(X - E(X))^2 = \frac{1}{n} \sum (x_i - \bar{x})^2 \\ \sigma_Y^2 &= E(Y - E(Y))^2 = \frac{1}{n} \sum (y_i - \bar{y})^2 \end{aligned} \right\}, \quad \dots (10.2)$$

the summation extending over  $i$  from 1 to  $n$ .

Another convenient form of the formula (10.2) for computational work is as follows :

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} (x_1 y_1 - x_1 \bar{y} - \bar{x} y_1 + \bar{x} \bar{y}) \\ &= \frac{1}{n} \sum x_i y_i - \bar{y} \frac{1}{n} \sum x_i - \bar{x} \frac{1}{n} \sum y_i + \bar{x} \bar{y} \\ \therefore \text{Cov}(X, Y) &= \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}, \sigma_X^2 = \frac{1}{n} \sum x_i^2 - \bar{x}^2 \\ \text{and } \sigma_Y^2 &= \frac{1}{n} \sum y_i^2 - \bar{y}^2 \end{aligned} \quad \dots (10.2a)$$

**Remarks 1.** Following are the figures of the standard data for  $r > 0$ ,  $< 0$ ,  $= 0$ , and  $r = \pm 1$ .



2. It may be noted that  $r(X, Y)$  provides a measure of *linear relationship* between  $X$  and  $Y$ . For nonlinear relationship, however, it is not very suitable.

3. Sometimes, we write :  $\text{Cov}(X, Y) = \sigma_{XY}$

4. Karl Pearson's correlation coefficient is also called *product-moment correlation coefficient*, since

$$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))] = \mu_{11}.$$

**10.3.1. Limits for Correlation Coefficient.** We have

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

$$\therefore r^2(X, Y) = \frac{(\sum a_i b_i)^2}{(\sum a_i^2)(\sum b_i^2)}, \text{ where } \begin{cases} a_i = x_i - \bar{x} \\ b_i = y_i - \bar{y} \end{cases} \quad \dots(*)$$

We have the Schwartz inequality which states that if  $a_i, b_i; i = 1, 2, \dots, n$  are real quantities then

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

the sign of equality holding if and only if

$$\frac{a_1}{b_1} = \frac{a_2}{b_2} = \dots = \frac{a_n}{b_n}$$

Using Schwartz inequality, we get from (\*)

$$r^2(X, Y) \leq 1 \text{ i.e., } |r(X, Y)| \leq 1 \Rightarrow -1 \leq r(X, Y) \leq 1 \quad \dots(10-3)$$

Hence correlation coefficient cannot exceed unity numerically. It always lies between -1 and +1. If  $r = +1$ , the correlation is perfect and positive and if  $r = -1$ , correlation is perfect and negative.

Aliter. If we write  $E(X) = \mu_X$  and  $E(Y) = \mu_Y$ , then we have

$$\begin{aligned} & E \left[ \left( \frac{X - \mu_X}{\sigma_X} \right) \pm \left( \frac{Y - \mu_Y}{\sigma_Y} \right) \right]^2 \geq 0 \\ \Rightarrow & E \left( \frac{X - \mu_X}{\sigma_X} \right)^2 + E \left( \frac{Y - \mu_Y}{\sigma_Y} \right)^2 \pm 2 \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} \geq 0 \\ \Rightarrow & 1 + 1 \pm 2r(X, Y) \geq 0 \\ \Rightarrow & -1 \leq r(X, Y) \leq 1. \end{aligned}$$

**Theorem 10-1.** Correlation coefficient is independent of change of origin and scale.

**Proof.** Let  $U = \frac{X - a}{h}, V = \frac{Y - b}{k}$ , so that  $X = a + hU$  and  $Y = b + kV$ , where  $a, b, h, k$  are constants;  $h > 0, k > 0$ .

We shall prove that  $r(X, Y) = r(U, V)$

Since  $X = a + hU$  and  $Y = b + kV$ , on taking expectations, we get

$$\begin{aligned} E(X) &= a + hE(U) \quad \text{and} \quad E(Y) = b + kE(V) \\ \therefore X - E(X) &= h[U - E(U)] \quad \text{and} \quad Y - E(Y) = k[V - E(V)] \\ \Rightarrow \text{Cov}(X, Y) &= E[(X - E(X))(Y - E(Y))] \\ &= E[h(U - E(U))(k(V - E(V)))] \\ &= hk E[(U - E(U))(V - E(V))] = hk \text{ Cov}(U, V) \quad \dots(10-4) \\ \sigma_X^2 &= E[(X - E(X))^2] = E[h^2(U - E(U))^2] = h^2 \sigma_U^2 \\ \Rightarrow \sigma_X &= h \sigma_U, (h > 0) \quad \dots(10-4a) \\ \text{and} \quad \sigma_Y^2 &= E[(Y - E(Y))^2] = E[k^2(V - E(V))^2] = k^2 \sigma_V^2 \\ \Rightarrow \sigma_Y &= k \sigma_V, (k > 0) \quad \dots(10-4b) \end{aligned}$$

Substituting from (10-4), (10-4a) and (10-4b) in (10-1), we get

$$r(X, Y) = \frac{\text{Cov}(\bar{X}, Y)}{\sigma_X \sigma_Y} = \frac{hk \cdot \text{Cov}(U, V)}{hk \cdot \sigma_U \sigma_V} = \frac{\text{Cov}(U, V)}{\sigma_U \sigma_V} = r(U, V)$$

This theorem is of fundamental importance in the numerical computation of the correlation coefficient.

**Corollary.** 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

$$r(aX + b, cY + d) = \frac{ac}{|ac|} r(X, Y)$$

**Proof.** With usual notations, we have

$$\text{Var}(aX + b) = a^2 \sigma_X^2; \quad \text{Var}(cY + d) = c^2 \sigma_Y^2;$$

$$\text{Cov}(aX + b, cY + d) = ac \sigma_{XY}$$

$$\therefore r(aX + b, cY + d) = \frac{\text{Cov}(aX + b, cY + d)}{[\text{Var}(aX + b) \text{Var}(cY + d)]^{1/2}} \\ = \frac{ac \sigma_{XY}}{|a| |c| \sigma_X \sigma_Y} = \frac{ac}{|ac|} r(X, Y)$$

If  $ac > 0$ , i.e., if  $a$  and  $c$  are of same signs, then  $ac/|ac| = +1$

If  $ac < 0$ , i.e., if  $a$  and  $c$  are of opposite signs, then  $ac/|ac| = -1$ .

**Theorem 10-2.** Two independent variables are uncorrelated.

**Proof.** If  $X$  and  $Y$  are independent variables, then

$$\text{Cov}(X, Y) = 0 \quad (\text{c.f. } \S\ 6-4)$$

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = 0$$

Hence two independent variables are uncorrelated.

But the converse of the theorem is not true, i.e., two uncorrelated variables may not be independent as the following example illustrates :

$X$	-3	-2	-1	1	2	3	Total $\sum X = 0$
$Y$	9	4	1	1	4	9	$\sum Y = 28$
$XY$	-27	-8	-1	1	8	27	$\sum XY = 0$

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

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = 0$$

Thus in the above example, the variables  $X$  and  $Y$  are uncorrelated. But on careful examination we find that  $X$  and  $Y$  are not independent but they are connected by the relation  $Y = X^2$ . Hence two uncorrelated variables need not necessarily be independent. A simple reasoning for this strange conclusion is that  $r(X, Y) = 0$ , merely implies the absence of any linear relationship between

the variables  $X$  and  $Y$ . There may, however, exist some other form of relationship between them, e.g., quadratic, cubic or trigonometric.

**Remarks.** 1. Following are some more examples where two variables are uncorrelated but not independent.

$$(i) X \sim N(0, 1) \text{ and } Y = X^2$$

$$\text{Since } X \sim N(0, 1), E(X) = 0 = E(X^3)$$

$$\therefore \text{Cov}(X, Y) = E(XY) - E(X)E(Y) \\ = E(X^3) - E(X)E(Y) = 0 \quad (\because Y = X^2)$$

$$\Rightarrow r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = 0$$

Hence  $X$  and  $Y$  are uncorrelated but not independent.

(ii) Let  $X$  be a r.v. with p.d.f.

$$f(x) = \frac{1}{2}, -1 \leq x \leq 1$$

and let  $Y = X^2$ . Here we shall get

$$E(X) = 0 \text{ and } E(XY) = E(X^3) = 0, \Rightarrow r(X, Y) = 0$$

2. However, the converse of the theorem holds in the following cases :

(a) If  $X$  and  $Y$  are jointly normally distributed with  $\rho = \rho(X, Y) = 0$ , then they are independent. If  $\rho = 0$ , then [c.f. § 10.10, Equation (10.25)]

$$f(x, y) = \frac{1}{\sigma_X \sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{X - \mu_X}{\sigma_X}\right)^2\right] \times \frac{1}{\sigma_Y \sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{Y - \mu_Y}{\sigma_Y}\right)^2\right]$$

$$\therefore f(x, y) = f_1(x)f_2(y)$$

$$\Rightarrow X \text{ and } Y \text{ are independent.}$$

(b) If each of the two variables  $X$  and  $Y$  takes two values, 0, 1 with positive probabilities, then  $r(X, Y) = 0 \Rightarrow X \text{ and } Y \text{ are independent.}$

**Proof.** Let  $X$  take the values 1 and 0 with positive probabilities  $p_1$  and  $q_1$  respectively and let  $Y$  take the values 1 and 0 with positive probabilities  $p_2$  and  $q_2$  respectively. Then

$$\begin{aligned} r(X, Y) &= 0 \Rightarrow \text{Cov}(X, Y) = 0 \\ \Rightarrow 0 &= E(XY) - E(X)E(Y) \\ &= 1 \cdot P(X = 1 \cap Y = 1) - [1 \cdot P(X = 1) \times 1 \cdot P(Y = 1)] \\ &= P(X = 1 \cap Y = 1) - p_1 p_2 \\ \Rightarrow P(X = 1 \cap Y = 1) &= p_1 p_2 = P(X = 1) \cdot P(Y = 1) \\ \Rightarrow X \text{ and } Y &\text{ are independent.} \end{aligned}$$

**10.3.2. Assumptions Underlying Karl Pearson's Correlation Coefficient.** Pearsonian correlation coefficient  $r$  is based on the following assumptions :

(i) The variables  $X$  and  $Y$  under study are linearly related. In other words, the scatter diagram of the data will give a straight line curve.

(ii) Each of the variables (series) is being affected by a large number of independent contributory causes of such a nature as to produce normal distribution. For example, the variables (series) relating to ages, heights, weight, supply, price, etc., conform to this assumption. In the words of Karl Pearson :

"The sizes of the complex of organs (something measurable) are determined by a great variety of independent contributory causes, for example, climate, nourishment, physical training and innumerable other causes which cannot be individually observed or their effects measured." Karl Pearson further observes, "The variations in intensity of the contributory causes are small as compared with their absolute intensity and these variations follow the normal law of distribution."

(iii) The forces so operating on each of the variable series are not independent of each other but are related in a causal fashion. In other word, cause and effect relationship exists between different forces operating on the items of the two variable series. These forces must be common to both the series. If the operating forces are entirely independent of each other and not related in any fashion, then there cannot be any correlation between the variables under study.

For example, the correlation coefficient between,

- (a) the series of heights and incomes of individuals over a period of time,
- (b) the series of marriage rate and the rate of agricultural production in a country over a period of time,

(c) the series relating to the size of the shoe and intelligence of a group of individuals,

should be zero, since the forces affecting the two variable series in each of the above cases are entirely independent of each other.

However, if in any of the above cases the value of  $r$  for a given set of data is not zero, then such correlation is termed as *chance correlation* or *spurious* or *nonsense correlation*.

**Example 10.1.** 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

**Solution.**

#### CALCULATIONS FOR CORRELATION COEFFICIENT

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
Total	544	37028	38132	37560

Correlation and Regression

$$\bar{X} = \frac{1}{n} \sum X = \frac{544}{8} = 68, \bar{Y} = \frac{1}{n} \sum Y = \frac{1}{8} \times 552 = 69$$

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\frac{1}{n} \sum XY - \bar{X} \bar{Y}}{\sqrt{\left(\frac{1}{n} \sum X^2 - \bar{X}^2\right) \left(\frac{1}{n} \sum Y^2 - \bar{Y}^2\right)}}$$

$$= \frac{\frac{1}{8} \times 37560 - 68 \times 69}{\sqrt{\left[\frac{37028}{8} - (68)^2\right] \left[\frac{38132}{8} - (69)^2\right]}}$$

$$= \frac{4695 - 4692}{\sqrt{(4628.5 - 4624)(4766.5 - 4761)}} = \frac{3}{\sqrt{4.5 \times 5.5}} = 0.603$$

**Aliter.**

(SHORT-CUT METHOD)

X	Y	$U = X - 68$	$V = Y - 69$	$U^2$	$V^2$	$UV$
65	67	-3	-2	9	4	6
66	68	-2	-1	4	1	2
67	65	-1	-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
Total		0	0	36	44	24

$$\bar{U} = \frac{1}{n} \sum U = 0, \bar{V} = \frac{1}{n} \sum V = 0$$

$$\text{Cov}(U, V) = \frac{1}{n} \sum UV - \bar{U} \bar{V} = \frac{1}{8} \times 24 = 3$$

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

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

$$\therefore r(U, V) = \frac{\text{Cov}(U, V)}{\sigma_U \sigma_V} = \frac{3}{\sqrt{4.5 \times 5.5}} = 0.603 = r(X, Y)$$

**Remark.** The reader is advised to calculate the correlation coefficient by arbitrary origin method rather than by the direct method; since the latter leads to much simpler arithmetical calculations.

**Example 10-2.** 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

X	Y
8	12
6	8

Obtain the correct value of correlation coefficient.

[Calcutta Univ. B.Sc. (Maths. Hons.), 1988, 1991]

**Solution.**

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

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

$$\text{Corrected } \sum X^2 = 650 - 6^2 - 8^2 + 8^2 + 6^2 = 650$$

$$\text{Corrected } \sum Y^2 = 460 - 14^2 - 6^2 + 12^2 + 8^2 = 436$$

$$\text{Corrected } \sum XY = 508 - 6 \times 14 - 8 \times 6 + 8 \times 12 + 6 \times 8 = 520$$

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

$$\text{Cov}(X, Y) = \frac{1}{n} \sum XY - \bar{X}\bar{Y} = \frac{1}{25} \times 520 - 5 \times 4 = \frac{4}{5}$$

$$\sigma_X^2 = \frac{1}{n} \sum X^2 - \bar{X}^2 = \frac{1}{25} \times 650 - (5)^2 = 1$$

$$\sigma_Y^2 = \frac{1}{n} \sum Y^2 - \bar{Y}^2 = \frac{1}{25} \times 436 - 16 = \frac{36}{25}$$

$$\therefore \text{Corrected } r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\frac{4}{5}}{1 \times \frac{6}{5}} = \frac{2}{3} = 0.67$$

**Example 10-3.** Show that if  $X'$ ,  $Y'$  are the deviations of the random variables  $X$  and  $Y$  from their respective means then

$$(i) \quad r = 1 - \frac{1}{2N} \sum_i \left( \frac{X'_i}{\sigma_X} - \frac{Y'_i}{\sigma_Y} \right)^2$$

$$(ii) \quad r = -1 + \frac{1}{2N} \sum_i \left( \frac{X'_i}{\sigma_X} + \frac{Y'_i}{\sigma_Y} \right)^2$$

Deduce that  $-1 \leq r \leq +1$ .

[Delhi Univ. B.Sc. Oct. 1992; Madras Univ. B.Sc., Nov. 1991]

**Solution.** (i) Here  $X'_i = (x_i - \bar{X})$  and  $Y'_i = (Y_i - \bar{Y})$

$$\text{R.H.S.} = 1 - \frac{1}{2N} \sum_i \left( \frac{X'_i}{\sigma_X} - \frac{Y'_i}{\sigma_Y} \right)^2$$

$$\begin{aligned}
 &= 1 - \frac{1}{2N} \sum_i \left[ \frac{X'_i}{\sigma_X^2} + \frac{Y'_i}{\sigma_Y^2} - \frac{2X'_i Y'_i}{\sigma_X \sigma_Y} \right] \\
 &= 1 - \frac{1}{2N} \left[ \frac{1}{\sigma_X^2} \sum_i X'^2 + \frac{1}{\sigma_Y^2} \sum_i Y'^2 - \frac{2}{\sigma_X \sigma_Y} \sum_i X'_i Y'_i \right] \\
 &= 1 - \frac{1}{2N} \left[ \frac{1}{\sigma_X^2} \sum_i (X_i - \bar{X})^2 + \frac{1}{\sigma_Y^2} \sum_i (Y_i - \bar{Y})^2 - \frac{2}{\sigma_X \sigma_Y} \sum_i (X'_i - \bar{X})(Y'_i - \bar{Y}) \right] \\
 &= 1 - \frac{1}{2} \left[ \frac{1}{\sigma_X^2} \cdot \sigma_X^2 + \frac{1}{\sigma_Y^2} \cdot \sigma_Y^2 - \frac{2}{\sigma_X \sigma_Y} \cdot r \sigma_X \sigma_Y \right] \\
 &= 1 - \frac{1}{2} [1 + 1 - 2r] = r
 \end{aligned}$$

(ii) Proceeding similarly, we will get

$$\text{R.H.S.} = -1 + \frac{1}{2}(1 + 1 + 2r) = r$$

**Deduction.** Since  $\left(\frac{X'_i}{\sigma_X} \pm \frac{Y'_i}{\sigma_Y}\right)^2$ , being the square of a real quantity is always non-negative,  $\sum_i \left(\frac{X'_i}{\sigma_X} \mp \frac{Y'_i}{\sigma_Y}\right)^2$  is also non-negative. From part (i), we get

$$r = 1 - (\text{some non-negative quantity}) \Rightarrow r \leq 1 \quad \dots(*)$$

Also from part (ii), we get

$$r = -1 + (\text{some non-negative quantity}) \Rightarrow -1 \leq r \quad \dots(**)$$

The sign of equality in (\*) and (\*\*) holds if and only if

$$\left. \begin{array}{l} \frac{X'_i}{\sigma_X} - \frac{Y'_i}{\sigma_Y} = 0 \\ \frac{X'_i}{\sigma_X} + \frac{Y'_i}{\sigma_Y} = 0 \end{array} \right\} \forall i = 1, 2, \dots, n$$

and

respectively.

From (\*) and (\*\*), we get

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

**Example 10.4.** The variables  $X$  and  $Y$  are connected by the equation  $aX + bY + c = 0$ . Show that the correlation between them is  $-1$  if the signs of  $a$  and  $b$  are alike and  $+1$  if they are different.

[Nagpur Univ. B.Sc. 1992; Delhi Univ. B.Sc. (Stat. Hons.) 1992]

**Solution.**  $aX + bY + c = 0 \Rightarrow aE(X) + bE(Y) + c = 0$

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

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

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

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

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

$$\therefore r = \frac{-\frac{b}{a} \cdot \sigma_Y^2}{\sqrt{\sigma_Y^2} \sqrt{\frac{b^2}{a^2} \cdot \sigma_Y^2}} = \frac{-\frac{b}{a} \sigma_Y^2}{\left| \frac{b}{a} \right| \sigma_Y^2}$$

$$= \begin{cases} +1, & \text{if } b \text{ and } a \text{ are of opposite signs.} \\ -1, & \text{if } b \text{ and } a \text{ are of same sign.} \end{cases}$$

**Example 10.5.** (a) If  $Z = aX + bY$  and  $r$  is the correlation coefficient between  $X$  and  $Y$ , show that

$$\sigma_Z^2 = a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab r \sigma_X \sigma_Y$$

(b) Show that the correlation coefficient  $r$  between two random variables  $X$  and  $Y$  is given by

$$r = (\sigma_X^2 + \sigma_Y^2 - \sigma_{X-Y}^2) / 2\sigma_X \sigma_Y$$

where  $\sigma_X$ ,  $\sigma_Y$  and  $\sigma_{X-Y}$  are the standard deviations of  $X$ ,  $Y$  and  $X - Y$  respectively.

[Calcutta Univ. B.Sc., 1992; M.S. Baroda Univ. B.Sc. 1992]

**Solution.** Taking expectation of both sides of  $Z = aX + bY$ , we get

$$E(Z) = aE(X) + bE(Y)$$

$$\therefore Z - E(Z) = a(X - E(X)) + b(Y - E(Y))$$

Squaring and taking expectation of both sides, we get

$$\begin{aligned} \sigma_Z^2 &= a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab \operatorname{Cov}(X, Y) \\ &= a^2 \sigma_X^2 + b^2 \sigma_Y^2 + 2ab r \sigma_X \sigma_Y \end{aligned}$$

(b) Taking  $a = 1$ ,  $b = -1$  in the above case, we have

$$Z = X - Y \text{ and } \sigma_{X-Y}^2 = \sigma_X^2 + \sigma_Y^2 - 2r \sigma_X \sigma_Y$$

$$\therefore r = \frac{\sigma_X^2 + \sigma_Y^2 - \sigma_{X-Y}^2}{2\sigma_X \sigma_Y}$$

**Remark.** In the above example, we have obtained

$$V(aX + bY) = a^2 V(X) + b^2 V(Y) + 2ab \operatorname{Cov}(X, Y)$$

Similarly, we could obtain the result

$$V(aX - bY) = a^2 V(X) + b^2 V(Y) - 2ab \operatorname{Cov}(X, Y)$$

The above results are useful in solving theoretical problems.

**Example 10.6.**  $X$  and  $Y$  are two random variables with variances  $\sigma_X^2$  and  $\sigma_Y^2$  respectively and  $r$  is the coefficient of correlation between them. If

$U = X + kY$  and  $V = X + \frac{\sigma_X}{\sigma_Y} Y$ , find the value of  $k$  so that  $U$  and  $V$  are uncorrelated.

[Delhi Univ. B.Sc. 1992; Andhra Univ. B.Sc. 1993]

**Solution.** Taking expectations of  $U = X + kY$  and  $V = X + \frac{\sigma_X}{\sigma_Y} Y$ , we get

$$E(U) = E(X) + kE(Y) \text{ and } E(V) = E(X) + \frac{\sigma_X}{\sigma_Y} E(Y)$$

$$U - E(U) = (X - E(X)) + k(Y - E(Y)) \text{ and}$$

$$V - E(V) = (X - E(X)) + \frac{\sigma_X}{\sigma_Y} (Y - E(Y))$$

$$\text{Cov}(U, V) = E[(U - E(U))(V - E(V))]$$

$$\begin{aligned} &= E[(X - E(X)) + k(Y - E(Y))] \times [(X - E(X)) + \frac{\sigma_X}{\sigma_Y} (Y - E(Y))] \\ &= \sigma_X^2 + \frac{\sigma_X}{\sigma_Y} \text{Cov}(X, Y) + k \text{Cov}(X, Y) + k \frac{\sigma_X}{\sigma_Y} \cdot \sigma_Y^2 \\ &= [\sigma_X^2 + k\sigma_X\sigma_Y] + \left[ \frac{\sigma_X}{\sigma_Y} + k \right] \text{Cov}(X, Y) \\ &= \sigma_X(\sigma_X + k\sigma_Y) + \left[ \frac{\sigma_X + k\sigma_Y}{\sigma_Y} \right] \text{Cov}(X, Y) \\ &= (\sigma_X + k\sigma_Y) \left[ \sigma_X + \frac{\text{Cov}(X, Y)}{\sigma_Y} \right] = (\sigma_X + k\sigma_Y)(1 + r)\sigma_X \end{aligned}$$

$U$  and  $V$  will be uncorrelated if

$$r(U, V) = 0 \Rightarrow \text{Cov}(U, V) = 0$$

$$\text{i.e., if } (\sigma_X + k\sigma_Y)(1 + r)\sigma_X = 0$$

$$\Rightarrow \sigma_X + k\sigma_Y = 0 \quad (\because \sigma_X \neq 0, r \neq -1)$$

$$\Rightarrow k = -\frac{\sigma_X}{\sigma_Y}$$

**Example 10.7.** The random variables  $X$  and  $Y$  are jointly normally distributed and  $U$  and  $V$  are defined by

$$U = X \cos \alpha + Y \sin \alpha,$$

$$V = Y \cos \alpha - X \sin \alpha$$

Show that  $U$  and  $V$  will be uncorrelated if

$$\tan 2\alpha = \frac{2r\sigma_X\sigma_Y}{\sigma_X^2 - \sigma_Y^2},$$

where  $r = \text{corr.}(X, Y)$ ,  $\sigma_X^2 = \text{Var}(X)$  and  $\sigma_Y^2 = \text{Var}(Y)$ . Are  $U$  and  $V$  then independent?

[Delhi Univ. B.Sc. (Stat. Hons.) 1989; (Maths. Hons.), 1990]

**Solution.** We have

$$\text{Cov}(U, V) = E[(U - E(U))(V - E(V))]$$

$$\begin{aligned} &= E[(X - E(X)) \cos \alpha + (Y - E(Y)) \sin \alpha] \\ &\quad \times [(Y - E(Y)) \cos \alpha - (X - E(X)) \sin \alpha] \end{aligned}$$

$$\begin{aligned}
 &= \cos^2 \alpha \operatorname{Cov}(X, Y) - \sin \alpha \cos \alpha \cdot \sigma_X^2 \\
 &\quad + \sin \alpha \cos \alpha \cdot \sigma_Y^2 - \sin^2 \alpha (\operatorname{Cov}(X, Y)) \\
 &= (\cos^2 \alpha - \sin^2 \alpha) \operatorname{Cov}(X, Y) - \sin \alpha \cos \alpha (\sigma_X^2 - \sigma_Y^2) \\
 &= \cos 2\alpha \operatorname{Cov}(X, Y) - \sin \alpha \cos \alpha (\sigma_X^2 - \sigma_Y^2)
 \end{aligned}$$

*U* and *V* will be uncorrelated if and only if

$$\begin{aligned}
 r(U, V) = 0, \text{ i.e., iff } \operatorname{Cov}(U, V) = 0 \\
 \text{i.e., if } \cos 2\alpha \operatorname{Cov}(X, Y) - \sin \alpha \cos \alpha (\sigma_X^2 - \sigma_Y^2) = 0
 \end{aligned}$$

$$\text{or if } \cos 2\alpha r \sigma_X \sigma_Y = \frac{\sin 2\alpha}{2} (\sigma_X^2 - \sigma_Y^2)$$

$$\text{or if } \tan 2\alpha = \frac{2r \sigma_X \sigma_Y}{\sigma_X^2 - \sigma_Y^2}$$

However,  $r(U, V) = 0$  does not imply that the variables *U* and *V* are independent. [For detailed discussion, see Theorem 10-2, page 10-4.]

**Example 10-8.** If *X*, *Y* are standardized random variables, and

$$r(aX + bY, bX + aY) = \frac{1 + 2ab}{a^2 + b^2} \quad \dots(*)$$

find  $r(X, Y)$ , the coefficient of correlation between *X* and *Y*.

[Sardar Patel Univ. B.Sc., 1993; Delhi Univ. B.Sc. (Stat. Hons.), 1989]

**Solution.** Since *X* and *Y* are standardised random variables, we have

$$\begin{aligned}
 \text{and } E(X) = E(Y) = 0 \\
 \text{and } \operatorname{Var}(X) = \operatorname{Var}(Y) = 1 \Rightarrow \bar{E}(X^2) = E(Y^2) = 1 \\
 \text{and } \operatorname{Cov}(X, Y) = E(XY) \Rightarrow E(XY) = r(X, Y) \cdot \sigma_X \sigma_Y \\
 &\qquad\qquad\qquad = r(X, Y) \quad \dots(**)
 \end{aligned}$$

Also we have

$$\begin{aligned}
 &r(aX + bY, bX + aY) \\
 &= \frac{E[(aX + bY)(bX + aY)] - E(aX + bY) E(bX + aY)}{[\operatorname{Var}(aX + bY) \cdot \operatorname{Var}(bX + aY)]^{1/2}} \\
 &= \frac{E[abX^2 + a^2 XY + b^2 YX + abY^2] - 0}{\{[a^2 \operatorname{Var}(X) + b^2 \operatorname{Var}(Y) + 2ab \operatorname{Cov}(X, Y)] \\
 &\quad \times [b^2 \operatorname{Var}(X) + a^2 \operatorname{Var}(Y) + 2ba \operatorname{Cov}(X, Y)]\}^{1/2}} \\
 &= \frac{ab \cdot 1 + a^2 r(X, Y) + b^2 r(X, Y) + ab \cdot 1}{\{[a^2 + b^2 + 2ab r(X, Y)][b^2 + a^2 + 2ba r(X, Y)]\}^{1/2}} \\
 &\qquad\qquad\qquad \text{[Using (**)]} \\
 &= \frac{2ab + (a^2 + b^2) \cdot r(X, Y)}{a^2 + b^2 + 2ab \cdot r(X, Y)}
 \end{aligned}$$

From (\*) and (\*\*), we get

$$\frac{1 + 2ab}{a^2 + b^2} = \frac{(a^2 + b^2) \cdot r(X, Y) + 2ab}{a^2 + b^2 + 2ab \cdot r(X, Y)}$$

Cross multiplying, we get

$$\begin{aligned}
 & (a^2 + b^2)(1 + 2ab) + 2ab \cdot r(X, Y)(1 + 2ab) = (a^2 + b^2)^2 \cdot r(X, Y) + 2ab(a^2 + b^2) \\
 \Rightarrow & (a^4 + b^4 + 2a^2b^2 - 2ab - 4a^2b^2) \cdot r(X, Y) = (a^2 + b^2)^2 \\
 \Rightarrow & [(a^2 - b^2)^2 - 2ab]r(X, Y) = a^2 + b^2 \\
 \Rightarrow & r(X, Y) = \frac{a^2 + b^2}{(a^2 - b^2)^2 - 2ab}
 \end{aligned}$$

**Example 10.9.** If  $X$  and  $Y$  are uncorrelated random variables with means zero and variances  $\sigma_1^2$  and  $\sigma_2^2$  respectively, show that

$$U = X \cos \alpha + Y \sin \alpha, V = X \sin \alpha - Y \cos \alpha$$

have a correlation coefficient  $\rho$  given by

$$\rho = \frac{\sigma_1^2 - \sigma_2^2}{[(\sigma_1^2 - \sigma_2^2)^2 + 4\sigma_1^2\sigma_2^2 \operatorname{cosec}^2 2\alpha]^{1/2}}$$

**Solution.** We are given that

$$r(X, Y) = 0 \Rightarrow \operatorname{Cov}(X, Y) = 0, \sigma_X^2 = \sigma_1^2 \text{ and } \sigma_Y^2 = \sigma_2^2 \quad \dots(1)$$

We have

$$\begin{aligned}
 \sigma_U^2 &= V(X \cos \alpha + Y \sin \alpha) \\
 &= \cos^2 \alpha V(X) + \sin^2 \alpha V(Y) + 2 \sin \alpha \cos \alpha \operatorname{Cov}(X, Y) \\
 &= \cos^2 \alpha \sigma_1^2 + \sin^2 \alpha \sigma_2^2 \quad [\text{Using (1)}]
 \end{aligned}$$

Similarly,

$$\sigma_V^2 = V(X \sin \alpha - Y \cos \alpha) = \sin^2 \alpha \cdot \sigma_1^2 + \cos^2 \alpha \cdot \sigma_2^2$$

$$\begin{aligned}
 \operatorname{Cov}(U, V) &= E[(U - E(U))(V - E(V))] \\
 &= E \left[ ((X - E(X)) \cos \alpha + (Y - E(Y)) \sin \alpha) \right. \\
 &\quad \times \left. ((X - E(X)) \sin \alpha - (Y - E(Y)) \cos \alpha) \right] \\
 &= \sin \alpha \cos \alpha V(X) - \cos^2 \alpha \operatorname{Cov}(X, Y) \\
 &\quad + \sin^2 \alpha \operatorname{Cov}(X, Y) - \sin \alpha \cos \alpha V(Y) \\
 &= (\sigma_1^2 - \sigma_2^2) \sin \alpha \cos \alpha \quad [\text{Using (1)}]
 \end{aligned}$$

$$\text{Now } \rho^2 = \frac{[\operatorname{Cov}(U, V)]^2}{\sigma_U^2 \sigma_V^2}$$

$$\begin{aligned}
 \text{where } \sigma_U^2 \sigma_V^2 &= (\cos^2 \alpha \sigma_1^2 + \sin^2 \alpha \sigma_2^2)(\sin^2 \alpha \sigma_1^2 + \cos^2 \alpha \sigma_2^2) \\
 &= \sin^2 \alpha \cos^2 \alpha (\sigma_1^4 + \sigma_2^4) + \sigma_1^2 \sigma_2^2 (\cos^4 \alpha + \sin^4 \alpha) \\
 &= \sin^2 \alpha \cos^2 \alpha (\sigma_1^4 + \sigma_2^4) + \sigma_1^2 \sigma_2^2 [(\sin^2 \alpha + \cos^2 \alpha)^2 - 2 \sin^2 \alpha \cos^2 \alpha] \\
 &= \sin^2 \alpha \cos^2 \alpha (\sigma_1^4 + \sigma_2^4 - 2\sigma_1^2 \sigma_2^2) + \sigma_1^2 \sigma_2^2 \\
 &= \sin^2 \alpha \cos^2 \alpha (\sigma_1^2 - \sigma_2^2)^2 + \sigma_1^2 \sigma_2^2
 \end{aligned}$$

$$\begin{aligned}
 \therefore \rho^2 &= \frac{(\sigma_1^2 - \sigma_2^2)^2 \cdot \sin^2 \alpha \cos^2 \alpha}{\sigma_1^2 \sigma_2^2 + \sin^2 \alpha \cos^2 \alpha (\sigma_1^2 - \sigma_2^2)^2} \\
 &= \frac{\frac{1}{4}(\sigma_1^2 - \sigma_2^2)^2 \sin^2 2\alpha}{\sigma_1^2 \sigma_2^2 + \sin^2 2\alpha \cdot \frac{1}{4}(\sigma_1^2 - \sigma_2^2)^2}
 \end{aligned}$$

$$= \frac{(\sigma_1^2 - \sigma_2^2)^2}{4\sigma_1^2\sigma_2^2 \cosec^2 2\alpha + (\sigma_1^2 - \sigma_2^2)^2}$$

$$\Rightarrow \rho = \frac{\sigma_1^2 - \sigma_2^2}{[(\sigma_1^2 - \sigma_2^2)^2 + 4\sigma_1^2\sigma_2^2 \cosec^2 2\alpha]^{1/2}}$$

**Example 10.10.** If  $U = aX + bY$  and  $V = cX + dY$ , where  $X$  and  $Y$  are measured from their respective means and if  $r$  is the correlation coefficient between  $X$  and  $Y$ , and if  $U$  and  $V$  are uncorrelated, show that

$$\sigma_U \sigma_V = (ad - bc) \sigma_X \sigma_Y (1 - r^2)^{1/2}$$

[Poona Univ. B.Sc., 1990; Delhi Univ. B.Sc. (Stat. Hons.), 1986]

**Solution.** We have

$$r = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} \Rightarrow 1 - r^2 = 1 - \frac{[\text{Cov}(X, Y)]^2}{\sigma_X^2 \sigma_Y^2}$$

$$\Rightarrow (1 - r^2) \sigma_X^2 \sigma_Y^2 = \sigma_X^2 \sigma_Y^2 - [\text{Cov}(X, Y)]^2 \quad \dots(*)$$

[This step is suggested by the answer]

$$U = aX + bY, V = cX + dY$$

Since  $X, Y$  are measured from their means,

$$E(X) = 0 = E(Y) \Rightarrow E(U) = 0 = E(V) \quad \left. \begin{array}{l} \\ \end{array} \right\} \quad \dots(**)$$

and  $\sigma_U^2 = E(U^2); \sigma_V^2 = E(V^2)$

$$\text{Also } aX + bY - U = 0 \text{ and } cX + dY - V = 0$$

$$\Rightarrow \frac{X}{-bV + dU} = \frac{Y}{-cU + aV} = \frac{1}{ad - bc}$$

$$\Rightarrow \left. \begin{array}{l} X = \frac{1}{ad - bc} (dU - bV) \\ Y = \frac{1}{ad - bc} (-cU + aV) \end{array} \right\} \quad \dots(***)$$

$$\therefore \text{Var}(X) = \frac{1}{(ad - bc)^2} [d^2 \sigma_U^2 + b^2 \sigma_V^2 - 2bd \text{Cov}(U, V)]$$

$$= \frac{1}{(ad - bc)^2} [d^2 \sigma_U^2 + b^2 \sigma_V^2]$$

[Since  $U, V$  are uncorrelated  $\Leftrightarrow \text{Cov}(U, V) = 0$ ]

Similarly, we have

$$\text{Var}(Y) = \frac{1}{(ad - bc)^2} (c^2 \sigma_U^2 + a^2 \sigma_V^2)$$

$$\text{Cov}(X, Y) = E(XY) - E(X) E(Y) = E(XY) \quad \left[ \because E(X) = 0 = E(Y) \right]$$

$$= \frac{1}{(ad - bc)^2} E[(dU - bV)(-cU + aV)] \quad \text{[From (***)]}$$

$$\begin{aligned}
 &= \frac{1}{(ad - bc)^2} [-cd \sigma_U^2 - ab \sigma_V^2] \\
 &\quad [\text{Using } (***) \text{ and } \text{Cov}(U, V) = 0, \text{ given}] \\
 &= \frac{-1}{(ad - bc)^2} [cd \sigma_U^2 + ab \sigma_V^2]
 \end{aligned}$$

Substituting in (\*), we get

$$\begin{aligned}
 (1 - r^2) \sigma_X^2 \sigma_Y^2 &= \frac{1}{(ad - bc)^4} \times [(d^2 \sigma_U^2 + b^2 \sigma_V^2)(c^2 \sigma_U^2 + a^2 \sigma_V^2) \\
 &\quad - (cd \sigma_U^2 + ab \sigma_V^2)^2] \\
 &= \frac{1}{(ad - bc)^4} \\
 &\quad \times [c^2 d^2 \sigma_U^4 + a^2 b^2 \sigma_V^4 + (a^2 d^2 + b^2 c^2) \sigma_U^2 \sigma_V^2 \\
 &\quad - c^2 d^2 \sigma_U^4 - a^2 b^2 \sigma_V^4 - 2abcd \sigma_U^2 \sigma_V^2] \\
 &= \frac{1}{(ad - bc)^4} [a^2 d^2 + b^2 c^2 - 2abcd] \sigma_U^2 \sigma_V^2 \\
 &= \frac{1}{(ad - bc)^4} (ad - bc)^2 \sigma_U^2 \sigma_V^2 \\
 &= \frac{\sigma_U^2 \sigma_V^2}{(ad - bc)^2}
 \end{aligned}$$

Cross multiplying and taking square root, we get the required result.

**Example 10.11. (a) Establish the formula :**

$$nr\sigma_X\sigma_Y = n_1r_1\sigma_{X_1}\sigma_{Y_1} + n_2r_2\sigma_{X_2}\sigma_{Y_2} + n_1dx_1dy_1 + n_2dx_2dy_2 \quad \dots(10.5)$$

where  $n_1$ ,  $n_2$  and  $n$  are respectively the sizes of the first, second and combined sample :  $(\bar{x}_1, \bar{y}_1)$ ,  $(\bar{x}_2, \bar{y}_2)$ ,  $(\bar{x}, \bar{y})$ , their means  $r_1$ ,  $r_2$  and  $r$  their coefficients of correlation;  $(\sigma_{X_1}, \sigma_{Y_1})$ ,  $(\sigma_{X_2}, \sigma_{Y_2})$ ,  $(\sigma_X, \sigma_Y)$  their standard deviations, and

$$dx_1 = \bar{x}_1 - \bar{x} \quad , \quad dy_1 = \bar{y}_1 - \bar{y}$$

$$dx_2 = \bar{x}_2 - \bar{x} \quad , \quad dy_2 = \bar{y}_2 - \bar{y}$$

**(b) Find the correlation co-efficient of combined sample given that**

	Sample I	Sample II
Sample size	100	150
Sample mean ( $\bar{x}$ )	80	72
Sample mean ( $\bar{y}$ )	100	118
Sample variance ( $\sigma_X^2$ )	10	12
Sample variance ( $\sigma_Y^2$ )	15	18
Correlation coefficient	0.6	0.4

**Solution.** (a) Let  $(x_{1i}, y_{1i})$ ;  $i = 1, 2, \dots, n_1$  and  $(x_{2j}, y_{2j})$ ;  $j = 1, 2, \dots, n_2$ , be the two samples of sizes  $n_1$  and  $n_2$  respectively from the bivariate population. Then with the given notations, we have

$$\begin{aligned}\bar{x} &= \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}, \quad \bar{y} = \frac{n_1 \bar{y}_1 + n_2 \bar{y}_2}{n_1 + n_2} \\ n\sigma_x^2 &= n_1 (\sigma_{x_1}^2 + dx_1^2) + n_2 (\sigma_{x_2}^2 + dx_2^2) \\ n\sigma_y^2 &= n_1 (\sigma_{y_1}^2 + dy_1^2) + n_2 (\sigma_{y_2}^2 + dy_2^2)\end{aligned}\quad \dots(1)$$

$$r_1 = \frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{x}_1)(y_{1i} - \bar{y}_1)}{n_1 \sigma_{x_1} \sigma_{y_1}}, \quad r_2 = \frac{\sum_{j=1}^{n_2} (x_{2j} - \bar{x}_2)(y_{2j} - \bar{y}_2)}{n_2 \sigma_{x_2} \sigma_{y_2}} \quad \dots(2)$$

For the pooled sample, we have

$$r = \frac{\sum_{i=1}^{n_1} (x_{1i} - \bar{x})(y_{1i} - \bar{y}) + \sum_{j=1}^{n_2} (x_{2j} - \bar{x})(y_{2j} - \bar{y})}{n \sigma_x \sigma_y} \quad \dots(3)$$

Now

$$\begin{aligned}\sum_{i=1}^{n_1} (x_{1i} - \bar{x})(y_{1i} - \bar{y}) &= \sum_{i=1}^{n_1} \left\{ \{(x_{1i} - \bar{x}_1) + (\bar{x}_1 - \bar{x})\} \{(y_{1i} - \bar{y}_1) + (\bar{y}_1 - \bar{y})\} \right\} \\ &= \sum_{i=1}^{n_1} (x_{1i} - \bar{x}_1)(y_{1i} - \bar{y}_1) + (\bar{y}_1 - \bar{y}) \sum_{i=1}^{n_1} (x_{1i} - \bar{x}_1) \\ &\quad + (\bar{x}_1 - \bar{x}) \sum_{i=1}^{n_1} (y_{1i} - \bar{y}_1) + n_1 (\bar{x}_1 - \bar{x})(\bar{y}_1 - \bar{y})\end{aligned}$$

$$\text{But } \sum_{i=1}^{n_1} (x_{1i} - \bar{x}_1) = 0 \text{ and } \sum_{i=1}^{n_1} (y_{1i} - \bar{y}_1) = 0,$$

being the algebraic sum of the deviations from the mean.

$$\therefore \sum_{i=1}^{n_1} (x_{1i} - \bar{x})(y_{1i} - \bar{y}) = n_1 r_1 \sigma_{x_1} \sigma_{y_1} + n_1 d x_1 d y_1 \quad [\text{Using (2)}]$$

Similarly, we will get

$$\sum_{j=1}^{n_2} (x_{2j} - \bar{x})(y_{2j} - \bar{y}) = n_2 r_2 \sigma_{x_2} \sigma_{y_2} + n_2 d x_2 d y_2$$

Substituting in (3), we get the required formula.

(b) Here we are given :

$$n_1 = 100, \bar{x}_1 = 80, \bar{y}_1 = 100, \sigma_{x_1}^2 = 10, \sigma_{y_1}^2 = 15, r_1 = 0.6$$

$$n_2 = 150, \bar{x}_2 = 72, \bar{y}_2 = 118, \sigma_{x_2}^2 = 12, \sigma_{y_2}^2 = 18, r_2 = 0.4$$

$$\therefore \bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2} = \frac{100 \times 80 + 150 \times 72}{100 + 150} = 75.2$$

$$\bar{y} = \frac{n_1 \bar{y}_1 + n_2 \bar{y}_2}{n_1 + n_2} = \frac{100 \times 100 + 150 \times 118}{100 + 150} = 110.8$$

$$dx_1 = \bar{x}_1 - \bar{x} = 4.8, \quad dy_1 = \bar{y}_1 - \bar{y} = 10.8$$

$$dx_2 = \bar{x}_2 - \bar{x} = 3.2, \quad dy_2 = \bar{y}_2 - \bar{y} = 7.2$$

$$n\sigma_x^2 = n_1 (\sigma_{x_1}^2 + dx_1^2) + n_2 (\sigma_{x_2}^2 + dx_2^2) = 6640$$

$$n\sigma_y^2 = n_1 (\sigma_{y_1}^2 + dy_1^2) + n_2 (\sigma_{y_2}^2 + dy_2^2) = 23640$$

Substituting these values in the formula and simplifying, we get

$$r = \frac{n_1 r_1 \sigma_{x_1} \sigma_{y_1} + n_2 r_2 \sigma_{x_2} \sigma_{y_2} + n_1 dx_1 dy_1 + n_2 dx_2 dy_2}{n \sigma_x \sigma_y} = 0.8186$$

**Example 10-12.** The independent variables  $X$  and  $Y$  are defined by :

$$\begin{aligned} f(x) &= 4ax, \quad 0 \leq x \leq r \\ &= 0, \quad \text{otherwise} \end{aligned} \quad \quad \quad \begin{aligned} f(y) &= 4by, \quad 0 \leq y \leq s \\ &= 0, \quad \text{otherwise} \end{aligned}$$

Show that :

$$\text{Cov}(U, V) = \frac{b-a}{b+a},$$

$$\text{where } U = X + Y \quad \text{and} \quad V = X - Y$$

[I.I.T. (B. Tech.), Nov. 1992]

**Solution.** Since the total area under probability curve is unity (one), we have :

$$\int_0^r f(x) dx = 4a \int_0^r x dx = 1 \Rightarrow 2ar^2 = 1 \Rightarrow a = \frac{1}{2r^2} \quad \dots(i)$$

$$\int_0^s f(y) dy = 4b \int_0^s y dy = 1 \Rightarrow 2bs^2 = 1 \Rightarrow b = \frac{1}{2s^2} \quad \dots(ii)$$

$$\therefore f(x) = 4ax = \frac{2x}{r^2}, \quad 0 \leq x \leq r; \quad \text{and} \quad f(y) = 4by = \frac{2y}{s^2}, \quad 0 \leq y \leq s \quad \dots(iii)$$

Since  $X$  and  $Y$  are independent variates,

$$\text{Cov}(X, Y) = 0 \Rightarrow \text{Cov}(X, Y) = 0 \quad \dots(iv)$$

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

$$= \text{Cov}(X, X) - \text{Cov}(X, Y) + \text{Cov}(Y, X) - \text{Cov}(Y, Y)$$

$$= \sigma_x^2 - \sigma_y^2 \quad [\text{Using (iv)}]$$

$$\text{Var}(U) = \text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2 \text{Cov}(X, Y)$$

$$= \sigma_x^2 + \sigma_y^2 \quad [\text{Using (iv)}]$$

$$\text{Var}(V) = \text{Var}(X - Y) = \text{Var}(X) + \text{Var}(Y) - 2 \text{Cov}(X, Y)$$

$$= \sigma_x^2 + \sigma_y^2 \quad [\text{Using (iv)}]$$

$$\therefore r(U, V) = \frac{\text{Cov}(U, V)}{\sigma_U \sigma_V} = \frac{\sigma_X^2 - \sigma_Y^2}{\sigma_X^2 + \sigma_Y^2} \quad \dots (v)$$

We have :

$$E(X) = \int_0^r x f(x) dx = \frac{2}{r^2} \int_0^r x^2 dx = \frac{2r}{3} \quad [\text{From (iii)}]$$

$$E(X^2) = \int_0^r x^2 f(x) dx = \frac{2}{r^2} \int_0^r x^3 dx = \frac{r^2}{2}$$

$$\therefore \text{Var}(X) = E(X^2) - [E(X)]^2 = \frac{r^2}{2} - \frac{4r^2}{9} = \frac{r^2}{18} = \frac{1}{36a} \quad [\text{From (i)}]$$

Similarly, we shall get

$$E(Y) = \frac{2s}{3}, E(Y^2) = \frac{s^2}{2} \text{ and } \text{Var}(Y) = \frac{s^2}{18} = \frac{1}{36b}$$

Substituting in (v), we get

$$r(U, V) = \frac{1/(36a) - 1/(36b)}{1/(36a) + 1/(36b)} = \frac{b-a}{b+a}$$

**Example 10.13.** Let the random variable  $X$  have the marginal density

$$f_1(x) = 1, -\frac{1}{2} < x < \frac{1}{2}$$

and let the conditional density of  $Y$  be

$$\left. \begin{aligned} f(y|x) &= 1, x < y < x+1, -\frac{1}{2} < x < 0 \\ &= 1, -x < y < 1-x, 0 < x < \frac{1}{2} \end{aligned} \right\} \quad (*)$$

Show that the variables  $X$  and  $Y$  are uncorrelated.

**Solution.** We have

$$E(X) = \int_{-\frac{1}{2}}^{\frac{1}{2}} x f_1(x) dx = \int_{-\frac{1}{2}}^{\frac{1}{2}} x \cdot 1 dx = \left[ \frac{x^2}{2} \right]_{-\frac{1}{2}}^{\frac{1}{2}} = 0$$

If  $f(x, y)$  is the joint p.d.f. of  $X$  and  $Y$ , then

$$f(x, y) = f(y|x) f_1(x) = f(y|x). \quad (***) \quad [\because f_1(x) = 1]$$

$$\begin{aligned} E(XY) &= \int_{-\frac{1}{2}}^0 \int_x^{x+1} xy f(x, y) dx dy + \int_0^{\frac{1}{2}} \int_{-x}^{1-x} xy f(x, y) dx dy \\ &= \int_{-\frac{1}{2}}^0 \left[ x \int_x^{x+1} y dy \right] dx + \int_0^{\frac{1}{2}} \left[ x \int_{-x}^{1-x} y dy \right] dx \quad [\text{From (*) and}] \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{2} \int_{-\frac{1}{2}}^0 \dot{x}(2x+1)dx + \frac{1}{2} \int_0^{\frac{1}{2}} x(1-2x) dx \\
 &= \frac{1}{2} \left[ \frac{2}{3}x^3 + \frac{x^2}{2} \right]_{-\frac{1}{2}}^0 + \frac{1}{2} \left[ \frac{x^2}{2} - \frac{2}{3}x^3 \right]_0^{\frac{1}{2}} \\
 &= \frac{1}{2} \left[ \frac{1}{12} - \frac{1}{8} - \frac{1}{12} + \frac{1}{8} \right] = 0
 \end{aligned}$$

$$\therefore \text{Cov } (XY) = E(XY) - E(X) E(Y) = 0 \Rightarrow r(X, Y) = 0$$

Hence the variables  $X$  and  $Y$  are uncorrelated.

### EXERCISE 10(a)

1. (a) Show that the co-efficient of correlation  $r$  is independent of a change of scale and origin of the variables. Also prove that for two independent variables  $r = 0$ . Show by an example that the converse is not true. State the limits between which  $r$  lies and give its proof.

[Delhi Univ. M.Sc. (O.R.), 1986]

- (b) Let  $\rho$  be the correlation coefficient between two jointly distributed random variables  $X$  and  $Y$ . Show that  $|\rho| \leq 1$  and that  $|\rho| = 1$  if and only if  $X$  and  $Y$  are linearly related. [Indian Forest Service, 1991]

2. (a) Calculate the coefficient of correlation between  $X$  and  $Y$  for the following :

X...	1	3	4	5	7	8	10
Y...	2	6	8	10	14	16	20

$$\text{Ans. } r(X, Y) = +1$$

(b) Discuss the statistical validity of the following statements :

- (i) "High positive coefficient of correlation between increase in the sale of newspapers and increase in the number of crimes leads to the conclusion that newspaper reading may be responsible for the increase in the number of crimes."

- (ii) "A high positive value of  $r$  between the increase in cigarette smoking and increase in lung cancer establishes that cigarette smoking is responsible for lung cancer."

- (c) (i) Do you agree with the statement that " $r = 0.8$  implies that 80% of the data are explained."

(ii) Comment on the following :

"The closeness of relationship between two variables is proportional to  $r$ ".

Hint. (a) No (b) Wrong.

- (d) By effecting suitable change of origin and scale, compute the product moment correlation coefficient for the following set of 5 observations on  $(X, Y)$  :

X :	-10	-5	0	5	10
Y :	5	9	7	11	13

$$\text{Ans. } r(X, Y) = 0.34$$

3. The marks obtained by 10 students in Mathematics and Statistics are given below. Find the coefficient of correlation between the two subjects.

Roll No.	1	2	3	4	5	6	7	8	9	10
Marks in Mathematics :	75	30	60	80	53	35	15	40	38	48
Marks in Statistics :	85	45	54	91	58	63	35	43	45	44

4. (a) The following table gives the number of blind per lakh of population in different age-groups. Find out the correlation between age and blindness.

Age in years	0—10	10—20	20—30	30—40	40—50
Number of blind per lakh	55	67	100	111	150
Age in year	50—60	60—70	70—80		
Number of blind per lakh	200	300.	500		

Ans. 0.89

(b) The following table gives the distribution of items of production and also the relatively defective items among them, according to size-groups. Is there any correlation between size and defect in quality?

Size-Group	15—16	16—17	17—18	18—19	19—20	20—21
No. of Items	200	270	340	360	400	300
No. of defective items	150	162	170	180	180	120

Hint. Here we have to find the correlation coefficient between the size-group ( $X$ ) and the percentage of defectives ( $Y$ ) given below.

$X$	15.5	16.5	17.5	18.5	19.5	20.5
$Y$	75	60	50	50	45	40

Ans.  $r = 0.94$ .

5. Using the formula

$$\sigma_{X-Y}^2 = \sigma_X^2 + \sigma_Y^2 - 2 r(X, Y) \sigma_X \sigma_Y$$

obtain the correlation coefficient between the heights of fathers ( $X$ ) and of the sons ( $Y$ ) from the following data :

$X$ :	65	66	67	68	69	70	71	67
$Y$ :	67	68	64	72	70	67	70	68

6. (a) From the following data, compute the co-efficient of correlation between  $X$  and  $Y$ .

	$X$ series	$Y$ series
No. of items	15	15
Arithmetic mean	25	18
Sum of squares of deviations from mean	136	138

Summation of product of deviations of  $X$  and  $Y$  series from the respective arithmetic means = 122.

**Ans.**  $r(X, Y) = 0.891$

(b) Coefficient of correlation between two variables  $X$  and  $Y$  is 0.32. Their covariance is 7.86. The variance of  $X$  is 10. Find the standard deviation of  $Y$  series.

(c) In two sets of variables  $X$  and  $Y$  with 50 observations each, the following data were observed :

$$\bar{X} = 10, \sigma_X = 3, \bar{Y} = 6, \sigma_Y = 2 \text{ and } r(X, Y) = 0.3$$

But on subsequent verification it was found that one value of  $X$  (= 10) and one value of  $Y$  (= 6) were inaccurate and hence weeded out. With the remaining 49 pairs of values, how is the original value of  $r$  affected ?

*(Nagpur Univ. B.Sc., 1990)*

**Hint.**  $\Sigma X = n\bar{X} = 500, \Sigma Y = n\bar{Y} = 300$

$$\Sigma X^2 = n(\sigma_X^2 + \bar{X}^2) = 5450, \Sigma Y^2 = 50(4 + 36) = 2000$$

$$r \sigma_X \sigma_Y = \text{Cov}(X, Y) = \frac{\Sigma XY}{n} - \bar{X} \bar{Y}$$

$$\Rightarrow 0.3 \times 3 \times 2 = \frac{\Sigma XY}{50} - 10 \times 6$$

$$\Rightarrow \Sigma XY = 50(1.8 + 60) = 3090$$

After weeding out the incorrect pair of observation, viz.,  $(X = 10, Y = 6)$ , the corrected values of  $\Sigma X$ ,  $\Sigma Y$ ,  $\Sigma X^2$ ,  $\Sigma Y^2$  and  $\Sigma XY$  for the remaining  $50 - 1 = 49$  pairs of observations are given below :

**Corrected Values :**

$$\Sigma X = 500 - 10 = 490; \Sigma Y = 300 - 6 = 294$$

$$\Sigma XY = 3090 - 10 \times 6 = 3090 - 60 = 3030$$

$$\Sigma X^2 = 5450 - 10^2 = 5350, \Sigma Y^2 = 2000 - 6^2 = 1964$$

$$\therefore r = \frac{\text{Corrected Cov}(X, Y)}{(\text{Corrected } \sigma_X) \times (\text{Corrected } \sigma_Y)} = \frac{90/49}{\sqrt{\frac{1450}{49} \times \frac{200}{49}}} = 0.3$$

Hence the correlation coefficient is invariant in this case.

(d) A prognostic test in Mathematics was given to 10 students who were about to begin a course in Statistics. The scores ( $X$ ) in their test were examined in relation to scores ( $Y$ ) in the final examination in Statistics. The following results were obtained :—

$$\Sigma X = 71, \Sigma Y = 70, \Sigma X^2 = 555, \Sigma Y^2 = 526 \text{ and } \Sigma XY = 527$$

Find the coefficient of correlation between  $X$  and  $Y$ .

*(Kerala Univ. B.Sc., 1990)*

7. (a)  $X_1$  and  $X_2$  are independent variables with means 5 and 10 and standard deviations 2 and 3 respectively. Obtain  $r(U, V)$  where

$$U = 3X_1 + 4X_2 \text{ and } V = 3X_1 - X_2$$

**Ans.** 0

*(Delhi Univ. B.Sc., 1988)*

(b) If  $X$  and  $Y$  are normal and independent with zero means and standard deviations 9 and 12 respectively, and if  $X + 2Y$  and  $kX - Y$  are non-correlated, find  $k$ .

(c)  $X, Y, Z$  are random variables each with expectation 10 and variances 1, 4 and 9 respectively. The correlation coefficients are

$$r(X, Y) = 0, r(Y, Z) = r(X, Y) = 1/4$$

Obtain the numerical values of :

- (i)  $E(X + Y - 2Z)$ , (ii)  $\text{Cov}(X + 3, Y + 3)$ , (iii)  $V(X - 2Z)$  and  
 (iv)  $\text{Cov}(3X, 5Z)$

**Ans.** (i) 0, (ii) 0, (iii) 34, and (iv) 45/4.

(d)  $X$  and  $Y$  are discrete random variables. If  $\text{Var}(X) = \text{Var}(Y) = \sigma^2$ ,

$\text{Cov}(X, Y) = \frac{\sigma^2}{2}$ , find (i)  $\text{Var}(2X - 3Y)$ , (ii)  $\text{Corr}(2X + 3, 2Y - 3)$ .

8. (a) Prove that :

$$V(aX \pm bY) = a^2V(X) + b^2V(Y) \pm 2ab \text{Cov}(X, Y)$$

Hence deduce that if  $X$  and  $Y$  are independent

$$V(X \pm Y) = V(X) + V(Y)$$

(b) Prove that correlation coefficient between  $X$  and  $Y$  is positive or negative according as

$$\sigma_{X+Y} > \text{or} < \sigma_{X-Y}$$

9. Show that if  $X$  and  $Y$  are two random variables each assuming only two values and the correlation co-efficient between them is zero, then they are independent. Indicate with justification whether the result is true in general.

Find the correlation coefficient between  $X$  and  $a - X$ , where  $X$  is any random variable and  $a$  is constant.

10. (a)  $X_i$  ( $i = 1, 2, 3$ ) are uncorrelated variables each having the same standard deviation. Obtain the correlation between  $X_1 + X_2$  and  $X_2 + X_3$ .

**Ans.** 1/2

(b) If  $X_i$  ( $i = 1, 2, 3$ ) are three uncorrelated variables having standard deviations  $\sigma_1, \sigma_2$  and  $\sigma_3$  respectively, obtain the coefficient of correlation between  $(X_1 + X_2)$  and  $(X_2 + X_3)$ .

**Ans.**  $\sigma_2^2 / \sqrt{(\sigma_1^2 + \sigma_2^2)(\sigma_2^2 + \sigma_3^2)}$

(c) Two random variables  $X$  and  $Y$  have zero means, the same variance  $\sigma^2$  and zero correlation. Show that

$$U = X \cos \alpha + Y \sin \alpha \quad \text{and} \quad V = X \sin \alpha - Y \cos \alpha$$

have the same variance  $\sigma^2$  and zero correlation.

**(Bangalore Univ. B.Sc., 1991)**

(d) Let  $X$  and  $Y$  be uncorrelated random variables. If  $U = X + Y$  and  $V = X - Y$ , prove that the coefficient of correlation between  $U$  and  $V$  is  $(\sigma_X^2 - \sigma_Y^2) / (\sigma_X^2 + \sigma_Y^2)$ , where  $\sigma_X^2$  and  $\sigma_Y^2$  are variances of  $X$  and  $Y$  respectively.

(e) Two independent random variables  $X$  and  $Y$  have the following variances :  $\sigma_X^2 = 36, \sigma_Y^2 = 16$ . Calculate the coefficient of correlation between

$$U = X + Y \text{ and } V = X - Y$$

(f) Random variables  $X$  and  $Y$  have zero means and non-zero variances  $\sigma_X^2$  and  $\sigma_Y^2$ . If  $Z = Y - X$ , then find  $\sigma_Z$  and the correlation coefficient  $\rho(X, Z)$  of  $X$  and  $Z$  in terms of  $\sigma_X$ ,  $\sigma_Y$  and the correlation coefficient  $\rho(X, Y)$  of  $X$  and  $Y$ .

(g) If the independent random variables  $X_1, X_2$  and  $X_3$  have the means 4, 9 and 3 and variances 3, 7, 5, respectively, obtain the mean and variance of

$$(i) \quad Y = 2X_1 - 3X_2 + 4X_3, \quad (ii) \quad Z = X_1 + 2X_2 - X_3, \text{ and}$$

(iii) Calculate the correlation between  $Y$  and  $Z$ .

[Delhi Univ. M.A.(Eco.), 1989]

11. (a)  $X_1, X_2, \dots, X_n$  are uncorrelated random variables, all with the same distribution and zero means. Let  $\bar{X} = \sum X_i/n$

Find the correlation coefficient between (i)  $X_i$  and  $\bar{X}$  and (ii)  $X_i - \bar{X}$  and  $\bar{X}$ .

[Delhi Univ. B.Sc. (Stat. Hons.), 1993]

Hint.  $r(X_i, \bar{X}) = \frac{\sigma^2/n}{\sqrt{\sigma^2 \cdot \sigma^2/n}} = \frac{1}{\sqrt{n}}$

$$\begin{aligned} \text{Cov}(X_i - \bar{X}, \bar{X}) &= \text{Cov}(X_i, \bar{X}) - \text{Var}(\bar{X}) \\ &= (\sigma^2/n) - (\sigma^2/n) = 0 \end{aligned}$$

$$\therefore r(X_i - \bar{X}, \bar{X}) = 0$$

(b)  $X_1, X_2, \dots, X_n$  are random variables each with the same expected value  $\mu$  and s.d.  $\sigma$ . The correlation coefficient between any two  $X$ 's is  $\rho$ . Show

that (i)  $\text{Var}(\bar{X}) = \frac{\sigma^2}{n} + \left(1 - \frac{1}{n}\right)\rho\sigma^2$ ,

$$(ii) \quad E \sum_{i=1}^n (X_i - \bar{X})^2 = (n-1)(1-\rho)\sigma^2, \text{ and } (iii) \quad \rho > -\frac{1}{n-1}$$

12. (a) If  $X$  and  $Y$  are independent random variables, show that

$$r(X+Y, X-Y) = r^2(X, X+Y) - r^2(Y, X+Y),$$

where  $r(X+Y, X-Y)$  denotes the coefficient of correlation between  $(X+Y)$  and  $(X-Y)$ .

(Meerut Univ. B.Sc., 1991)

(b) Let  $X$  and  $Y$  be random variables having mean 0, variance 1 and correlation  $r$ . Show that  $X - rY$  and  $Y$  are uncorrelated and that  $X - rY$  has mean zero and variance  $1 - r^2$ .

13.  $X_1$  and  $X_2$  are two variables with zero means, variances  $\sigma_1^2$  and  $\sigma_2^2$  respectively and  $r$  is the correlation coefficient between them. Determine the values of the constants  $a$  and  $b$  which are independent of  $r$  such that  $X_1 + aX_2$  and  $X_1 + bX_2$  are uncorrelated.

14. (a) If  $X_1$  and  $X_2$  are two random variables with means  $\mu_1$  and  $\mu_2$ , variances  $\sigma_1^2, \sigma_2^2$  and correlation coefficient  $r$ , find the correlation coefficient between

$$U = a_1X_1 + a_2X_2 \text{ and } V = b_1X_1 + b_2X_2,$$

where  $a_1, a_2$  and  $b_1, b_2$  are constants.

(b) Let  $X_1, X_2$  be independent random variables with means  $\mu_1, \mu_2$  and non-zero variances  $\sigma_1^2, \sigma_2^2$  respectively. Let  $U = X_1 - X_2$  and  $V = X_1 X_2$ . Find the correlation coefficient between (i)  $X_1$  and  $U$ , (ii)  $X_1$  and  $V$ , in terms of  $\mu_1, \mu_2, \sigma_1^2, \sigma_2^2$ .

15. (a) If  $U = aX + bY$  and  $V = bX - aY$ , where  $X$  and  $Y$  are measured from their respective means and if  $U$  and  $V$  are uncorrelated,  $r$  the co-efficient of correlation between  $X$  and  $Y$  is given by the equation.

$$\sigma_U \sigma_V = (a^2 + b^2) \sigma_X \sigma_Y (1 - r^2)^{1/2} \quad (\text{Utkal Univ. B. Sc., 1993})$$

(b) Let  $U = aX + bY$  and  $V = aX - bY$  where  $X, Y$  represent deviations from the means of two measurements on the same individual. The coefficient of correlation between  $X$  and  $Y$  is  $\rho$ . If  $U, V$  are uncorrelated, show that

$$\sigma_U \sigma_V = 2ab\sigma_X \sigma_Y (1 - r^2)^{1/2}$$

16. Show that, if  $a$  and  $b$  are constants and  $r$  is the correlation coefficient between  $X$  and  $Y$ , then the correlation coefficient between  $aX$  and  $bY$  is equal to  $r$  if the signs of  $a$  and  $b$  are alike, and to  $-r$  if they are different.

Also show that, if constants  $a, b$  and  $c$  are positive, the correlation coefficient between  $(aX + bY)$  and  $cY$  is equal to

$$(ar\sigma_X + b\sigma_Y) / \sqrt{(a^2\sigma_X^2 + b^2\sigma_Y^2 + 2ab\sigma_X\sigma_Y)}$$

17. If  $X_1, X_2$  and  $X_3$  are three random variables measured from their respective means as origin and of equal variances, find the coefficient of correlation between  $X_1 + X_2$  and  $X_2 + X_3$  in terms of  $r_{12}, r_{13}$  and  $r_{23}$  and show that it is equal to

$$(i) \frac{r_{12} + 1}{2}, \text{ if } r_{13} = r_{23} = 0, \text{ and } (ii) \frac{r_{12} + 3}{4}, \text{ if } r_{13} = r_{23} = 1$$

18. (a) For a weighted distribution  $(x_i, w_i)$ , ( $i = 1, 2, \dots, n$ ) show that the weighted arithmetic mean  $\bar{x}_w = \sum w_i x_i / \sum w_i >$  or  $<$  the unweighted mean  $\bar{x} = \sum x_i / n$  according as  $r_{xw} >$  or  $< 0$ .

(b) Given  $N$  values  $x_1, x_2, \dots, x_N$  of variable  $X$  and weights  $w_1, w_2, \dots, w_N$ , express the coefficient of correlation between  $X$  and  $W$  in terms involving the difference between the arithmetic mean and the weighted mean of  $X$ .

19. (a) A coin is tossed  $n$  times. If  $X$  and  $Y$  denote the (random) number of heads and number of tails turned up respectively, show that  $r(X, Y) = -1$ .

**Hint.** Note that  $X + Y = n \Rightarrow Y = n - X$

$$\therefore r(X, Y) = r(X, n - X) = r(X, -X) = -r(X, X) = -1.$$

(b) Two dice are thrown, their scores being  $a$  and  $b$ . The first die is left on the table while the second is picked up and thrown again giving the score  $c$ . Suppose the process is repeated a large number of times. What is the correlation coefficient between  $X = a + b$  and  $Y = a + c$ ?

$$\text{Ans. } r(X, Y) = \frac{1}{2}$$

20. (a) If  $X$  and  $Y$  are independent random variables with means  $\mu_1$  and  $\mu_2$  and variances  $\sigma_1^2, \sigma_2^2$  respectively, show that the correlation coefficient between  $U = X$  and  $V = X - Y$  in terms of  $\mu_1, \mu_2, \sigma_1^2$  and  $\sigma_2^2$  is  $\sigma_1 / \sqrt{\sigma_1^2 + \sigma_2^2}$ .

(b) If  $X$  and  $Y$  are independent random variables with non-zero variances, show that the correlation coefficient between  $U = XY$  and  $V = X$  in terms of mean and variance of  $X$  and  $Y$  is given by

$$\mu_2\sigma_1/\sqrt{\sigma_1^2\sigma_2^2 + \mu_1^2\sigma_2^2 + \mu_2^2\sigma_1^2}$$

[Delhi Univ. B.Sc. (Stat Hons.), 1987]

21. If  $X_i$ ,  $Y_j$  and  $Z_k$  are all independent random variables with mean zero and unit variance, find the correlation coefficient between

$$U = \sum_{i=1}^m X_i + \sum_{j=1}^n Y_j \text{ and } V = \sum_{i=1}^m X_i + \sum_{k=1}^n Z_k$$

Ans.  $r(U, V) = m/(m + n)$  (Bombay Univ., B.Sc, 1990)

22. (a) Find the value of  $l$  so that the correlation coefficient between  $(X - lY)$  and  $(X + Y)$  is maximum, where  $X, Y$  are independent random variables each with mean zero and variance 1. [Ans.  $l = -1$ ]

Hint.  $U = X - lY$ ;  $V = X + Y$ . Now find  $l$  so that  $r(U, V) = 1$ .

(b) If  $U = X + kY$  and  $V = X + mY$  and  $r$  is the correlation coefficient between  $X$  and  $Y$ , find the correlation coefficient between  $U$  and  $V$ . Show that  $U$  and  $V$  are uncorrelated if

$$k = \frac{-\sigma_X(\sigma_X + rm\sigma_Y)}{\sigma_Y(r\sigma_X + m\sigma_Y)}$$

and further if  $m = \frac{\sigma_X}{\sigma_Y}$ , then  $k = -\frac{\sigma_X}{\sigma_Y}$ . (Gujarat Univ. M.A., 1993)

23.  $X_1, X_2, X_3$  are three variables, each with variance  $\sigma^2$  and the correlation coefficient between any two of them is  $r$ . If  $\bar{X} = (X_1 + X_2 + X_3)/3$ , show that

$$\text{Var}(\bar{X}) = \frac{\sigma^2}{3}(1 + 2r)$$

Deduce that  $r \geq -1/2$ .

24. (a) If  $U = aX + bY$  and  $V = bX - aY$ , show that  $U$  and  $V$  are uncorrelated if

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

where  $\rho$  is the coefficient of correlation between  $X$  and  $Y$ . Show further that, in this case

$$\sigma_U^2 + \sigma_V^2 = (a^2 + b^2)(\sigma_X^2 + \sigma_Y^2) \text{ and } \sigma_U\sigma_V = (a^2 + b^2)\sigma_X\sigma_Y\sqrt{1 - \rho^2}$$

(b) If  $u = aX + bY$ ,  $v = cX + dY$ , show that

$$\begin{vmatrix} \text{var}(u) & \text{cov}(u, v) \\ \text{cov}(u, v) & \text{var}(v) \end{vmatrix} = \begin{vmatrix} a & b \\ c & d \end{vmatrix}^2 \begin{vmatrix} \text{var}(X) & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \text{var}(Y) \end{vmatrix}$$

25. If  $X$  is a standard normal variate and  $Y = a + bX + cX^2$ ,

where  $a, b, c$  are constants, find the correlation coefficient between  $X$  and  $Y$ . Hence or otherwise obtain the conditions when (i)  $X$  and  $Y$  are uncorrelated and (ii)  $X$  and  $Y$  are perfectly correlated.

26. (a) If  $X \sim N(0, 1)$ , find  $\text{corr}(X, Y)$  where  $Y = a + bX + cX^2$ .

[Delhi Univ. B.Sc. (Maths. Hons.), 1985]

**Ans.**  $r(X, Y) = \frac{b}{\sqrt{b^2 + 2c^2}}$

- (b) If  $X$  has Laplace distribution with parameters  $(\lambda, 0)$  and  $Y = a + bX + cX^2$ , find  $\rho(X, Y)$

[*Delhi Univ. B.A. (Stat. Hons. Spl. Course), 1989*]

**Hint.**  $p(x) = \frac{1}{2} \lambda \exp[-\lambda |x|], -\infty < x < \infty.$

$$E(X^{2k+1}) = 0 = \mu_{2k+1}; E(X^{2k}) = \mu_{2k} = (2k)! / \lambda^{2k}$$

$$\rho_{XY} = \frac{\lambda b}{\sqrt{b^2 \lambda^2 + 10c^2}}$$

27. In a sample of  $n$  random observations from exponential distribution with parameter  $\lambda$ , the number of observations in  $(0, 1/\lambda)$  and  $(1/\lambda, 2/\lambda)$ , denoted by  $X$  and  $Y$  are noted. Find  $\rho(X, Y)$ .

**Hint.**  $p_1 = p(0 < X < 1/\lambda) = \int_0^{1/\lambda} \lambda e^{-\lambda x} dx = \frac{e - 1}{e}$

$$p_2 = p(1/\lambda < Y < 2/\lambda) = \int_{1/\lambda}^{2/\lambda} \lambda e^{-\lambda y} dy = \frac{e - 1}{e^2}$$

Then  $(X, Y)$  has a trinomial distribution with parameters  $(n = 3, p_1, p_2, p_3 = 1 - p_1 - p_2)$ .

Hence we have

$$\rho(X, Y) = - \left[ \frac{p_1 p_2}{(1 - p_1)(1 - p_2)} \right]^{1/2} = - \frac{e - 1}{\sqrt{e^2 - e + 1}}.$$

28. Prove that :

$$r(X, Y + Z) = \frac{\sigma_Y}{\sigma_{Y+Z}} \cdot r(X, Y) + \frac{\sigma_Z}{\sigma_{Y+Z}} \cdot r(X, Z)$$

29. If  $X$  and  $Y$  are independent random variables, find  $\text{Corr}(X, XY)$ . Deduce the value of  $\text{Corr}(X, X/Y)$ .

**Ans.**  $r(X, XY) = \sigma_X \mu_Y / [\sigma_X^2 \sigma_Y^2 + \mu_X^2 \sigma_Y^2 + \mu_Y^2 \sigma_X^2]^{1/2}$

30. Prove or Disprove :

(a)  $r(X, Y) = 0 \Rightarrow r(|X|, Y) = 0$

(b)  $r(X, Y) = 0, r(Y, Z) = 0 \Rightarrow r(X, Z) = 0$ .

**Ans.** (a) False, unless  $X$  and  $Y$  are independent.

(b) Hint. Let  $Z = X$ , and  $X$  and  $Y$  be independent. Then

$$r(X, Y) = 0 = r(Y, Z). \text{ But } r(X, Z) = r(X, X) = 1.$$

31. Let random variable  $X$  have a p.d.f.  $f(\cdot)$  with distribution function  $F(\cdot)$ , mean  $\mu$  and variance  $\sigma^2$ . Define  $Y = \alpha + \beta X$ , where  $\alpha$  and  $\beta$  are constants satisfying  $-\infty < \alpha < \infty$ , and  $\beta > 0$ .

(a) Select  $\alpha$  and  $\beta$  so that  $Y$  has mean 0 and variance 1.

(b) What is the correlation coefficient between  $X$  and  $Y$ ?

32. Let  $(X, Y)$  be jointly discrete random variables such that each  $X$  and  $Y$  have at most two mass points. Prove or disprove :  $X$  and  $Y$  are independent if and only if they are uncorrelated.

Ans. True.

33. If the variables  $X_1, X_2, \dots, X_{2n}$  all have the same variance  $\sigma^2$  and the correlation coefficient between  $X_i$  and  $X_j$  ( $i \neq j$ ) has the same value, show that the correlation between  $\sum_{i=1}^n X_i$  and  $\sum_{j=n+1}^{2n} X_j$  is given by  $[n\rho/(1 + (n - 1)\rho)]$ .

34. The means of independent r.v's  $X_1, X_2, \dots, X_n$  are zero and variances are equal, say unity. The correlation coefficients between the sum of selected  $t$  ( $< n$ ) variables out of these variables and the sum of all  $n$  variables are found out. Prove that the sum of squares of all these correlation coefficients is  $n^{-1}C_{t-1}$ .

[Burdwan Univ. B.Sc. (Hons.), 1989]

35. Two variables  $U$  and  $V$  are made up of the sum of a number of terms as follows :

$$U = X_1 + X_2 + \dots + X_a + Y_1 + Y_2 + \dots + Y_b,$$

$$V = X_1 + X_2 + \dots + X_a + Z_1 + Z_2 + \dots + Z_b,$$

where  $a$  and  $b$  are all suffixes and where  $X$ 's,  $Y$ 's and  $Z$ 's are all uncorrelated standardised random variables. Show that the correlation coefficient between

$U$  and  $V$  is  $\frac{n}{\sqrt{(n+a)(n+b)}}$ . Show further that

$$\begin{aligned} \xi &= \sqrt{(n+b)} U + \sqrt{(n+a)} V \\ \eta &= \sqrt{(n+b)} U - \sqrt{(n+a)} V \end{aligned} \quad \dots (*)$$

are uncorrelated

[South Gujarat Univ. B.Sc., 1989]

36. (a) Let the random variables  $X$  and  $Y$  have the joint p.d.f.

$$f(x, y) = 1/3 ; (x, y) = (0, 0), (1, 1) (2, 0)$$

Compute  $E(X)$ ,  $V(X)$ ,  $E(Y)$ ,  $V(Y)$  and  $r(X, Y)$ . Are  $X$  and  $Y$  stochastically independent ? Give reasons.

(b) Let  $(X, Y)$  have the probability distribution :

$$f(0, 0) = 0.45, f(0, 1) = 0.05, f(1, 0) = 0.35, f(1, 1) = 0.15.$$

Evaluate  $V(X)$ ,  $V(Y)$  and  $r(X, Y)$ .

Show that while  $X$  and  $Y$  are correlated,  $X$  and  $X - 5Y$  are uncorrelated. Are  $X$  and  $X - 5Y$  independent ?

(c) Given the bivariate probability distribution :

$$f(-1, 0) = 1/15, \quad f(-1, 1) = 3/15, \quad f(-1, 2) = 2/15$$

$$f(0, 0) = 2/15, \quad f(0, 1) = 2/15, \quad f(0, 2) = 1/15$$

$$f(1, 0) = 1/15, \quad f(1, 1) = 1/15, \quad f(1, 2) = 2/15$$

$$f(x, y) = 0, \text{ elsewhere.}$$

Obtain :

(i) The marginal distributions of  $X$  and  $Y$ .

- (ii) The conditional distributions of  $Y$  given  $X = 0$ .
- (iii)  $E(Y|X = 0)$ .
- (iv) The product moment correlation coefficient between  $X$  and  $Y$ .  
Are  $X$  and  $Y$  independently distributed?

37. If  $X$  and  $Y$  are standardised variates with correlation coefficient  $\rho$ , prove that  $E[\max(X^2, Y^2)] \leq 1 + \sqrt{1 - \rho^2}$

**Hint.**  $\max(X^2, Y^2) = \frac{1}{2}|X^2 - Y^2| + \frac{1}{2}(X^2 + Y^2)$  ...(\*)

$$E(X) = E(Y) = 0; E(X^2) = E(Y^2) = 1; E(XY) = \rho$$

and  $[E|X - Y|, |X + Y|]^2 \leq E(X - Y)^2 \cdot E(X + Y)^2$

(By Cauchy-Schwartz Inequality)

38. The joint p.d.f. of two variates  $X$  and  $Y$  is given by

$$f(x, y) = k[(x+y) - (x^2 + y^2)]; 0 < (x, y) < 1$$

= 0, otherwise.

Show that  $X$  and  $Y$  are uncorrelated but not independent.

39(a). If the random variables  $X$  and  $Y$  have the joint p.d.f.,

$$f(x, y) = \begin{cases} x + y; & 0 < x < 1, 0 < y < 1 \\ 0, & \text{elsewhere} \end{cases}$$

then show that the correlation coefficient between  $X$  and  $Y$  is  $-\frac{1}{11}$ .

[Madras Univ. B.Sc., Oct., 1990]

(b) The density function  $f$  of a random variable  $X$  is given by

$$f(x) = \begin{cases} kx^2, & \text{if } -1 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

- (i) What is the value of  $k$ ? What is the distribution function of  $X$ ?
- (ii) Obtain the density function of the random variable  $Y = X^2$ .
- (iii) Obtain the correlation coefficient between  $X$  and  $Y$ .
- (iv) Are  $X$  and  $Y$  independently distributed?

40(a). If  $f(x, y) = \frac{6-x-y}{8}; 0 \leq x \leq 2, 2 \leq y \leq 4$ ,

find (i)  $\text{Var}(X)$ , (ii)  $\text{Var}(Y)$  (iii)  $r(X, Y)$ .

**Ans.** (i)  $\frac{11}{36}$ , (ii)  $\frac{11}{36}$ , (iii)  $-\frac{1}{11}$ .

(b) Given the joint density of random variables  $X, Y, Z$  as :

$$f(x, y, z) = k x \exp[-(y+z)], 0 < x < 2, y \geq 0, z \geq 0$$

= 0, elsewhere

Find

- (i)  $k$ ,
- (ii) the marginal density function,
- (iii) conditional expectation of  $Y$ , given  $X$  and  $Z$ , and
- (iv) the product moment correlation between  $X$  and  $Y$ .

[Madras Univ. B.Sc. (Main Stat.), 1988]

(c) Suppose that the two dimensional random variable  $(X, Y)$  has p.d.f. given by  $f(x, y) = ke^{-y}$ ,  $0 < x < y < 1$   
 $= 0$ , elsewhere

Find the correlation coefficient  $r_{XY}$ . [Delhi Univ. M.C.A., 1991]

41. The joint density of  $(X, Y)$  is :

$$f(x, y) = \frac{1}{8}(x + y), \quad 0 \leq x \leq 2, 0 \leq y \leq 2.$$

Find  $\mu'_{rs} = E(X^r Y^s)$  and hence find  $\text{Corr}(X, Y)$ .

$$\text{Ans. } \mu'_{rs} = 2^{r+s} \left[ \frac{1}{(r+2)(s+1)} + \frac{1}{(r+1)(s+2)} \right]; r = -\frac{1}{11}.$$

(b) Find the m.g.f. of the bivariate distribution :

$$f(x, y) = 1, \quad 0 < (x, y) < 1 \\ = 0, \text{ otherwise}$$

and hence find  $r(X, Y)$ .

$$\text{Ans. } M(t_1, t_2) = (e^{t_1} - 1)(e^{t_2} - 1)/(t_1 t_2); t_1 \neq 0, t_2 \neq 0, r(X, Y) = 0.$$

42. Let  $(X, Y)$  have joint density :

$$f(x, y) = e^{-(x+y)} I_{(0, \infty)}(x) \cdot I_{(0, \infty)}(y)$$

Find  $\text{Corr}(X, Y)$ . Are  $X$  and  $Y$  independent?

Ans.  $\text{Corr}(X, Y) = 0$ :  $X$  and  $Y$  are independent.

43. A bivariate distribution in two discrete random variables  $X$  and  $Y$  is defined by the probability generating function :

$$\exp[a(u-1) + b(v-1) + c(u-1)(v-1)],$$

simultaneous probability of  $X = r \cap Y = s$ , where  $r$  and  $s$  are integers being the coefficient of  $u^r v^s$ . Find the correlation coefficient between  $X$  and  $Y$ .

Hint. Put  $u = e^{t_1}$  and  $v = e^{t_2}$  in  $\exp[a(u-1) + b(v-1) + c(u-1)(v-1)]$ , the result will be the m.g.f. of a bivariate distribution and is given by

$$M(t_1, t_2) = \exp[a(e^{t_1} - 1) + b(e^{t_2} - 1) + c(e^{t_1} - 1)(e^{t_2} - 1)]$$

$$\text{We have } \left[ \frac{\partial M}{\partial t_1} \right]_{t_1=0, t_2=0} = a, \quad \left[ \frac{\partial^2 M}{\partial t_1^2} \right]_{t_1=0, t_2=0} = a(a+1).$$

$$\left[ \frac{\partial^2 M}{\partial t_1 \partial t_2} \right]_{t_1=0, t_2=0} = ab + c, \quad \left[ \frac{\partial M}{\partial t_2} \right]_{t_1=0, t_2=0} = b, \quad \left[ \frac{\partial^2 M}{\partial t_2^2} \right]_{t_1=0, t_2=0} = b(b+1)$$

So we have

$$E(X) = a, E(X^2) = a(a+1), E(Y) = b, E(Y^2) = b(b+1) \text{ and } E(XY) = ab + c$$

$$\therefore r(X, Y) = \frac{E(XY) - E(X)E(Y)}{\sqrt{[E(X^2) - \{E(X)\}^2][E(Y^2) - \{E(Y)\}^2]}} = \frac{c}{\sqrt{ab}}$$

44. Let the number  $X$  be chosen at random from among the integers 1, 2, 3, 4 and the number  $Y$  be chosen from among those at least as large as  $X$ . Prove that  $\text{Cov}(X, Y) = 5/8$ . Find also the regression line of  $Y$  on  $X$ .

[Delhi Univ. B.Sc. (Maths. Hons.), 1990]

Hint.  $P(X = k) = \frac{1}{4}; k = 1, 2, 3, 4$  and  $Y \geq X$ .

$$\hat{P}(Y=y \mid X=1) = \frac{1}{4}; y = 1, 2, 3, 4 (\because y \geq x);$$

$$P(Y=y \mid X=2) = \frac{1}{3}, y = 2, 3, 4$$

$$P(Y=y \mid X=3) = \frac{1}{2}, y = 3, 4; P(Y=y \mid X=4) = 1, y = 4.$$

The joint probability distribution can be obtained on using :

$$P(X=x, Y=y) = P(X=x) \cdot P(Y=y \mid X=x).$$

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{5/8}{\sqrt{(5/4) \times (41/48)}} = \sqrt{\frac{15}{41}}$$

$$\text{Regression line of } Y \text{ on } X : Y - E(Y) = \frac{r \sigma_Y}{\sigma_X} [X - E(X)]$$

**45.** Two ideal dice are thrown. Let  $X_1$  be the score on the first dice and  $X_2$ , the score on the second dice. Let  $Y = \max\{X_1, X_2\}$ . Obtain the joint distribution of  $Y$  and  $X_1$  and show that

$$\text{Corr}((Y, X_1)) = \frac{3}{2 \sqrt{73}}$$

**46.** Consider an experiment of tossing two tetrahedra. Let  $X$  be the number of the down turned face of first tetrahedron and  $Y$ , the larger of the two numbers. Obtain the joint distribution of  $X$  and  $Y$  and hence  $\rho(X, Y)$ .

$$\text{Ans. } \rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{5/8}{\sqrt{5/4} \cdot \sqrt{55/64}} = \frac{2}{\sqrt{11}}$$

**47.** Three fair coins are tossed. Let  $X$  denote the number of heads on the first two coins and let  $Y$  denote the number of tails on the last two coins.

(a) Find the joint distribution of  $X$  and  $Y$ .

(b) Find the conditional distribution of  $Y$  given that  $X = 1$ .

(c) Find  $\text{Cov}(X, Y)$ .

$$\text{Ans. Cov}(X, Y) = -1/4.$$

**48.** For the trinomial distribution of two random variables  $X$  and  $Y$ :

$$f(x, y) = \frac{n!}{x! y! (n-x-y)!} p^x q^y (1-p-q)^{n-x-y}$$

for  $x, y = 0, 1, 2, \dots, n$  and  $x+y \leq n$ ,  $p \geq 0$ ,  $q \geq 0$  and  $p+q \leq 1$ .

(a) Obtain the marginal distribution of  $Y$ .

(b) Obtain  $E(X \mid Y=y)$ .

(c) Find  $\rho(X, Y)$ .

$$\text{Ans. (a) } X \sim B(n, p), Y \sim B(n, q)$$

$$(b) (X \mid Y=y) \sim B\left(n-y, \frac{p}{1-q}\right)$$

(Note :  $p+q \neq 1$ )

$$\therefore E(X \mid Y=y) = (n-y) \left( \frac{p}{1-q} \right)$$

$$(c) \text{ Cov}(X, Y) = -npq; \rho(X, Y) = -\left[ \frac{pq}{(1-p)(1-q)} \right]^{1/2}$$

### OBJECTIVE TYPE QUESTIONS

I. Comment on the following :

- (i)  $r_{XY} = 0 \Rightarrow X$  and  $Y$  are independent.
- (ii) If  $r_{XY} > 0$  then  $r_{X^2, Y} > 0$ ,  $r_{X^2, Y^2} > 0$  and  $r_{X^2, -Y} > 0$
- (iii)  $r_{XY} > 0 \Rightarrow E(XY) > E(X)E(Y)$
- (iv) Pearson's coefficient of correlation is independent of origin but not of scale.
- (v) The numerical value of product moment correlation coefficient ' $r$ ' between two variables  $X$  and  $Y$  cannot exceed unity.
- (vi) If the correlation coefficient between the variables  $X$  and  $Y$  is zero then the correlation coefficient between  $X^2$  and  $Y^2$  is also zero.
- (vii) If  $r > 0$ , then as  $X$  increases,  $Y$  also increases.
- (viii) "The closeness of relationship between two variables is proportional to  $r$ ."
- (ix)  $r$  measures every type of relationship between the two variables.

II. Comment on the following values of ' $r$ ' (correlation coefficient) :

1, -0.95, 0, -1.64, 0.87, 0.32, -1, 2.4.

III. (i) If  $\rho_{XY} = -0.9$ , then for large values of  $X$ , what sort of values do we expect for  $Y$ ?

(ii) If  $\rho_{XY} = 0$ , what is the value of  $\text{cov}(X, Y)$  and how are  $X$  and  $Y$  related?

IV. Indicate the correct answer :

- (i) The coefficient of correlation will have positive sign when
  - (a)  $X$  is increasing,  $Y$  is decreasing, (b) both  $X$  and  $Y$  are increasing,
  - (c)  $X$  is decreasing,  $Y$  is increasing, (d) there is no change in  $X$  and  $Y$ .
- (ii) The coefficient of correlation (a) can take any value between -1 and +1
  - (b) is always less than -1, (c) is always more than +1, (d) cannot be zero.
- (iii) The coefficient of correlation (a) cannot be positive, (b) cannot be negative, (c) is always positive, (d) can be both positive as well as negative.
- (iv) Probable error of  $r$  is

$$(a) 0.6475 \frac{1-r^2}{\sqrt{n}}, (b) 0.6754 \frac{1+r^2}{\sqrt{n}}, (c) 0.6547 \frac{1-r^2}{n},$$

$$(d) 0.6754 \frac{1-r^2}{n}.$$

(v) The coefficient of correlation between  $X$  and  $Y$  is 0.6. Their covariance is 4.8. The variance of  $X$  is 9. Then the S.D. of  $Y$  is

$$(a) \frac{4.8}{3 \times 0.6}, (b) \frac{0.6}{4.8 \times 3}, (c) \frac{3}{4.8 \times 0.6}, (d) \frac{4.8}{9 \times 0.6}.$$

Probable error also enables us to find the limits within which the population correlation coefficient can be expected to vary. The limits are  $r \pm P.E.(r)$ .

**10.6. Rank Correlation.** Let us suppose that a group of  $n$  individuals is arranged in order of merit or proficiency in possession of two characteristics  $A$  and  $B$ . These ranks in the two characteristics will, in general, be different. For example, if we consider the relation between intelligence and beauty, it is not necessary that a beautiful individual is intelligent also. Let  $(x_i, y_i); i = 1, 2, \dots, n$  be the ranks of the  $i$ th individual in two characteristics  $A$  and  $B$  respectively. Pearsonian coefficient of correlation between the ranks  $x_i$ 's and  $y_i$ 's is called the rank correlation coefficient between  $A$  and  $B$  for that group of individuals.

Assuming that no two individuals are bracketed equal in either classification, each of the variables  $X$  and  $Y$  takes the values  $1, 2, \dots, n$ .

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

$$\begin{aligned}\sigma_x^2 &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 = \frac{1}{n} (1^2 + 2^2 + \dots + n^2) - \left(\frac{n+1}{2}\right)^2 \\ &= \frac{n(n+1)(2n+1)}{6n} - \left(\frac{n+1}{2}\right)^2 = \frac{n^2-1}{12}\end{aligned}$$

$$\therefore \sigma_x^2 = \frac{n^2-1}{12} = \sigma_y^2$$

In general  $x_i \neq y_i$ . Let  $d_i = x_i - y_i$

$$\therefore d_i = (x_i - \bar{x}) - (y_i - \bar{y}) \quad (\because \bar{x} = \bar{y})$$

Squaring and summing over  $i$  from 1 to  $n$ , we get

$$\begin{aligned}\sum d_i^2 &= \sum ((x_i - \bar{x}) - (y_i - \bar{y}))^2 \\ &= \sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2 - 2\sum (x_i - \bar{x})(y_i - \bar{y})\end{aligned}$$

Dividing both sides by  $n$ , we get

$$\frac{1}{n} \sum d_i^2 = \sigma_x^2 + \sigma_y^2 - 2 \operatorname{Cov}(X, Y) = \sigma_x^2 + \sigma_y^2 - 2\rho \sigma_x \sigma_y,$$

where  $\rho$  is the rank correlation coefficient between  $A$  and  $B$ .

$$\begin{aligned}\therefore \frac{1}{n} \sum d_i^2 &= 2\sigma_x^2 - 2\rho \sigma_x^2 \Rightarrow 1 - \rho = \frac{\sum d_i^2}{2n\sigma_x^2} \\ \Rightarrow \rho &= 1 - \frac{\sum d_i^2}{2n\sigma_x^2} = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2-1)} \quad \dots(10.7)\end{aligned}$$

which is the *Spearman's formula for the rank correlation coefficient*.

**Remark.** We always have

$$\sum d_i = \sum (x_i - y_i) = \sum x_i - \sum y_i = n(\bar{x} - \bar{y}) = 0 \quad (\because \bar{x} = \bar{y})$$

This serves as a check on the calculations.

**10-6-1. Tied Ranks.** If some of the individuals receive the same rank in a ranking of merit, they are said to be tied. Let us suppose that  $m$  of the individuals, say,  $(k+1)^{th}$ ,  $(k+2)^{th}$ , ...,  $(k+m)^{th}$  are tied. Then each of these  $m$  individuals is assigned a common rank, which is the arithmetic mean of the ranks  $k+1, k+2, \dots, k+m$ .

*Derivation of  $\rho(X, Y)$ :* We have :

$$\rho(X, Y) = \frac{\sum (X - \bar{X})(Y - \bar{Y})}{[\sum (X - \bar{X})^2 \cdot \sum (Y - \bar{Y})^2]^{1/2}} = \frac{\sum xy}{\sqrt{\sum x^2 \cdot \sum y^2}} \quad \dots (*)$$

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

If  $X$  and  $Y$  each takes the values  $1, 2, \dots, n$ , then we have

$$\bar{X} = (n+1)/2 = \bar{Y}$$

$$\text{and } n\sigma_X^2 = \sum x^2 = \frac{n(n^2-1)}{12} \text{ and } n\sigma_Y^2 = \sum y^2 = \frac{n(n^2-1)}{12} \quad \dots (**)$$

$$\text{Also } \sum d^2 = \sum (X - Y)^2 = \sum [(X - \bar{X}) - (Y - \bar{Y})]^2 = \sum (x - y)^2$$

$$\Rightarrow \sum d^2 = \sum x^2 + \sum y^2 - 2\sum xy$$

$$\Rightarrow \sum xy = \frac{1}{2} [\sum x^2 + \sum y^2 - \sum d^2] \quad \dots (***)$$

We shall now investigate the effect of common ranking, (in case of ties), on the sum of squares of the ranks. Let  $S^2$  and  $S_1^2$  denote the sum of the squares of untied and tied ranks respectively.

Then we have :

$$\begin{aligned} S^2 &= (k+1)^2 + (k+2)^2 + \dots + (k+m)^2 \\ &= mk^2 + (1^2 + 2^2 + \dots + m^2) + 2k(1+2+\dots+m) \\ &= mk^2 + \frac{m(m+1)(2m+1)}{6} + mk(m+1) \end{aligned}$$

$$\begin{aligned} S_1^2 &= m(\text{Average rank})^2 \\ &= m \left[ \frac{(k+1) + (k+2) + \dots + (k+m)}{m} \right]^2 \\ &= m \left( k + \frac{m+1}{2} \right)^2 = mk^2 + \frac{m(m+1)^2}{4} + m k (m+1) \end{aligned}$$

$$\therefore S^2 - S_1^2 = \frac{m(m+1)}{12} [2(2m+1) - 3(m+1)] = \frac{m(m^2-1)}{12}$$

Thus the effect of tying  $m$  individuals (ranks) is to reduce the sum of the squares by  $m(m^2-1)/12$ , though the mean value of the ranks remains the same, viz.,  $(n+1)/2$ .

Suppose that there are  $s$  such sets of ranks to be tied in the  $X$ -series so that the total sum of squares due to them is

$$\frac{1}{12} \sum_{i=1}^s m_i (m_i^2 - 1) = \frac{1}{12} \sum_{i=1}^s (m_i^3 - m_i) = T_X, \text{ (say)} \quad \dots (10.7a)$$

Similarly suppose that there are  $t$  such sets of ranks to be tied with respect to the other series  $Y$  so that sum of squares due to them is :

$$\frac{1}{12} \sum_{j=1}^t m_j' \cdot (m_j'^2 - 1) = \frac{1}{12} \sum_{j=1}^t (m_j'^3 - m_j') = T_Y, \text{ (say)} \quad \dots(10.7b)$$

Thus, in the case of ties, the new sums of squares are given by :

$$n \operatorname{Var}'(X) = \sum x^2 - T_X = \frac{n(n^2 - 1)}{12} - T_X$$

$$n \operatorname{Var}'(Y) = \sum y^2 - T_Y = \frac{n(n^2 - 1)}{12} - T_Y$$

$$\text{and } n \operatorname{Cov}'(X, Y) = \frac{1}{2} [\sum x^2 - T_X + \sum y^2 - T_Y - \sum d^2] \quad [\text{From } (***)]$$

$$= \frac{1}{2} \left[ \frac{n(n^2 - 1)}{12} - T_X + \frac{n(n^2 - 1)}{12} - T_Y - \sum d^2 \right]$$

$$= \frac{n(n^2 - 1)}{12} - \frac{1}{2} [(T_X + T_Y) + \sum d^2]$$

$$\rho(X, Y) = \frac{\frac{n(n^2 - 1)}{12} - \frac{1}{2} [T_X + T_Y + \sum d^2]}{\left[ \frac{n(n^2 - 1)}{12} - T_X \right]^{1/2} \left[ \frac{n(n^2 - 1)}{12} - T_Y \right]^{1/2}}$$

$$= \frac{\frac{n(n^2 - 1)}{6} - [\sum d^2 + T_X + T_Y]}{\left[ \frac{n(n^2 - 1)}{6} - 2T_X \right]^{1/2} \left[ \frac{n(n^2 - 1)}{6} - 2T_Y \right]^{1/2}}$$

...(10.7c)

where  $T_X$  and  $T_Y$  are given by (10.7a) and (10.7b).

**Remark.** If we adjust only the covariance term i.e.,  $\sum xy$  and not the variances  $\sigma_X^2$  (or  $\sum x^2$ ) and  $\sigma_Y^2$  (or  $\sum y^2$ ) for ties, then the formula (10.7c) reduces to :

$$\begin{aligned} \rho(X, Y) &= \frac{\frac{n(n^2 - 1)}{6} - (\sum d^2 + T_X + T_Y)}{n(n^2 - 1)/6} \\ &= 1 - \frac{6 [\sum d^2 + T_X + T_Y]}{n(n^2 - 1)}, \end{aligned} \quad \dots(10.7d)$$

a formula which is commonly used in practice for numerical problems. For illustration, see Example 10-18.

**Example 10-16.** The ranks of same 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 proficiencies of this group in Mathematics and Physics.

**Solution.**

Ranks in Maths. (X)	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	Total
Ranks in Physics (Y)	1	10	3	4	5	7	2	6	8	11	15	9	14	12	16	13	
$d = X - Y$	0	-8	0	0	0	-1	5	2	1	-1	-4	3	-1	2	-1	3	0
$d^2$	0	64	0	0	0	1	25	4	1	1	16	9	1	4	1	9	136

Rank correlation coefficient is given by

$$\rho = 1 - \frac{6 \sum d^2}{n(n^2 - 1)} = 1 - \frac{6 \times 136}{16 \times 255} = 1 - \frac{1}{5} = \frac{4}{5} = 0.8$$

**Example 10-17.** Ten competitors in a musical test were ranked by the three judges A, B and C in the following order :

Ranks by A : 1 6 5 10 3 2 4 9 7 8

Ranks by B : 3 5 8 4 7 10 2 1 6 9

Ranks by C : 6 4 9 8 1 2 3 10 5 7

Using rank correlation method, discuss which pair of judges has the nearest approach to common likings in music.

**Solution.** Here  $n = 10$

Ranks by A (X)	Ranks by B (Y)	Ranks by C (Z)	$d_1 = X - Y$	$d_2 = X - Z$	$d_3 = Y - Z$	$d_1^2$	$d_2^2$	$d_3^2$
1	3	6	-2	-5	-3	4	25	9
6	5	4	1	2	1	1	4	1
5	8	9	-3	-4	-1	9	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
Total			$\sum d_1 = 0$	$\sum d_2 = 0$	$\sum d_3 = 0$	$\sum d_1^2 = 200$	$\sum d_2^2 = 60$	$\sum d_3^2 = 214$

$$\rho(X, Y) = 1 - \frac{6 \sum d_1^2}{n(n^2 - 1)} = 1 - \frac{6 \times 200}{10 \times 99} = 1 - \frac{40}{33} = -\frac{7}{33}$$

$$\rho(X, Z) = 1 - \frac{6 \sum d_2^2}{n(n^2 - 1)} = 1 - \frac{6 \times 60}{10 \times 99} = 1 - \frac{4}{11} = \frac{7}{11}$$

$$\rho(Y, Z) = 1 - \frac{6 \sum d_3^2}{n(n^2 - 1)} = 1 - \frac{6 \times 214}{10 \times 99} = -\frac{49}{165}$$

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

**10-6-2. Repeated Ranks (Continued).** If any two or more individuals are bracketed equal in any classification with respect to characteristics A and B, or if there is more than one item with the same value in the series, then the Spearman's formula (10-7) for calculating the rank correlation coefficient breaks down, since in this case each of the variables X and Y does not assume the values 1, 2, ..., n and consequently,  $\bar{x} \neq \bar{y}$ .

In this case, common ranks are given to the repeated items. This common rank is the average of the ranks which these items would have assumed if they were slightly different from each other and the next item will get the rank next to the ranks already assumed. As a result of this, following adjustment or correction is made in the rank correlation formula [c.f. (10-7c) and (10-7d)].

In the formula, we add the factor  $\frac{m(m^2 - 1)}{12}$  to  $\sum d^2$ , where m is the number of times an item is repeated. This correction factor is to be added for each repeated value in both the X-series and Y-series.

**Example 10-18.** Obtain the rank correlation coefficient for the following data :

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

**Solution.**

#### CALCULATIONS FOR RANK CORRELATION

X	Y	Rank X (x)	Rank Y (y)	$d = x - y$	$d^2$
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
$\Sigma d = 0$					$\Sigma d^2 = 72$

In the X-series we see that the value 75 occurs 2 times. The common rank given to these values is 2.5 which is the average of 2 and 3, the ranks which these values would have taken if they were different. The next value 68, then gets the next rank which is 4. Again we see that value 64 occurs thrice. The common rank given to it is 6 which is the average of 5, 6 and 7. Similarly in

the  $Y$ -series, the value 68 occurs twice and its common rank is 3.5 which is the average of 3 and 4. As a result of these common rankings, the formula for ' $\rho$ ' has to be corrected. To  $\sum d^2$  we add  $\frac{m(m^2 - 1)}{12}$  for each value repeated, where  $m$  is the number of times a value occurs. In the  $X$ -series the correction is to be applied twice, once for the value 75 which occurs twice ( $m = 2$ ) and then for the value 64 which occurs thrice ( $m = 3$ ). The total correction for the  $X$ -series is

$$\frac{2(4 - 1)}{12} + \frac{3(9 - 1)}{12} = \frac{5}{2}$$

Similarly, this correction for the  $Y$ -series is  $\frac{2(4 - 1)}{12} = \frac{1}{2}$ , as the value 68 occurs twice.

$$\text{Thus } \rho = 1 - \frac{6 \left[ \sum d^2 + \frac{5}{2} + \frac{1}{2} \right]}{n(n^2 - 1)} = 1 - \frac{6(72 + 3)}{10 \times 99} = 0.545$$

**10.6.3. Limits for the Rank Correlation Coefficient.** Spearman's rank correlation coefficient is given by

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

' $\rho$ ' is maximum, if  $\sum_{i=1}^n d_i^2$  is minimum, i.e., if each of the deviations  $d_i$  is minimum. But the minimum value of  $d_i$  is zero in the particular case  $x_i = y_i$ , i.e., if the ranks of the  $i$ th individual in the two characteristic are equal. Hence the maximum value of  $\rho$  is +1, i.e.,  $\rho \leq 1$ .

' $\rho$ ' is minimum, if  $\sum_{i=1}^n d_i^2$  is maximum, i.e., if each of the deviations  $d_i$  is maximum which is so if the ranks of the  $n$  individuals in the two characteristics are in the opposite directions as given below :

$x$	1	2	3	...	...	$n - 1$	$n$	
$y$	$n$	$n - 1$	$n - 2$	...	...	2	1	...(*)

**Case 1.** Suppose  $n$  is odd and equal to  $(2m + 1)$  then the values of  $d$  are :

$d : 2m, 2m - 2, 2m - 4, \dots, 2, 0, -2, -4, \dots, -(2m - 2), -2m$ .

$$\therefore \sum_{i=1}^n d_i^2 = 2 \{ (2m)^2 + (2m - 2)^2 + \dots + 4^2 + 2^2 \}$$

$$= 8 \{ m^2 + (m - 1)^2 + \dots + 1^2 \} = \frac{8m(m + 1)(2m + 1)}{6}$$

$$\text{Hence } \rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)} = 1 - \frac{8m(m+1)(2m+1)}{(2m+1)\{(2m+1)^2 - 1\}} \\ = \frac{8m(m+1)}{(4m^2 + 4m)} = 1 - \frac{8m(m+1)}{4m(m+1)} = -1$$

**Case II.** Let  $n$  be even and equal to  $2m$ , (say).

Then the values of  $d$  are

$$(2m-1), (2m-3), \dots, 1, -1, -3, \dots, -(2m-3), -(2m-1) \\ \therefore \sum d_i^2 = 2\{(2m-1)^2 + (2m-3)^2 + \dots + 1^2\} \\ = 2[(\{(2m)^2 + (2m-1)^2 + (2m-2)^2 + \dots + 2^2 + 1^2\}) \\ - \{(2m)^2 + (2m-2)^2 + \dots + 4^2 + 2^2\}] \\ = 2[1^2 + 2^2 + \dots + (2m)^2 - \{2^2m^2 + 2^2(m-1)^2 + \dots + 2^2\}] \\ = 2\left[\frac{2m(2m+1)(4m+1)}{6} - 4\frac{m(m+1)(2m+1)}{6}\right] \\ = \frac{2m}{3}[(2m+1)(4m+1) - 2(m+1)(2m+1)] \\ = \frac{2m}{3}[(2m+1)(4m+1) - 2m(m+1) - 2(m+1)] \\ = \frac{2m}{3}(2m+1)(2m-1) = \frac{2m(4m^2-1)}{3} \\ \therefore \rho = 1 - \frac{6\sum d_i^2}{n(n^2 - 1)} = 1 - \frac{4m(4m^2-1)}{2m(4m^2-1)} = -1$$

Thus the limits for rank correlation coefficient are given by  $-1 \leq \rho \leq 1$ .

**Aliter.** For an alternate and simpler proof for obtaining the minimum value of  $\rho$ , from (\*) onward, proceed as in Hint to Question Number 9 of Exercise 10(c).

#### Remarks on Spearman's Rank Correlation Coefficient.

1.  $\sum d = \sum x - \sum y = n(\bar{x} - \bar{y}) = 0$ , which provides a check for numerical calculations.

2. Since Spearman's rank correlation coefficient  $\rho$  is nothing but Pearsonian correlation coefficient between the ranks, it can be interpreted in the same way as the Karl Pearson's correlation coefficient.

3. Karl Pearson's correlation coefficient assume that the parent population from which sample observations are drawn is normal. If this assumption is violated then we need a measure which is distribution free (or non-parametric). A distribution-free measure is one which does not make any assumptions about the parameters of the population. Spearman's  $\rho$  is such a measure (*i.e.*, distribution-free), since no strict assumptions are made about the form of the population from which sample observations are drawn.

4. Spearman's formula is easy to understand and apply as compared with Karl Pearson's formula. The value obtained by the two formulae, *viz.*, Pearsonian  $r$  and Spearman's  $\rho$ , are generally different. The difference arises due to the fact that when ranking is used instead of full set of observations, there is

always some loss of information. Unless many ties exist, the coefficient of rank correlation should be only slightly lower than the Pearsonian coefficient.

5. Spearman's formula is the only formula to be used for finding correlation coefficient if we are dealing with qualitative characteristics which cannot be measured quantitatively but can be arranged serially. It can also be used where actual data are given. In case of extreme observations, Spearman's formula is preferred to Pearson's formula.

6. Spearman's formula has its limitations also. It is not practicable in the case of bivariate frequency distribution (Correlation Table). For  $n > 30$ , this formula should not be used-unless the ranks are given, since in the contrary case the calculations are quite time-consuming.

### EXERCISE 10(c)

1. Prove that Spearman's rank correlation coefficient is given by

$$1 - \frac{6 \sum d_i^2}{n^3 - n}$$
, where  $d_i$  denotes the difference between the ranks of  $i$ th individual.

2. (a) Explain the difference between product moment correlation coefficient and rank correlation coefficient.

- (b) The rankings of ten students in two subjects  $A$  and  $B$  are as follows :

$A$ :	3	5	8	4	7	10	2	1	6	9
$B$ :	6	4	9	8	1	2	3	10	5	7

Find the correlation coefficient.

3. (a) Calculate the coefficient of correlation for ranks from the following data :

$$(X, Y) : (5, 8), (10, 3), (6, 2), (3, 9), (19, 12), (5, 3), (6, 17), (12, 18), (8, 22), (2, 12), (10, 17), (19, 20).$$

[Calicut Univ. B.Sc. (Subs. Stat.), Oct. 1991]

- (b) Ten recruits were subjected to a selection test to ascertain their suitability for a certain course of training. At the end of training they were given a proficiency test.

The marks secured by recruits in the selection test ( $X$ ) and in the proficiency test ( $Y$ ) are given below :—

Serial No. :	1	2	3	4	5	6	7	8	9	10
$X$ :	10	15	12	17	13	16	24	14	22	20
$Y$ :	30	42	45	46	33	34	40	35	39	38

Calculate product moment correlation coefficient and rank correlation coefficient. Why are two coefficients different ?

4. (a) The I.Q.'s of a group of 6 persons were measured, and they then sat for a certain examination. Their I.Q.'s and examination marks were as follows :

Person :	$A$	$B$	$C$	$D$	$E$	$F$
I.Q. :	110	100	140	120	80	90
Exam. marks :	70	60	80	60	10	20

Compute the coefficients of correlation and rank correlation. Why are the correlation figures obtained different ?

Ans. 0.882 and 0.9.

The difference arises due to the fact that when ranking is used instead of the full set of observations, there is always some loss of information.

(b) The value of ordinary correlation ( $r$ ) for the following data is 0.636 :—

$X$  : .05 .14 .24 .30 .47 .52 .57 .61 .67 .72

$Y$  : 1.08 1.15 1.27 1.33 1.41 1.46 1.54 2.72 4.01 9.63

(i) Calculate Spearman's rank-correlation ( $\rho$ ) for this data.

(ii) What advantage of  $\rho$  was brought out in this example ?

4. Ten competitors in a beauty contest are ranked by three judges as follows :

Competitors										
Judges	1	2	3	4	5	6	7	8	9	10
A	6	5	3	10	2	4	9	7	8	1
B	5	8	4	7	10	2	1	6	9	3
C	4	9	8	1	2	3	10	5	7	6

Discuss which pair of judges has the nearest approach to common tastes of beauty.

5. A sample of 12 fathers and their eldest sons gave the following data about their height in inches :

Father : 65 63 67 64 68 62 70 66 68 67 69 71

Son : 68 66 68 65 69 66 68 65 71 67 68 70

Calculate coefficient of rank correlation. (Ans. 0.7220)

6. The coefficient of rank correlation between marks in Statistics and marks in Mathematics obtained by a certain group of students is 0.8. If the sum of the squares of the difference in ranks is given to be 33, find the number of student in the group (Ans. 10). [Madras Univ. B.Sc., 1990]

7. The coefficient of rank correlation of the marks obtained by 10 students in Maths and Statistics was found to be 0.5. It was later discovered that the difference in ranks in two subjects obtained by one of the students was wrongly taken as 3 instead of 7. Find the correct coefficient of rank correlation.

$$\text{Hint.} \quad 0.5 = 1 - \frac{6 \sum d^2}{10 \times 99}$$

$$\Rightarrow \quad \sum d^2 = \frac{990}{6 \times 2} = 82.5$$

Since one difference was wrongly taken as 3 instead of 7, the correct value of  $\sum d^2$  is given by

$$\text{Corrected } \sum d^2 = 82.5 - (3)^2 + (7)^2 = 122.5$$

$$\therefore \quad \text{Corrected } \rho = 1 - \frac{6 \times 122.5}{10 \times 99} = 0.2576$$

8. If  $d_i$  be the difference in the ranks of the  $i$ th individual in two different characteristics, then show that the maximum value of  $\sum_{i=1}^n d_i^2$  is  $\frac{1}{3}(n^3 - n)$ .

Hence or otherwise, show that rank correlation coefficient lies between  $-1$  and  $+1$ .  
[Delhi Univ. B.Sc. (Maths. Hons.), 1986]

9. Let  $x_1, x_2, \dots, x_n$  be the ranks of  $n$  individuals according to a character  $A$  and  $y_1, y_2, \dots, y_n$  be the ranks of the same individuals according to other character  $B$ . Obviously  $(x_1, x_2, \dots, x_n)$  and  $(y_1, y_2, \dots, y_n)$  are permutations of  $1, 2, \dots, n$ . It is given that  $x_i + y_i = 1 + n$ , for  $i = 1, 2, \dots, n$ . Show that the value of the rank correlation coefficient  $\rho$  between the characters  $A$  and  $B$  is  $-1$ .

**Hint.** We are given  $x_i + y_i = n + 1 \forall i = 1, 2, \dots, n$

$$\text{Also } x_i - y_i = d_i$$

$$\therefore 2x_i = n + 1 + d_i \Rightarrow d_i = 2x_i - (n + 1)$$

$$\therefore \sum_{i=1}^n d_i^2 = \sum_{i=1}^n [4x_i^2 + (n+1)^2 - 2(n+1)2x_i]$$

$$= 4 \frac{n(n+1)(2n+1)}{6} + n(n+1)^2 - \frac{4(n+1)n(n+1)}{2}$$

$$= \frac{n(n^2-1)}{3}$$

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

**Remark.** From Spearman's formula we note that  $\rho$  will be minimum if  $\sum d_i^2$  is maximum, which will be so if the ranks  $X$  and  $Y$  are in opposite directions as given below :

$x$	1	2	3	$\dots$	$n$
$y$	$n$	$n-1$	$n-2$	$\dots$	1

This gives us

$$x_i + y_i = n + 1, i = 1, 2, \dots, n.$$

Hence the value of  $\rho = -1$  obtained above is minimum value of  $\rho$ .

10. Show that in a ranked bivariate distribution in which no ties occur and in which the variables are independent

(a)  $\sum_i d_i^2$  is always even, and

(b) there are not more than  $\frac{1}{6}(n^3 - n) + 1$  possible values of  $r$ .

**11.** Show that if  $X, Y$  be identically distributed with common probability mass function :  $P(X = k) = \frac{1}{N}$ , for  $k = 1, 2, \dots, N; N > 1$ ,

then  $\rho_{X,Y}$ , the correlation coefficient between  $X$  and  $Y$ , is given by

$$1 - \frac{6E(X - Y)^2}{N^2 - 1}$$

[Delhi Univ. B.Sc. (Maths Hons.), 1992]

**10.7. Regression.** The term “regression” literally means “stepping back towards the average”. It was first used by a British biometrician Sir Francis Galton (1822—1911), in connection with the inheritance of stature. Galton found that the offsprings of abnormally tall or short parents tend to “regress” or “step back” to the average population height. But the term “regression” as now used in Statistics is only a convenient term without having any reference to biometry.

**Definition.** *Regression analysis is a mathematical measure of the average relationship between two or more variables in terms of the original units of the data.*

In regression analysis there are two types of variables. The variable whose value is influenced or is to be predicted is called *dependent variable* and the variable which influences the values or is used for prediction, is called *independent variable*. In regression analysis independent variable is also known as *regressor or predictor or explanatory variable* while the dependent variable is also known as *regressed or explained variable*.

**10.7.1. Lines of 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*.

The line of regression is the line which gives the best estimate to the value of one variable for any specific value of the other variable. Thus the line of regression is the line of “best fit” and is obtained by the *principles of least squares*.

Let us suppose that in the bivariate distribution  $(x_i, y_i); i = 1, 2, \dots, n$ ;  $Y$  is dependent variable and  $X$  is independent variable. Let the line of regression of  $Y$  on  $X$  be  $Y = a + bX$ .

According to the principle of least squares, the normal equations for estimating  $a$  and  $b$  are (c.f. (9.2a))

$$\sum_{i=1}^n y_i = na + b \sum_{i=1}^n x_i \quad \dots(10-8)$$

and  $\sum_{i=1}^n x_i y_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2 \quad \dots(10-9)$

From (10-8) on dividing by  $n$ , we get

$$\bar{y} = a + b\bar{x} \quad \dots(10-10)$$

Thus the line of regression of  $Y$  on  $X$  passes through the point  $(\bar{x}, \bar{y})$ .

Now

$$\mu_{11} = \text{Cov}(X, Y) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} \Rightarrow \frac{1}{n} \sum_{i=1}^n x_i y_i = \mu_{11} + \bar{x} \bar{y} \quad \dots(10-11)$$

$$\text{Also } \sigma_X^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \Rightarrow \frac{1}{n} \sum_{i=1}^n x_i^2 = \sigma_X^2 + \bar{x}^2 \quad \dots(10-11a)$$

Dividing (10-9) by  $n$  and using (10-11) and (10-11a), we get

$$\mu_{11} + \bar{x} \bar{y} = a\bar{x} + b(\sigma_X^2 + \bar{x}^2) \quad \dots(10-12)$$

Multiplying (10-10) by  $\bar{x}$  and then subtracting from (10-12), we get

$$\mu_{11} = b\sigma_X^2 \Rightarrow b = \frac{\mu_{11}}{\sigma_X^2} \quad \dots(10-13)$$

Since ' $b$ ' is the slope of the line of regression of  $Y$  on  $X$  and since the line of regression passes through the point  $(\bar{x}, \bar{y})$ , its equation is

$$Y - \bar{y} = b(X - \bar{x}) = \frac{\mu_{11}}{\sigma_X^2} (X - \bar{x}) \quad \dots(10-14)$$

$$\Rightarrow Y - \bar{y} = r \frac{\sigma_Y}{\sigma_X} (X - \bar{x}) \quad \dots(10-14a)$$

Starting with the equation  $X = A + BY$  and proceeding similarly or by simply interchanging the variables  $X$  and  $Y$  in (10-14) and (10-14a), the equation of the line of regression of  $X$  on  $Y$  becomes

$$X - \bar{x} = \frac{\mu_{11}}{\sigma_Y^2} (Y - \bar{y}) \quad \dots(10-15)$$

$$\Rightarrow X - \bar{x} = r \frac{\sigma_X}{\sigma_Y} (Y - \bar{y}) \quad \dots(10-15a)$$

Aliter. The straight line defined by

$$Y = a + bX \quad \dots(i)$$

and satisfying the residual (least square) condition

$$S = E [(Y - a - bX)^2] = \text{Minimum} \quad \dots(10-16)$$

for variations in  $a$  and  $b$ , is called the line of regression of  $Y$  on  $X$ .

The necessary and sufficient conditions for a minima of  $S$ , subject to variations in  $a$  and  $b$  are :

$$(i) \frac{\partial S}{\partial a} = 0, \quad \frac{\partial S}{\partial b} = 0 \quad \text{and} \quad \dots (*)$$

$$(ii) \Delta = \begin{vmatrix} \frac{\partial^2 S}{\partial a^2} & \frac{\partial^2 S}{\partial a \partial b} \\ \frac{\partial^2 S}{\partial b \partial a} & \frac{\partial^2 S}{\partial b^2} \end{vmatrix} > 0 \quad \text{and} \quad \frac{\partial^2 S}{\partial a^2} > 0 \quad \dots (**)$$

Using (\*), we get

$$\frac{\partial S}{\partial a} = -2 E [Y - a - bX] = 0 \quad \dots (iii)$$

$$\frac{\partial S}{\partial b} = -2 E [X(Y - a - bX)] = 0 \quad \dots (iv)$$

$$\Rightarrow E(Y) = a + bE(X) \dots (v) \quad \text{and} \quad E(XY) = aE(X) + bE(X^2) \quad \dots (vi)$$

Equation (v) implies that the line (i) of regression of  $Y$  on  $X$  passes through the mean value  $[E(X), E(Y)]$ .

Multiplying (v) by  $E(X)$  and subtracting from (vi), we get

$$E(XY) - E(X)E(Y) = b[E(X^2) - \{E(X)\}^2]$$

$$\Rightarrow \text{Cov}(X, Y) = b \sigma_X^2 \Rightarrow b = \frac{\text{Cov}(X, Y)}{\sigma_X^2} = \frac{r \sigma_Y}{\sigma_X} \quad \dots (vii)$$

Subtracting (v) from (i) and using (vii), we obtain the equation of line of regression of  $Y$  on  $X$  as :

$$Y - E(Y) = \frac{\text{Cov}(X, Y)}{\sigma_X^2} [X - E(X)] \Rightarrow Y - E(Y) = \frac{r \sigma_Y}{\sigma_X} [X - E(X)]$$

Similarly, the straight line defined by  $X = A + BY$   
and satisfying the residual condition

$$E[X - A - BY]^2 = \text{Minimum},$$

is called the line of regression of  $X$  on  $Y$ .

**Remarks 1.** We note that

$$\frac{\partial^2 S}{\partial a^2} = 2 > 0, \text{ and}$$

$$\frac{\partial^2 S}{\partial b^2} = 2E(X^2) \quad \text{and} \quad \frac{\partial^2 S}{\partial a \partial b} = 2E(X)$$

Substituting in (\*\*), we have

$$\begin{aligned} \Delta &= \frac{\partial^2 S}{\partial a^2} \cdot \frac{\partial^2 S}{\partial b^2} - \left( \frac{\partial^2 S}{\partial a \partial b} \right)^2 \\ &= 4 [E(X^2) - \{E(X)\}^2] = 4 \cdot \sigma_X^2 > 0 \end{aligned}$$

Hence the solution of the least square equations (iii) and (iv), in fact, provides a minima of  $S$ .

2. The regression equation (10-14a) implies that the line of regression of  $Y$  on  $X$  passes through the point  $(\bar{x}, \bar{y})$ . Similarly (10-15a) implies that the line of regression of  $X$  on  $Y$  also passes through the point  $(\bar{x}, \bar{y})$ . Hence both the lines of regression pass through the point  $(\bar{x}, \bar{y})$ . In other words, the mean

values ( $\bar{x}, \bar{y}$ ) can be obtained as the point of intersection of the two regression lines.

**3. Why two lines of Regression ?** There are always two lines of regression, one of  $Y$  on  $X$  and the other of  $X$  on  $Y$ . The line of regression of  $Y$  on  $X$  (10.14a) is used to estimate or predict the value of  $Y$  for any given value of  $X$ , i.e., when  $Y$  is a dependent variable and  $X$  is an independent variable. 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 squares. We can also obtain an estimate of  $X$  for any given value of  $Y$  by using equation (10.14a) but the estimate so obtained will not be best since (10.14a) is obtained on minimising the sum of the squares of errors of estimates in  $Y$  and not in  $X$ . Hence to estimate or predict  $X$  for any given value of  $Y$ , we use the regression equation of  $X$  on  $Y$  (10.15a) which is derived on minimising the sum of the squares of errors of estimates in  $X$ . Here  $X$  is a dependent variable and  $Y$  is an independent variable. The two regression equations are not reversible or interchangeable because of the simple reason that the basis and assumptions for deriving these equations are quite different. The regression equation of  $Y$  on  $X$  is obtained on minimising the sum of the squares of the errors parallel to the  $Y$ -axis while the regression equation of  $X$  on  $Y$  is obtained on minimising the sum of squares of the errors parallel to the  $X$ -axis.

In a particular case of perfect correlation, positive or negative, i.e.,  $r \pm 1$ , the equation of line of regression of  $Y$  on  $X$  becomes :

$$\begin{aligned} Y - \bar{y} &= \pm \frac{\sigma_x}{\sigma_y} (X - \bar{x}) \\ \Rightarrow \quad \frac{Y - \bar{y}}{\sigma_y} &= \pm \left( \frac{X - \bar{x}}{\sigma_x} \right) \end{aligned} \quad \dots(10.16)$$

Similarly, the equation of the line of regression of  $X$  on  $Y$  becomes :

$$\begin{aligned} X - \bar{x} &= \pm \frac{\sigma_x}{\sigma_y} (Y - \bar{y}) \\ \Rightarrow \quad \frac{X - \bar{x}}{\sigma_x} &= \pm \left( \frac{Y - \bar{y}}{\sigma_y} \right), \end{aligned}$$

which is same as (10.16).

Hence in case of perfect correlation, ( $r = \pm 1$ ), both the lines of regression coincide. Therefore, in general, we always have two lines of regression except in the particular case of perfect correlation when both the lines coincide and we get only one line.

**10.7.2. Regression Curves.** In modern terminology, the conditional mean  $E(Y | X = x)$  for a continuous distribution is called the regression function of  $Y$  on  $X$  and the graph of this function of  $x$  is known as the regression curve of  $Y$  on  $X$  or sometimes the regression curve for the mean of  $Y$ . Geometrically, the regression function represents the  $y$  co-ordinate of the centre of mass of the bivariate probability mass in the infinitesimal vertical strip bounded by  $x$  and  $x + dx$ .

Similarly, the regression function of  $X$  on  $Y$  is  $E(X|Y=y)$  and the graph of this function of  $y$  is called the regression curve (of the mean) of  $X$  on  $Y$ .

In case a regression curve is a straight line, the corresponding regression is said to be *linear*. If one of the regressions is linear, it does not however follow that the other is also linear. For illustration, See Example 10.21.

**Theorem 10.4.** Let  $(X, Y)$  be a two-dimensional random variable with  $E(X) = \bar{X}$ ,  $E(Y) = \bar{Y}$ ,  $V(X) = \sigma_x^2$ ,  $V(Y) = \sigma_y^2$  and let  $r = r(X, Y)$  be the correlation coefficient between  $X$  and  $Y$ . If the regression of  $Y$  on  $X$  is linear then

$$E(Y|X) = \bar{Y} + r \frac{\sigma_y}{\sigma_x} (X - \bar{X}) \quad \dots(10.16a)$$

Similarly, if the regression of  $X$  on  $Y$  is linear, then

$$E(X|Y) = \bar{X} + r \frac{\sigma_x}{\sigma_y} (Y - \bar{Y}) \quad \dots(10.16b)$$

**Proof.** Let the regression equation of  $Y$  on  $X$  be

$$E(Y|x) = a + bx \quad \dots(1)$$

But by definition,

$$\begin{aligned} E(Y|x) &= \int_{-\infty}^{\infty} y f(y|x) dy = \int_{-\infty}^{\infty} y \frac{f(x,y)}{f_X(x)} dy \\ \therefore \quad \frac{1}{f_X(x)} \int_{-\infty}^{\infty} y f(x,y) dy &= a + bx \end{aligned} \quad \dots(2)$$

Multiplying both sides of (2) by  $f_X(x)$  and integrating w.r.t.  $x$ , we get

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f(x,y) dy dx &= a \int_{-\infty}^{\infty} f_X(x) dx + b \int_{-\infty}^{\infty} x f_X(x) dx \\ \Rightarrow \quad \int_{-\infty}^{\infty} y \left[ \int_{-\infty}^{\infty} f(x,y) dx \right] dy &= a + bE(X) \\ \Rightarrow \quad \int_{-\infty}^{\infty} y f_Y(y) dy &= a + bE(X) \end{aligned}$$

$$\text{i.e., } E(Y) = a + bE(X) \Rightarrow \bar{Y} = a + b\bar{X} \quad \dots(3)$$

Multiplying both sides of (2) by  $x f_X(x)$  and integrating w.r.t.  $x$ , we get

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f(x,y) dy dx &= a \int_{-\infty}^{\infty} x f_X(x) dx + b \int_{-\infty}^{\infty} x^2 f_X(x) dx \\ \Rightarrow \quad E(XY) &= a E(X) + b E(X^2) \\ \Rightarrow \quad \mu_{11} + \bar{X} \bar{Y} &= a\bar{X} + b(\sigma_x^2 + \bar{X}^2) \end{aligned} \quad \dots(4)$$

$$\therefore \mu_{11} = E(XY) - E(X)E(Y) = E(XY) - \bar{X}\bar{Y};$$

$$\sigma_X^2 = E(X^2) - \{E(X)\}^2 = E(X^2) - \bar{X}^2$$

Solving (3) and (4) simultaneously, we get

$$\therefore b = \frac{\mu_{11}}{\sigma_X^2} \text{ and } a = \bar{Y} - \frac{\mu_{11}}{\sigma_X^2} \bar{X}$$

Substituting in (1) and simplifying, we get the required equation of the line of regression of  $Y$  on  $X$  as

$$\begin{aligned} E(Y|x) &= \bar{Y} + \frac{\mu_{11}}{\sigma_X^2} (x - \bar{X}) \\ \Rightarrow E(Y|X) &= \bar{Y} + \frac{\mu_{11}}{\sigma_X^2} (X - \bar{X}) \\ \Rightarrow E(Y|X) &= \bar{Y} + r \frac{\sigma_Y}{\sigma_X} (X - \bar{X}) \end{aligned}$$

By starting with the line  $E(X|y) = A + By$  and proceeding similarly we shall obtain the equation of the line of regression of  $X$  on  $Y$  as

$$E(X|Y) = \bar{X} + \frac{\mu_{11}}{\sigma_Y^2} (Y - \bar{Y}) = \bar{X} + r \frac{\sigma_X}{\sigma_Y} (Y - \bar{Y})$$

### Example 10.19. Given

$$f(x, y) = xe^{-x(y+1)}, x \geq 0, y \geq 0,$$

find the regression curve of  $Y$  on  $X$ . [B.H. Univ. M.Sc., 1989]

**Solution.** Marginal p.d.f. of  $X$  is given by

$$\begin{aligned} f_1(x) &= \int_0^{\infty} f(x, y) dy = \int_0^{\infty} xe^{-x(y+1)} dy \\ &= xe^{-x} \int_0^{\infty} e^{-xy} dy = xe^{-x} \left[ \frac{e^{-xy}}{-x} \right]_0^{\infty} \\ &= e^{-x}, x \geq 0 \end{aligned}$$

Conditional p.d.f. of  $Y$  on  $X$  is given by

$$f(y|x) = \frac{f(x, y)}{f_1(x)} = \frac{xe^{-x(y+1)}}{e^{-x}} = xe^{-xy}, y \geq 0.$$

The regression curve of  $Y$  on  $X$  is given by

$$y = E(Y|X=x) = \int_0^{\infty} y f(y|x) dy = \int_0^{\infty} yxe^{-xy} dy$$

$$= x \left[ \left| \frac{ye^{-xy}}{-x} \right|_0^\infty + \int_0^\infty \frac{e^{-xy}}{x} dy \right] = 0 + \left| \frac{e^{-xy}}{-x} \right|_0^\infty = \frac{1}{x}$$

i.e.,  $y = \frac{1}{x} \Rightarrow xy = 1.$

which is the equation of a rectangular hyperbola. Hence the regression of  $Y$  on  $X$  is not linear.

**Example 10-20.** Obtain the regression equation of  $Y$  on  $X$  for the following distribution :

$$f(x, y) = \frac{y}{(1+x)^4} \exp\left(-\frac{y}{1+x}\right); x, y \geq 0$$

**Solution.** Marginal p.d.f. of  $X$  is given by

$$\begin{aligned} f_1(x) &= \int_0^\infty f(x, y) dy = \frac{1}{(1+x)^4} \int_0^\infty y e^{-y/(1+x)} dy \\ &= \frac{1}{(1+x)^4} \cdot \Gamma 2 \cdot (1+x)^2 \quad (\text{Using Gamma Integral}) \\ &= \frac{1}{(1+x)^2}; x \geq 0 \end{aligned}$$

The conditional p.d.f. of  $Y$  (for given  $X$ ) is

$$f(y|x) = \frac{f(x, y)}{f_1(x)} = \frac{y}{(1+x)^2} \exp\left(-\frac{y}{1+x}\right); y \geq 0$$

Regression equation of  $Y$  on  $X$  is given by

$$\begin{aligned} y = E(Y|X) &= \int_0^\infty y f(y|x) dy = \frac{1}{(1+x)^2} \int_0^\infty y^2 e^{-y/(1+x)} dx \\ &= \frac{1}{(1+x)^2} \cdot \Gamma 3 \cdot (1+x)^3 \quad [\text{Using Gamma Integral}] \\ \Rightarrow y &= 2(1+x) \quad [\because \Gamma 3 = 2! = 2] \end{aligned}$$

Hence the regression of  $Y$  on  $X$  is linear.

**Example 10-21.** Let  $(X, Y)$  have the joint p.d.f. given by

$$f(x, y) = \begin{cases} 1, & \text{if } |y| < x, 0 < x < 1 \\ 0, & \text{otherwise} \end{cases}$$

Show that the regression of  $Y$  on  $X$  is linear but regression of  $X$  on  $Y$  is not linear.

**Solution.**  $|y| < x \Rightarrow -x < y < x$  and  $x > |y|$ .

The marginal p.d.f.'s  $f_1(\cdot)$  of  $X$  and  $f_2(\cdot)$  of  $Y$  are given by :

$$f_1(x) = \int_{-x}^x f(x, y) dy = \int_{-x}^x 1 dy = 2x; 0 < x < 1$$

$$f_2(y) = \int_{-|y|}^1 f(x, y) dx = \int_{-|y|}^1 1 dx = 1 - |y|; -1 < y < 1$$

$$\therefore f_1(x|y) = \frac{f(x,y)}{f_2(y)} = \frac{1}{1-|y|}; -1 \leq y < 1, 0 < x < 1$$

$$= \begin{cases} \frac{1}{1-y}, & 0 < y < 1; 0 < x < 1 \\ \frac{1}{1+y}, & -1 < y < 0; 0 < x < 1 \end{cases}$$

$$f_2(y|x) = \frac{f(x,y)}{f_1(x)} = \frac{1}{2x}, 0 < x < 1; |y| < x$$

$$E(Y|X=x) = \int_{-x}^x y \cdot f_2(y|x) dy = \int_{-x}^x \frac{y}{2x} dy = \frac{1}{4x} \cdot |y^2| \Big|_{-x}^x = 0$$

Hence the curve of regression of  $Y$  on  $X$  is  $y = 0$ , which is a straight line.

$$E(X|Y=y) = \int x f_1(x|y) dx$$

$$\therefore E(X|Y=y) = \int_0^1 x \left( \frac{1}{1-y} \right) dx = \frac{1}{2(1-y)}, 0 < y < 1$$

$$\text{and } E(X|Y=y) = \int_0^1 x \left( \frac{1}{1+y} \right) dx = \frac{1}{2(1+y)}, -1 < y < 0$$

Hence the curve of regression of  $X$  on  $Y$  is

$$x = \begin{cases} \frac{1}{2(1-y)}, & 0 < y < 1 \\ \frac{1}{2(1+y)}, & -1 < y < 0, \end{cases}$$

which is not a straight line.

**Example 10-22.** Variables  $X$  and  $Y$  have the joint p.d.f.

$$f(x,y) = \frac{1}{3}(x+y), 0 \leq x \leq 1, 0 \leq y \leq 2.$$

*Find :*

- (i)  $r(X, Y)$
- (ii) The two lines of regression
- (iii) The two regression curves for the means.

**Solution.** The marginal p.d.f.'s of  $X$  and  $Y$  are given by :

$$f_1(x) = \int_0^2 f(x,y) dy = \frac{1}{3} \int_0^2 (x+y) dy = \frac{1}{3} \left| xy + \frac{y^2}{2} \right|_0^2$$

$$\Rightarrow f_1(x) = \frac{2}{3}(1+x); 0 \leq x \leq 1 \quad \dots(1)$$

$$f_2(y) = \int_0^1 f(x,y) dx = \frac{1}{3} \int_0^1 (x+y) dx = \frac{1}{3} \left| \frac{x^2}{2} + xy \right|_0^1$$

$$\Rightarrow f_2(y) = \frac{1}{3} \left( \frac{1}{2} + y \right); 0 \leq y \leq 2 \quad \dots(2)$$

The conditional distributions are given by :

$$f_3(y|x) = \frac{f(x,y)}{f_1(x)} = \frac{1}{2} \left( \frac{x+y}{1+x} \right)$$

$$f_4(x|y) = \frac{f(x,y)}{f_2(y)} = \frac{2(x+y)}{1+2y} \quad . \quad \dots(3)$$

$$\begin{aligned} E(Y|x) &= \int_0^2 y \cdot f_3(y|x) dy = \frac{1}{2(1+x)} \int_0^2 y(x+y) dy \\ &= \frac{1}{2(1+x)} \left[ \frac{xy^2}{2} + \frac{y^3}{3} \right]_{y=0}^{y=2} = \frac{3x+4}{3(x+1)} \end{aligned}$$

Similarly, we shall get

$$E(X|y) = \int_0^1 x f_4(x|y) dx = \frac{2}{1+2y} \int_0^1 (x^2 + xy) dx = \frac{2+3y}{3(1+2y)}$$

(iii) Hence the regression curves for means are :

$$y = E(Y|x) = \frac{3x+4}{3(x+1)} \quad \text{and} \quad x = E(X|y) = \frac{2+3y}{3(1+2y)}.$$

From the marginal distributions we shall get

$$E(X) = \int_0^1 x f_1(x) dx = \frac{5}{9}, \quad E(X^2) = \int_0^1 x^2 f_1(x) dx = \frac{7}{18}$$

$$\Rightarrow \text{Var}(X) = \sigma_X^2 = \frac{7}{18} - \left(\frac{5}{9}\right)^2 = \frac{13}{162}$$

Similarly, we shall get

$$E(Y) = \frac{11}{9}, \quad E(Y^2) = \frac{16}{9}; \quad \sigma_Y^2 = \frac{16}{9} - \left(\frac{11}{9}\right)^2 = \frac{23}{81}$$

$$\begin{aligned} \text{Also } E(XY) &= \int_0^1 \int_0^2 xy f(x,y) dx dy = \frac{1}{3} \int_0^1 \int_0^2 (x^2 y + xy^2) dx dy \\ &= \frac{1}{3} \left\{ \left( \int_0^1 x^2 dx \right) \left( \int_0^2 y dy \right) + \left( \int_0^1 x dx \right) \left( \int_0^2 y^2 dy \right) \right\} \\ &= \frac{1}{3} \left[ \frac{1}{3} \times 2 + \frac{1}{2} \times \frac{8}{3} \right] = \frac{2}{3} \end{aligned}$$

$$\therefore \text{Cov}(X, Y) = E(XY) - E(X)E(Y) = \frac{2}{3} - \frac{5}{9} \times \frac{11}{9} = -\frac{1}{81}$$

$$(i) \quad r(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \cdot \sigma_Y} = \frac{-\frac{1}{81}}{\sqrt{\frac{13}{162} \times \frac{23}{81}}} = -\left(\frac{2}{299}\right)^{1/2}$$

(ii) The two lines of regression are :

$$Y - E(Y) = \frac{\text{Cov}(X, Y)}{\sigma_X^2} [X - E(X)] \Rightarrow Y - \frac{11}{9} = -\frac{2}{13} \left( X - \frac{5}{9} \right)$$

10-58

## Fundamentals of Mathematical Statistics

$$\text{and } X - E(X) = \frac{\text{Cov}(X, Y)}{\sigma_Y^2} [Y - E(Y)] \Rightarrow X - \frac{5}{9} = -\frac{1}{23} \left( Y - \frac{11}{9} \right)$$

**10-7-3. Regression Coefficients.** 'b', the slope of the line of regression of  $Y$  on  $X$  is also called the coefficient of regression of  $Y$  on  $X$ . It represents the increment in the value of dependent variable  $Y$  corresponding to a unit change in the value of independent variable  $X$ . More precisely, we write

$$b_{YX} = \text{Regression coefficient of } Y \text{ on } X = \frac{\mu_{11}}{\sigma_X^2} = r \frac{\sigma_Y}{\sigma_X} \quad \dots(10-17)$$

Similarly, the coefficient of regression of  $X$  on  $Y$  indicates the change in the value of variable  $X$  corresponding to a unit change in the value of variable  $Y$  and is given by

$$b_{XY} = \text{Regression coefficient of } X \text{ on } Y = \frac{\mu_{11}}{\sigma_Y^2} = r \frac{\sigma_X}{\sigma_Y} \quad \dots(10-17a)$$

**10-7-4. Properties of Regression Coefficients.**

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

**Proof.** Multiplying (10-17) and (10-17a), we get

$$b_{XY} \times b_{YX} = r \frac{\sigma_X}{\sigma_Y} \times r \frac{\sigma_Y}{\sigma_X} = r^2$$

$$\therefore r = \pm \sqrt{b_{XY} \times b_{YX}} \quad \dots(10-18)$$

**Remark.** We have

$$r = \frac{\mu_{11}}{\sigma_X \cdot \sigma_Y}, \quad b_{YX} = \frac{\mu_{11}}{\sigma_X^2} \quad \text{and} \quad b_{XY} = \frac{\mu_{11}}{\sigma_Y^2}$$

It may be noted that the sign of correlation coefficient is the same as that of regression coefficients, since the sign of each depends upon the co-variance term  $\mu_{11}$ . Thus if the regression coefficients are positive, 'r' is positive and if the regression coefficients are negative 'r' is negative.

From (10-18), we have

$$r = \pm \sqrt{b_{XY} \times b_{YX}}$$

the sign to be taken before the square root is that of the regression coefficients.

(b) 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, then we have to show that  $b_{XY} < 1$ .

$$\text{Now } b_{YX} > 1 \Rightarrow \frac{1}{b_{YX}} < 1 \quad \dots(*)$$

$$\text{Also } r^2 \leq 1 \Rightarrow b_{YX} \cdot b_{XY} \leq 1$$

$$\text{Hence } b_{XY} \leq \frac{1}{b_{YX}} < 1 \quad [\text{From } (*)]$$

(c) Arithmetic mean of the regression coefficients is greater than the correlation coefficient  $r$ , provided  $r > 0$ .

**Proof.** We have to prove that  $\frac{1}{2}(b_{YX} + b_{XY}) \geq r$

$$\text{or } \frac{1}{2} \left( r \frac{\sigma_Y}{\sigma_X} + r \frac{\sigma_X}{\sigma_Y} \right) \geq r \quad \text{or } \frac{\sigma_Y}{\sigma_X} + \frac{\sigma_X}{\sigma_Y} \geq 2 \quad (\because r > 0)$$

$$\Rightarrow \sigma_Y^2 + \sigma_X^2 - 2\sigma_X\sigma_Y \geq 0 \quad \text{i.e., } (\sigma_Y - \sigma_X)^2 \geq 0$$

which is always true, since the square of a real quantity is  $\geq 0$ .

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

**Proof.** Let  $U = \frac{X - a}{h}$ ,  $V = \frac{Y - b}{k} \Rightarrow X = a + hU$ ,  $Y = b + kV$ ,

where  $a, b, h (> 0)$  and  $k (> 0)$  are constants.

Then  $\text{Cov}(X, Y) = hk \text{Cov}(U, V)$ ,  $\sigma_X^2 = h^2\sigma_U^2$  and  $\sigma_Y^2 = k^2\sigma_V^2$

$$\begin{aligned} b_{YX} &= \frac{\mu_{11}}{\sigma_X^2} = \frac{hk \text{cov}(U, V)}{h^2\sigma_U^2} \\ &= \frac{k}{h} \cdot \frac{\text{cov}(U, V)}{\sigma_U^2} = \frac{k}{h} b_{UV} \end{aligned}$$

Similarly, we can prove that

$$b_{XY} = (h/k) b_{UV}$$

**10.7.5. Angle Between Two Lines of Regression.** Equations of the lines of regression of  $Y$  on  $X$ , and  $X$  on  $Y$  are

$$Y - \bar{y} = r \cdot \frac{\sigma_Y}{\sigma_X} (X - \bar{x}) \quad \text{and} \quad X - \bar{x} = r \cdot \frac{\sigma_X}{\sigma_Y} (Y - \bar{y})$$

Slopes of these lines are  $r \cdot \frac{\sigma_Y}{\sigma_X}$  and  $\frac{\sigma_Y}{r\sigma_X}$  respectively. If  $\theta$  is the angle between the two lines of regression then

$$\begin{aligned} \tan \theta &= \frac{r \frac{\sigma_Y}{\sigma_X} - \frac{\sigma_Y}{r\sigma_X}}{1 + r \frac{\sigma_Y}{\sigma_X} \cdot \frac{\sigma_Y}{r\sigma_X}} = \frac{r^2 - 1}{r} \left( \frac{\sigma_X\sigma_Y}{\sigma_X^2 + \sigma_Y^2} \right) \\ &= \frac{1 - r^2}{r} \left( \frac{\sigma_X\sigma_Y}{\sigma_X^2 + \sigma_Y^2} \right) \quad (\because r^2 \leq 1) \\ \therefore \theta &= \tan^{-1} \left\{ \frac{1 - r^2}{r} \left( \frac{\sigma_X\sigma_Y}{\sigma_X^2 + \sigma_Y^2} \right) \right\} \quad \dots (10.19) \end{aligned}$$

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

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

**Case (ii).** ( $r = \pm 1$ ). If  $r = \pm 1$ ,  $\tan \theta = 0 \Rightarrow \theta = 0$  or  $\pi$ .

In this case the two lines of regression either coincide or they are parallel to each other. But since both the 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, the two lines of regression coincide.

**Remarks 1.** Whenever two lines intersect, there are two angles between them, one acute angle and the other obtuse angle. Further  $\tan \theta > 0$  if  $0 < \theta < \pi/2$ , i.e.,  $\theta$  is an acute angle and  $\tan \theta < 0$  if  $\pi/2 < \theta < \pi$ , i.e.,  $\theta$  is an obtuse angle and since  $0 < r^2 < 1$ , the acute angle ( $\theta_1$ ) and obtuse angle  $\theta_2$  between the two lines of regression are given by

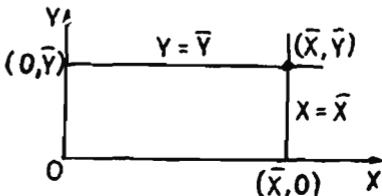
$$\theta_1 = \text{Acute angle} = \tan^{-1} \left\{ \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \cdot \frac{1 - r^2}{r} \right\}, r > 0$$

and  $\theta_2 = \text{Obtuse angle} = \tan^{-1} \left\{ \frac{\sigma_x \cdot \sigma_y}{\sigma_x^2 + \sigma_y^2} \cdot \frac{r^2 - 1}{r} \right\}, r > 0$

**2.** When  $r = 0$ , i.e., variables  $X$  and  $Y$  are uncorrelated, then the lines of regressions of  $Y$  on  $X$  and  $X$  on  $Y$  are given respectively by : [From (10-14a) and (10-15a)]

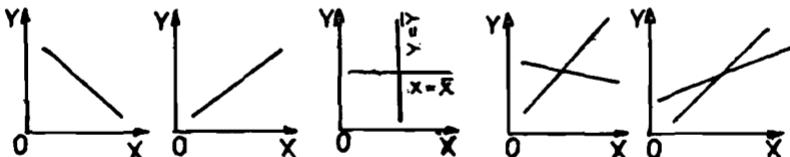
$$Y = \bar{Y} \text{ and } X = \bar{X},$$

as shown in the adjoining diagram. Hence, in this case ( $r = 0$ ), the lines of regression are perpendicular to each other and are parallel to  $X$ -axis and  $Y$ -axis respectively.



**3.** The fact that if  $r = 0$  (variables uncorrelated), the two lines of regression are perpendicular to each and if  $r = \pm 1$ ,  $\theta = 0$ , i.e., the two lines coincide, leads us to the conclusion that for higher degree of correlation between the variables, the angle between the lines is smaller, i.e., the two lines of regression are nearer to each other. On the other hand, if the lines of regression make a larger angle, they indicate a poor degree of correlation between the variables and ultimately for  $\theta = \pi/2$ ,  $r = 0$ , i.e., the lines become perpendicular if no correlation exists between the variables. Thus by plotting the lines of regression on a graph paper, we can have an approximate idea about the degree of correlation between the two variables under study. Consider the following illustrations :

TWO LINES COINCIDE ( $r = -1$ )	TWO LINES COINCIDE ( $r = +1$ )	TWO LINES PERPENDICULAR ( $r = 0$ )	TWO LINES APART (LOW DEGREE OF CORRELATION)	TWO LINES APART (HIGH DEGREE OF CORRELATION)
------------------------------------	------------------------------------	--	---	--



**10-7-6. Standard Error of Estimate or Residual Variance.** The equation of the line of regression of  $Y$  on  $X$  is

$$Y = \bar{Y} + r \frac{\sigma_Y}{\sigma_X} (X - \bar{X})$$

$$\Rightarrow \frac{Y - \bar{Y}}{\sigma_Y} = r \left( \frac{X - \bar{X}}{\sigma_X} \right)$$

The residual variance  $s_Y^2$  is the expected value of the squares of deviations of the observed values of  $Y$  from the expected values as given by the line of regression of  $Y$  on  $X$ . Thus

$$\begin{aligned}s_Y^2 &= E[(Y - (\bar{Y} + r\sigma_Y(X - \bar{X})/\sigma_X))^2] \\&= \sigma_Y^2 E \left[ \frac{Y - \bar{Y}}{\sigma_Y} - r \left( \frac{X - \bar{X}}{\sigma_X} \right) \right]^2 = \sigma_Y^2 E(Y^* - rX^*)^2\end{aligned}$$

where  $X^*$  and  $Y^*$  are standardised variates so that

$$E(X^*^2) = 1 = E(Y^*^2) \text{ and } E(X^* Y^*) = r.$$

$$\therefore s_Y^2 = \sigma_Y^2 [E(Y^*^2) + r^2 E(X^*^2) - 2r E(X^* Y^*)] = \sigma_Y^2 (1 - r^2)$$

$$\Rightarrow s_Y = \sigma_Y (1 - r^2)^{1/2}$$

Similarly, the standard error of estimate of  $X$  is given by

$$s_X = \sigma_X (1 - r^2)^{1/2}$$

**Remarks 1.** Since  $s_X^2$  or  $s_Y^2 \geq 0$ , it follows that

$$(1 - r^2) \geq 0 \Rightarrow |r| \leq 1$$

Hence

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

2. If  $r = \pm 1$ ,  $s_X = s_Y = 0$  so that each deviation is zero, and the two lines of regression are coincident.

3. Since, as  $r^2 \rightarrow 1$ ,  $s_X$  and  $s_Y \rightarrow 0$ , it follows that departure of the value  $r^2$  from unity indicates the departure of the relationship between the variables  $X$  and  $Y$  from linearity.

4. From the definition of linear regression, the minima condition implies that  $s_Y$  or  $s_X$  is the minimum variance.

**10.7.7. Correlation Coefficient between Observed and Estimated Value.** Here we will find the correlation between  $Y$  and

$$\hat{Y} = \bar{Y} + r \frac{\sigma_Y}{\sigma_X} (X - \bar{X})$$

where  $\hat{Y}$  is the estimated value of  $Y$  as given by the line of regression of  $Y$  on  $X$ , which is given by

$$r(Y, \hat{Y}) = \frac{\text{Cov}(Y, \hat{Y})}{\sigma_Y \sigma_{\hat{Y}}}$$

We have

$$E(\hat{Y}) = E \left[ \bar{Y} + r \frac{\sigma_Y}{\sigma_X} (X - \bar{X}) \right] = \bar{Y} + r \frac{\sigma_Y}{\sigma_X} E(X - \bar{X}) = \bar{Y}$$

$$\therefore \text{Var}(\hat{Y}) = E[\hat{Y} - E(\hat{Y})]^2 = E \left[ r \frac{\sigma_Y}{\sigma_X} (X - \bar{X}) \right]^2 = r^2 \sigma_Y^2$$

$$\Rightarrow \hat{\sigma}_Y = r \sigma_Y$$

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

$$= E\left[\{b(X - E(X))\} \left\{r \frac{\sigma_Y}{\sigma_X} (X - \bar{X})\right\}\right]$$

$$= br \frac{\sigma_Y}{\sigma_X} E[(X - E(X))^2] = \left(r \frac{\sigma_Y}{\sigma_X}\right)^2 \sigma_X^2 = r^2 \sigma_Y^2$$

$$\therefore r(Y, \hat{Y}) = \frac{r^2 \sigma_Y^2}{\sigma_Y r \sigma_Y} = r = r(X, Y)$$

Hence the correlation coefficient between observed and estimated value of  $Y$  is the same as the correlation coefficient between  $X$  and  $Y$ .

**Example 10-23.** Obtain the equations of the lines of regression for the data in Example 10-1. Also obtain the estimate of  $X$  for  $Y = 70$ .

**Solution.** Let  $U = X - 68$  and  $V = Y - 69$ , then

$$\bar{U} = 0, \bar{V} = 0, \sigma_U^2 = 4.5, \sigma_V^2 = 5.5, \text{Cov}(U, V) = 3 \text{ and } r(U, V) = 0.6$$

Since correlation coefficient is independent of change of origin, we get

$$r = r(X, Y) = r(U, V) = 0.6$$

We know that if  $U = \frac{X - a}{h}$ ,  $V = \frac{Y - b}{k}$ , then

$$\bar{X} = a + h\bar{U}, \bar{Y} = b + k\bar{V}, \sigma_X = h\sigma_U \text{ and } \sigma_Y = k\sigma_V$$

In our case  $h = k = 1$ ,  $a = 68$  and  $b = 69$ .

$$\text{Thus } \bar{X} = 68 + 0 = 68, \bar{Y} = 69 + 0 = 69$$

$$\sigma_X = \sigma_U = \sqrt{4.5} = 2.12 \text{ and } \sigma_Y = \sigma_V = \sqrt{5.5} = 2.35$$

Equation of line of regression of  $\hat{Y}$  on  $X$  is

$$Y - \bar{Y} = r \frac{\sigma_Y}{\sigma_X} (X - \bar{X})$$

$$\text{i.e., } Y = 69 + 0.6 \times \frac{2.35}{2.12} (X - 68) \Rightarrow Y = 0.665 X + 23.78$$

Equation of line of regression of  $X$  on  $Y$  is

$$X - \bar{X} = r \frac{\sigma_X}{\sigma_Y} (Y - \bar{Y})$$

$$\Rightarrow X = 68 + 0.6 \times \frac{2.12}{2.35} (Y - 69) \text{ i.e., } X = 0.54Y + 30.74$$

To estimate  $X$  for given  $Y$ , we use the line of regression of  $X$  on  $Y$ . If  $Y = 70$ , estimated value of  $X$  is given by

$$\hat{X} = 0.54 \times 70 + 30.74 = 68.54,$$

where  $\hat{X}$  is estimate of  $X$ .

**Example 10-24.** In a partially destroyed laboratory record of an analysis of correlation data, the following results only are legible :

Variance of  $X = 9$ .

Regression equations :  $8X - 10Y + 66 = 0$ ,  $40X - 18Y = 214$ .

What were (i) the mean values of  $X$  and  $Y$ ,

(ii) the correlation coefficient between  $X$  and  $Y$ , and

(iii) the standard deviation of  $Y$ ?

[Punjab Univ. B.Sc. (Hons.), 1993]

**Solution** (i) Since both the lines of regression pass through the point  $(\bar{X}, \bar{Y})$ , we have  $8\bar{X} - 10\bar{Y} + 66 = 0$ , and  $40\bar{X} - 18\bar{Y} = 214$ .

Solving, we get  $\bar{X} = 13$ ,  $\bar{Y} = 17$ .

(ii) Let  $8X - 10Y + 66 = 0$  and  $40X - 18Y = 214$  be the lines of regression of  $Y$  on  $X$  and  $X$  on  $Y$  respectively. These equations can be put in the form :

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

$$\therefore b_{YX} = \text{Regression coefficient of } Y \text{ on } X = \frac{8}{10} = \frac{4}{5}$$

$$\text{and } b_{XY} = \text{Regression coefficient of } X \text{ on } Y = \frac{18}{40} = \frac{9}{20}$$

$$\text{Hence } r^2 = b_{YX} \cdot b_{XY} = \frac{4}{5} \cdot \frac{9}{20} = \frac{9}{25}$$

$$\therefore r = \pm \frac{3}{5} = \pm 0.6$$

But since both the regression coefficients are positive, we take  $r = +0.6$

$$(iii) \text{ We have } b_{YX} = r \cdot \frac{\sigma_Y}{\sigma_X} \Rightarrow \frac{4}{5} = \frac{3}{5} \times \frac{\sigma_Y}{3} [\because \sigma_X^2 = 9 \text{ (Given)}]$$

$$\text{Hence } \sigma_Y = 4$$

**Remarks.** 1. It can be verified that the values of  $\bar{X} = 13$  and  $\bar{Y} = 17$  as obtained in part (i) satisfy both the regression equations. In numerical problems of this type, this check should invariably be applied to ascertain the correctness of the answer.

2. If we had assumed that  $8X - 10Y + 66 = 0$ , is the equation of the line of regression of  $X$  on  $Y$  and  $40X - 18Y = 214$  is the equation of line of regression of  $Y$  on  $X$ , then we get respectively :

$$8X = 10Y - 66 \text{ and } 18Y = 40X - 214$$

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

$$\Rightarrow b_{XY} = \frac{18}{8} \text{ and } b_{YX} = \frac{40}{18}$$

$$\therefore r^2 = b_{XY} \cdot b_{YX} = \frac{10}{8} \times \frac{40}{18} = 2.78$$

But since  $r^2$  always lies between 0 and 1, our supposition is wrong.

**Example 10-25.** Find the most likely price in Bombay corresponding to the price of Rs. 70 at Calcutta from the following :

	Calcutta	Bombay
Average price	65	67
Standard deviation	2.5	3.5

Correlation coefficient between the prices of commodities in the two cities is 0.8: [Nagpur Univ. B.Sc., 1993; Sri Venkateswara Univ. B.Sc. (Oct.) 1990]

**Solution.** Let the prices, (in Rupees), in Bombay and Calcutta be denoted by  $Y$  and  $X$  respectively. Then we are given

$\bar{X} = 65$ ,  $\bar{Y} = 67$ ,  $\sigma_X = 2.5$ ,  $\sigma_Y = 3.5$  and  $r = r(X, Y) = 0.8$ . We want  $Y$  for  $X = 70$ .

Line of regression of  $Y$  on  $X$  is

$$Y - \bar{Y} = r \frac{\sigma_Y}{\sigma_X} (X - \bar{X})$$

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

$$\text{When } X = 70, \quad \hat{Y} = 67 + 0.8 \times \frac{3.5}{2.5} (70 - 65) = 72.6$$

**Example 10-26.** Can  $Y = 5 + 2.8X$  and  $X = 3 - 0.5Y$  be the estimated regression equations of  $Y$  on  $X$  and  $X$  on  $Y$  respectively? Explain your answer with suitable theoretical arguments. [Delhi Univ. M.A.(Eco.), 1986]

**Solution.** Line of regression of  $Y$  on  $X$  is :

$$Y = 5 + 2.8X \Rightarrow b_{YX} = 2.8 \quad \dots(*)$$

Line of regression of  $X$  on  $Y$  is :

$$X = 3 - 0.5Y \Rightarrow b_{XY} = -0.5 \quad \dots(**)$$

This is not possible, since each of the regression coefficients  $b_{YX}$  and  $b_{XY}$  must have the same sign, which is same as that of  $\text{Cov}(X, Y)$ . If  $\text{Cov}(x, y)$  is positive, then both the regression coefficients are positive and if  $\text{Cov}(X, Y)$  is negative, then both the regression coefficients are negative. Hence (\*) and (\*\*) cannot be the estimated regression equations of  $Y$  on  $X$  and  $X$  on  $Y$  respectively.

### EXERCISE .10 (d)

1. (a) Explain what are regression lines. Why are there two such lines? Also derive their equations.

(b) Define (i) Line of regression, (ii) Regression coefficient. Find the equations to the lines of regression and show that the coefficient of correlation is the geometric mean of coefficients of regression.

(c) What equation is the equivalent mathematical statement for the following words?

"If the respective deviations in each series,  $X$  and  $Y$ , from their means were expressed in units of standard deviations, i.e., if each were divided by the

standard deviation of the series; to which it belongs and plotted to a scale of standard deviations, the slope of a straight line best describing the plotted points would be the correlation coefficient  $r$ .

**2(a)** Obtain the equation of the line of regression of  $Y$  on  $X$  and show that the angle  $\theta$ , between the two lines of regression is given by

$$\tan \theta = \frac{1 - \rho^2}{\rho} \times \frac{\sigma_1 \sigma_2}{\sigma_1^2 + \sigma_2^2}$$

where  $\sigma_1, \sigma_2$  are the standard deviations of  $X$  and  $Y$  respectively, and  $\rho$  is the correlation coefficient. **(Delhi Univ. B.Sc. (Maths. Hons.), 1989)**

Interpret the cases when  $\rho = 0$  and  $\rho = \pm 1$ .

**(Bangalore Univ. B.Sc. 1990)**

**(b)** If  $\theta$  is the acute angle between the two regression lines with correlation coefficient  $r$ , show that  $\sin \theta \leq 1 - r^2$ .

**3. (a)** Explain the term "regression" by giving examples. Assuming that the regression of  $Y$  on  $X$  is linear, outline a method for the estimation of the coefficients in the regression line based on the random paired sample of  $X$  and  $Y$ , and show that the variance of the error of the estimate for  $Y$  for the regression line is  $\sigma_Y^2(1 - \rho^2)$ , where  $\sigma_Y^2$  is the variance of  $Y$  and  $\rho$  is the correlation coefficient between  $X$  and  $Y$ .

**(b)** Prove that  $X$  and  $Y$  are linearly related if and only if  $\rho_{XY}^2 = 1$ . Further show that the slope of the regression line is positive or negative according as  $\rho = +1$  or  $\rho = -1$ .

**(c)** Let  $X$  and  $Y$  be two variates. Define  $X^* = \frac{X - a}{b}$ ,  $Y^* = \frac{Y - c}{d}$  for some constants  $a, b, c$  and  $d$ . Show that the regression line (least square) of  $Y$  on  $X$  can be obtained from that of  $Y^*$  on  $X^*$ .

**(d)** Show that the coefficient of correlation between the observed and the estimated values of  $Y$  obtained from the line of regression of  $Y$  on  $X$ , is the same as that between  $X$  and  $Y$ .

**4.** Two variables  $X$  and  $Y$  are known to be related to each other by the relation  $Y = X/(aX + b)$ . How is the theory of linear regression to be employed to estimate the constants  $a$  and  $b$  from a set of  $n$  pairs of observations  $(x_i, y_i)$ ,  $i = 1, 2, \dots, n$  ?

**Hint.**

$$\frac{1}{Y} = \frac{aX + b}{X} = a + \frac{b}{X}$$

**Put**

$$\frac{1}{X} = U \text{ and } \frac{1}{Y} = V$$

$\therefore$

$$V = a + bU$$

**5.** Derive the standard error of estimate of  $Y$  obtained from the linear regression equation of  $Y$  on  $X$ . What does this standard error measure ?

**6. (a)** Calculate the coefficient of correlation from the following data :

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

Also obtain the equations of the lines of regression and obtain an estimate of  $Y$  which should correspond on the average to  $X = 6.2$ .

**Ans.**  $r = 0.95$ ,  $Y - 12 = 0.95(X - 5)$ ,  $X - 5 = 0.95(Y - 12)$ ,  $13.14$

(b) Why do we have, in general, two lines of regression ? Obtain the regression of  $Y$  on  $X$ , and  $X$  on  $Y$  from the following table and estimate the blood pressure when the age is 45 years :

Age in years (X)	Blood pressure (Y)	Age in years (X)	Blood pressure (Y)
56	147	55	150
42	125	49	145
72	160	38	115
36	118	42	140
63	149	68	152
47	128	60	155

**Ans.**  $Y = 1.138X + 80.778$ ,  $Y = 131.988$  for  $X = 45$ .

(c) Suppose the observations on  $X$  and  $Y$  are given as :

X :	59	65	45	52	60	62	70	55	45	49
Y :	75	70	55	65	60	69	80	65	59	61

where  $N = 10$  students, and  $Y$  = Marks in Maths,  $X$  = Marks in Economics. Compute the least square regression equations of  $Y$  on  $X$  and of  $X$  on  $Y$ .

If a student gets 61 marks in Economics, what would you estimate his marks in Maths to be ?

7. (a) In a correlation analysis on the ages of wives and husbands, the following data were obtained. Find

(i) the value of the correlation coefficient, and (ii) the lines of regression.

Estimate the age of husband whose wife's age is 31 years. Estimate the age of wife whose husband is 40 years old.

Age of wife		15—25	25—35	35—45	45—55	55—65
Age of Husband ↓	→	30	6	3	—	—
		18	32	15	12	8
15—30	2	28	40	16	9	
30—45	—	4	9	10	8	
45—60	—	—	—	—	—	
60—75	—	—	—	—	—	

(b) The following table gives the distribution of total cultivable area ( $X$ ) and area under cultivation ( $Y$ ) in a district of 69 villages.

Calculate (i) the linear regression of  $Y$  on  $X$ ,

**Correlation and Regression**

10-67

(ii) the correlation coefficient  $r(X, Y)$ , and (iii) the average area under wheat corresponding to total area of 1,000 Bighas.

		Total area (in Bighas)				
		0—500	500—1000	1000—1500	1500—2000	2000—2500
Area under wheat	0—200	12	6	...	...	...
	200—400	2	18	4	2	1
	400—600	...	4	7	3	...
	600—800	...	1	...	2	1
	800—1000	...	...	1	2	3

$$\text{Ans. (i)} \quad Y = 0.7641X - 455.3854, \quad \text{(ii)} \quad r(X, Y) = 0.756$$

$$\text{(iii)} \quad Y = 308.7146 \text{ for } X = 1000$$

8. (a) Compare and contrast the roles of correlation and regression in studying the inter-dependence of two variates.

For 10 observations on price ( $X$ ) and supply ( $Y$ ) the following data were obtained (in appropriate units).

$$\Sigma X = 130, \quad \Sigma Y = 220, \quad \Sigma X^2 = 2288, \quad \Sigma Y^2 = 5506 \text{ and } \Sigma XY = 3467$$

Obtain the line of regression of  $Y$  on  $X$  and estimate the supply when the price is 16 units, and find out the standard error of the estimate.

$$\text{Ans. } Y = 8.8 + 1.015X, \quad 25.04$$

(b) If a number  $X$  is chosen at random from among the integers 1, 2, 3, 4 and a number  $Y$  is chosen from among those at least as large as  $X$ , prove that

$$\text{Cov}(X, Y) = \frac{5}{8}$$

Find also the regression line of  $X$  on  $Y$ .

(c) Calculate the correlation coefficient from the following data:—

$$N = 100, \quad \Sigma X = 12500 \quad \Sigma Y = 8000$$

$$\Sigma X^2 = 1585000, \quad \Sigma Y^2 = 648100 \quad \Sigma XY = 1007425.$$

Also obtain the regression equation of  $Y$  on  $X$ .

9. (a) The means of a bivariate frequency distribution are at (3, 4) and  $r = 0.4$ . The line of regression of  $Y$  on  $X$  is parallel to the line  $Y = X$ . Find the two lines of regression and estimate the mean of  $X$  when  $Y = 1$ .

(b) For certain data,  $Y = 1.2X$  and  $X = 0.6Y$ , are the regression lines. Compute  $\rho(X, Y)$  and  $\sigma_X/\sigma_Y$ . Also compute  $\rho(X, Z)$ , if  $Z = Y - X$ .

(c) The equations of two regression lines obtained in a correlation analysis are as follows :

$$3X + 12Y = 19, \quad 3Y + 9X = 46$$

Obtain (i) the value of correlation coefficient,

(ii) mean values of  $X$  and  $Y$ , and

(iii) the ratio of the coefficient of variability of  $X$  to that of  $Y$ .

Ans. (i)  $-\frac{1}{2}\sqrt{3}$ , (ii)  $\bar{X} = 5$ ,  $\bar{Y} = 1/3$ .

(d) For an army personnel of strength 25, the regression of weight of kidneys ( $Y$ ) on weight of heart ( $X$ ), both measured in ounces is

$$Y - 0.399X - 6.934 = 0$$

and the regression of weight of heart on weight of kidney is

$$X - 1.212Y + 2.461 = 0$$

Find the correlation coefficient between  $X$  and  $Y$  and their mean values. Can you find out the standard deviation of  $X$  and  $Y$  as well?

Ans.  $r(X, Y) = 0.70$ ,  $\bar{X} = 11.5086$ ,  $\bar{Y} = 11.5261$ , No.

(e) Find the coefficient of correlation for distribution in which

$$\text{S.D. of } X = 3.0 \text{ units}$$

$$\text{S.D. of } Y = 1.4 \text{ units}$$

Coefficient of regression of  $Y$  on  $X = 0.28$ .

10. (a) Given that  $X = 4Y + 5$  and  $Y = kX + 4$ , are the lines of regression of  $X$  on  $Y$  and  $Y$  on  $X$  respectively, show that  $0 < 4k < 1$ . If  $k = \frac{1}{16}$ , find the means of the two variables and coefficient of correlation between them.

[Punjab Univ. B.Sc. (Hons.), 1989]

Hint.  $X = 4Y + 5 \Rightarrow b_{XY} = 4$

$$Y = kx + 4 \Rightarrow b_{YX} = k$$

$$\therefore r^2 = 4k \quad \dots(*)$$

$$\text{But } 0 \leq r^2 \leq 1 \Rightarrow 0 \leq 4k \leq 1.$$

$$\text{If } k = \frac{1}{16}, \text{ then from } (*), \text{ we get}$$

$$r^2 = 4 \times \frac{1}{16} \Rightarrow r = +\frac{1}{2} \quad [\text{Since both the regression coefficient are positive}]$$

$$\text{For } k = \frac{1}{16}, \text{ the two lines of regression become}$$

$$X = 4Y + 5 \text{ and } Y = \frac{1}{16}X + 4$$

Solving the two equations, we get  $\bar{Y} = 5.75$ ,  $\bar{X} = 28$ .

(b) For 50 students of a class the regression equation of marks in Statistics ( $X$ ) on marks in Mathematics ( $Y$ ) is  $3Y - 5X + 180 = 0$ . The mean marks in Mathematics is 44 and variance of marks in Statistics is  $9/16$ th of the variance of marks in Mathematics. Find the mean marks in Statistics and the coefficient of correlation between marks in two subjects.

[Bangalore Univ. B.Sc., 1989]

Hint. We are given  $n = 50$ ,  $\bar{Y} = 44$

$$\text{and } \sigma_X^2 = \frac{9}{16} \sigma_Y^2 \Rightarrow \frac{\sigma_X}{\sigma_Y} = \frac{3}{4} \quad \dots(*)$$

The equation of the line of regression of  $X$  on  $Y$  is given to be

$$3Y - 5X + 180 = 0 \Rightarrow X = \frac{3}{5}Y + \frac{180}{5}$$

$$\therefore b_{XY} = r \frac{\sigma_X}{\sigma_Y} = \frac{3}{5} \Rightarrow r \cdot \frac{3}{4} = \frac{3}{5} \quad \text{or} \quad r = 0.8$$

Since the lines of regression pass through the point  $(\bar{X}, \bar{Y})$ , we get

$$\bar{X} = \frac{3}{5}\bar{Y} + \frac{180}{5} = \frac{3}{5} \times 44 + 36 = 62.4$$

(c) Out of the two lines of regression given by

$$X + 2Y - 5 = 0 \text{ and } 2X + 3Y - 8 = 0,$$

which one is the regression line of  $X$  on  $Y$ ?

Use the equations to find the mean of  $X$  and the mean of  $Y$ . If the variance of  $X$  is 12, calculate the variance of  $Y$ .

Ans.  $\bar{X} = 1, \bar{Y} = 2, \sigma_Y^2 = 4$

(d) The lines of regression in a bivariate distribution are :

$$X + 9Y = 7 \text{ and } Y + 4X = \frac{49}{3}$$

Find (i) the coefficient of correlation, (iii) the ratios  $\sigma_X^2 : \sigma_Y^2 : \text{Cov}(X, Y)$ , (iii) the means of the distribution and (iv)  $E(X | Y = 1)$ .

(e) Estimate  $X$  when  $Y = 10$ , if the two lines of regression are :

$$X = -\frac{1}{18}Y + \lambda \text{ and } Y = -2x + \mu,$$

$(\lambda, \mu)$  being unknown and the mean of the distribution is at  $(-1, 2)$ . Also compute  $r, \lambda$  and  $\mu$ . [Gujarat Univ. B.Sc., Oct. 1992]

11. (a) The following results were obtained in the analysis of data on yield of dry bark in ounces ( $Y$ ) and age in years ( $X$ ) of 200 cinchona plants :

	X	Y
Average	9.2	16.5
Standard deviation	2.1	4.2
Correlation coefficient	+0.84	

Construct the two lines of regression and estimate the yield of dry bark of a plant of age 8 years. [Patna Univ. B.Sc., 1991],

(b) The following data pertain to the marks in subjects  $A$  and  $B$  in a certain examination :

Mean marks in  $A = 39.5$

Mean marks in  $B = 47.5$

Standard deviation of marks in  $A = 10.8$

Standard deviation of marks in  $B = 16.8$

Coefficient of correlation between marks in  $A$  and marks in  $B = 0.42$ .

Draw the two lines of regression and explain why there are two regression equations. Give the estimate of marks in  $B$  for candidates who secured 50 marks in  $A$ .

Ans.  $Y = 0.65X + 21.825, X = 0.27Y + 26.675$  and  $Y = 54.342$  for  $X = 50$

(c) You are given the following information about advertising expenditure and sales :

	<i>Advertising Expenditure (X)</i> (Rs. lakhs)	<i>Sales (Y)</i> (Rs. lakhs)
Mean	10	90
s.d.	3	12

Correlation coefficient = 0.8

What should be the advertising budget if the company wants to attain sales target of Rs. 120 lakhs ? [Delhi Univ. M.C.A., 1990]

12. Twenty-five pairs of value of variates  $X$  and  $Y$  led to the following results :

$$N = 25, \sum X = 127, \sum Y = 100, \sum X^2 = 760, \sum Y^2 = 449 \text{ and } \sum XY = 500$$

A subsequent scrutiny showed that two pairs of values were copied down as :

<i>X</i>	<i>Y</i>
8	14
8	6

<i>X</i>	<i>Y</i>
8	12
6	8

(i) Obtain the correct value of the correlation coefficient.

(ii) Hence or otherwise, find the correct equations of the two lines of regression.

(iii) Find the angle between the regression lines.

Ans. (i)  $r(X, Y) = -(0.64 \times 0.15)^{1/2}$ ,

$$(ii) X = -0.64Y + 7.56, Y = -0.15X + 4.75.$$

13. Suppose you have  $n$  observations :

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

on two variables  $X$  and  $Y$ , and you have fitted a linear regression  $Y = a + bX$  by the method of least squares. Denote the 'expected' value of  $Y$  by  $Y^*$ , and the residual  $Y - Y^*$  by  $e$ . Find means and variances of  $Y^*$  and  $e$ , and the correlation co-efficient between (i)  $X$  and  $e$ , (ii)  $Y$  and  $e$  and (iii)  $Y$  and  $Y^*$ . Use these results to bring out the significance and limitations of the correlation coefficient.

Ans.  $r(X, e) = 0, r(Y, e) = 0$  and  $r(Y, Y^*) = r(X, Y)$ .

14. (a) The regression lines of  $Y$  on  $X$  and of  $X$  on  $Y$  are respectively  $Y = aX + b$  and  $X = cY + d$ . Show that

(i) Means are  $\bar{X} = (bc + d)/(1 - ac)$  and  $\bar{Y} = (ad + b)/(1 - ac)$

(ii) Correlation coefficient between  $X$  and  $Y$  is  $\sqrt{ac}$ .

(iii) The ratio of the standard deviations of  $X$  and  $Y$  is  $\sqrt{c/a}$ .

(b) For two random variables  $X$  and  $Y$  with the same mean, the two regression equations are  $Y = aX + b$  and  $X = \alpha Y + \beta$ . Show that  $\frac{b}{\beta} = \frac{1 - \alpha}{1 - \alpha}$ . Find also the common mean.

[Punjab Univ.B.Sc. (Maths Hons.), 1992]

(c) If the lines of regression of  $Y$  on  $X$  and  $X$  on  $Y$  are respectively  $a_1X + b_1Y + c_1 = 0$  and  $a_2X + b_2Y + c_2 = 0$ , prove that  $a_1b_2 \leq a_2b_1$ .

(Delhi Univ. B.Sc. (Stat. Hons.), 1989)

**Hint.**  $r^2 = b_{YX} \cdot b_{XY} \leq 1 \Rightarrow \left(-\frac{a_1}{b_1}\right) \times \left(-\frac{b_2}{a_2}\right) = \frac{a_1b_2}{a_2b_1} \leq 1$

15. (a) By minimising  $\sum_{i=1}^n f_i(x_i \cos \alpha + y_i \sin \alpha - p)^2$  for variations in  $\alpha$  and  $p$ , show that there are two straight lines passing through the mean of the distribution for which the sum of squares of normal deviations has an extreme value. Prove also that their slopes are given by

$$\tan 2\alpha = \frac{2\mu_{11}}{\sigma_x^2 - \sigma_y^2}$$

**Hint.** We have to minimize

$$S = \sum_{i=1}^n f_i(x_i \cos \alpha + y_i \sin \alpha - p)^2 \quad \dots(1)$$

Equating to zero, the partial derivatives of (1) w.r.t.  $\alpha$  and  $p$ , we have

$$\frac{\partial S}{\partial \alpha} = 0 = 2 \sum_{i=1}^n f_i(x_i \cos \alpha + y_i \sin \alpha - p)(-x_i \sin \alpha + y_i \cos \alpha) \quad \dots(2)$$

$$\frac{\partial S}{\partial p} = 0 = -2 \sum_{i=1}^n f_i(x_i \cos \alpha + y_i \sin \alpha - p) \quad \dots(3)$$

Equation (3) can be written as

$$\sum_{i=1}^n f_i(x_i \cos \alpha + y_i \sin \alpha - p) = 0 \Rightarrow \bar{x} \cos \alpha + \bar{y} \sin \alpha - p = 0 \quad \dots(4)$$

From equation (2), we get a quadratic equation which shows that there are two straight lines for extreme values of  $E$ .

From equation (4), it becomes clear that both the straight lines pass through the point  $(\bar{x}, \bar{y})$ .

Again equation (2) can be written as :

$$\begin{aligned} & \sum_{i=1}^n f_i(x_i \cos \alpha + y_i \sin \alpha - p)(y_i \cos \alpha - x_i \sin \alpha) = 0 \\ \Rightarrow & \sum_{i=1}^n f_i [\cos \alpha (x_i - \bar{x}) + \sin \alpha (y_i - \bar{y})] [y_i \cos \alpha - x_i \sin \alpha] = 0 \\ & \quad [\text{Using (4)}] \\ \Rightarrow & \cos^2 \alpha \sum_{i=1}^n f_i y_i (x_i - \bar{x}) - \sin \alpha \cos \alpha \sum_{i=1}^n f_i x_i (x_i - \bar{x}) \\ & + \sin \alpha \cos \alpha \sum_{i=1}^n f_i y_i (y_i - \bar{y}) - \sin^2 \alpha \sum_{i=1}^n f_i x_i (y_i - \bar{y}) = 0 \quad \dots(5) \end{aligned}$$

We have  $\mu_{11} = \frac{1}{N} \sum_i f_i (x_i - \bar{x})(y_i - \bar{y})$

$$= \frac{1}{N} \sum_i f_i x_i (y_i - \bar{y}) - \bar{x} \cdot \frac{1}{N} \sum_i f_i (y_i - \bar{y}) = \frac{1}{N} \sum_i f_i x_i (y_i - \bar{y})$$

Similarly,  $\mu_{11} = \frac{1}{N} \sum_i f_i y_i (x_i - \bar{x})$

$$\sigma_x^2 = \frac{1}{N} \sum_i f_i x_i (x_i - \bar{x}) \text{ and } \sigma_y^2 = \frac{1}{N} \sum_i f_i y_i (y_i - \bar{y})$$

Substituting these values in (5), we get the required result.

(b) If the straight line defined by

$$Y = a + bX$$

satisfies the condition  $E[(Y - a - bX)^2] = \text{minimum}$ , show that the regression line of the random variable  $Y$  on the random variable  $X$  is

$$Y - \bar{Y} = r \frac{\sigma_Y}{\sigma_X} (X - \bar{X}), \text{ where } \bar{X} = E(X), \bar{Y} = E(Y)$$

16. (a) Define Curve of regression of  $Y$  on  $X$ .

The joint density function of  $X$  and  $Y$  is given by :

$$\begin{aligned} f(x, y) &= x + y, 0 < x < 1, 0 < y < 1 \\ &= 0, \text{ otherwise} \end{aligned}$$

Find

- (i) the correlation coefficient between  $X$  and  $Y$ ,
- (ii) the regression curve of  $Y$  on  $X$ , and
- (iii) the regression curve of  $X$  on  $Y$ .

Ans.  $\rho(X, Y) = -\frac{1}{11}$ . [Madras Univ. B.Sc., Stat. (Main), 1992]

(b) Let  $f(x_1, x_2) = \frac{2}{a^2}; 0 < x_1 < x_2, 0 < x_2 < a$   
 $= 0, \text{ elsewhere}$

be the joint p.d.f. of  $X_1$  and  $X_2$ .

Find conditional means and variances. Also show that  $\rho = \frac{1}{2}$ .

17. If the joint density of  $X$  and  $Y$  is given by

$$f(x, y) = \begin{cases} (x + y)/3, & \text{for } 0 < x < 1, 0 < y < 2 \\ 0, & \text{otherwise} \end{cases}$$

obtain the regressions (i) of  $Y$  on  $X$  and (ii) of  $X$  on  $Y$ .

Are the regressions linear? Find the correlation coefficient between  $X$  and  $Y$ . (Allahabad Univ. B.Sc. 1992)

$$\text{Ans. } y = E(Y|x) = \frac{3x + 4}{3(x + 1)}; x = E(X|y) = \frac{2 + 3y}{3(1 + 2y)}$$

$$\text{Corr. } (X, Y) = -\left(\frac{2}{299}\right)^{1/2}$$

18. Let the joint density function of  $X$  and  $Y$  be given by

$$\begin{aligned} f(x, y) &= 8xy, 0 < x < y < 1 \\ &= 0, \text{ otherwise} \end{aligned}$$

Find: (i)  $E(Y|X = x)$ , (ii)  $E[XY|X = x]$ , (iii)  $\text{Var}[Y|X = x]$

[Delhi Univ. BSc. (Maths Hons.), 1988]

Ans. (i)  $E(Y|x) = \frac{2}{3} \left( \frac{1+x+x^2}{1+x} \right)$ ;  $E(XY|x) = x E(Y|x)$ , (iii)  $E(Y^2|x) = \frac{1+x^2}{2}$

19. Give an example to show that it is possible to have the regression of  $Y$  on  $X$  constant (does not depend on  $X$ ), but the regression of  $X$  on  $Y$  is not constant (does depend on  $Y$ ).

Hint. See Example 10.21

20. Prove or disprove :

$$E(Y|X = x) = \text{constant} \Rightarrow r(X,Y) = 0$$

Ans. True

21. If  $f(x,y) = \frac{1}{3}x^2 \exp[-y(1+x)]$ ,  $x \geq 0, y \geq 0$ , is the joint p.d.f. of  $(X,Y)$ , obtain the equation of regression of  $Y$  on  $X$ .

Ans.  $y = E(Y|x) = 1/(1+x)$ .

22. Variables  $(X,Y)$  have joint p.d.f.

$$f(x,y) = 6(1-x-y), x > 0, y > 0, x+y < 1, \\ = 0, \text{ otherwise.}$$

Find  $f_X(x), f_Y(y)$  and  $\text{Cov}(X,Y)$ . Are  $X$  and  $Y$  independent? Obtain the regression curves for the means.

[Calcutta Univ. B.Sc. (Maths Hons.), 1986]

Ans.  $f_1(x) = 3(1-x)^2, 0 < x < 1; f_2(y) = 3(1-y)^2, 0 < y < 1$ .

$X$  and  $Y$  are not independent.

Regression curves for the means are:

$$y = E(Y|x) = \frac{1}{3}(1-x) \text{ and } x = E(X|y) = \frac{1}{3}(1-y).$$

23. For the joint p.d.f.

$$f(x,y) = 3x^2 - 8xy + 6y^2, 0 \leq x, y \leq 1,$$

find the least square regression lines and the regression curves for the means.

[Calcutta Univ. B.Sc. (Maths, Hons.), 1987]

Ans. Regression lines :

$$y - \frac{2}{3} = -\frac{10}{67}\left(x - \frac{5}{12}\right); \quad x - \frac{5}{12} = -\frac{25}{32}\left(y - \frac{2}{3}\right)$$

Regression curves for means are :

$$y = E(Y|x) = \frac{9x^2 - 16x + 9}{6(3x^2 - 4x + 2)}; \quad x = E(X|y) = \frac{36y^2 - 32y - 9}{12(6y^2 - 4y + 1)}$$

24. Let  $(X, Y)$  be jointly distributed with p.d.f.

$$f(x,y) = e^{-y}, 0 < x < y < \infty$$

$$= 0, \text{ otherwise}$$

Prove that :

$$E(Y|X = x) = x + \bar{x} \quad \text{and} \quad E(X|Y = y) = y/2.$$

Hence prove that  $r(X, Y) = \sqrt{1/2}$ .

25. Let  $f(x, y) = e^{-y} (1 - e^{-x})$ ,  $0 < x < y ; 0 < y < \infty$   
 $= e^{-x} (1 - e^{-y})$ ,  $0 < y < x ; 0 < x < \infty$

- (a) Show that  $f(x, y)$  is a p.d.f.
- (b) Find marginal distributions of  $X$  and  $Y$ .
- (c) Find  $E(Y|X = x)$  for  $x > 0$ .
- (d) Find  $P(X \leq 2, Y \leq 2)$ .
- (e) Find the correlation coefficient  $r(X, Y)$ .
- (f) Find another joint p.d.f. having the same marginals.

Ans. (b)  $f_1(x) = xe^{-x}$ ,  $0 < x < \infty$ ;  $f_2(y) = ye^{-y}$ ,  $0 < y < \infty$ .

(c)  $E(Y|x) = \frac{1-e^x}{x}[x-1] + \frac{1}{x}\left(\frac{x^2}{2} + xe^x + e^{-x} - 1\right)$

(d)  $1 - \frac{1}{e^4} - \frac{4}{e^2}$ ; (e)  $r(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{1}{\sqrt{2}} \frac{1}{\sqrt{2}} = \frac{1}{2}$

(f) Hint.  $f(x, y, \alpha) = f_1(x)f_2(y)[1 + \alpha(2F(x)-1)(2F(y)-1)]$ ,  $|\alpha| < 1$ , has the same marginals  $f_1(x)$  and  $f_2(y)$ .

26. Obtain regression equation of  $Y$  on  $X$  for the distributions :

(a)  $f(x, y) = \frac{9}{2} \cdot \frac{1+x+y}{(1+x)^4(1+y)^4}$ ;  $x, y \geq 0$

(b)  $f(x, y) = \frac{4}{5}(x+3y)e^{-x-2y}$ ;  $x, y \geq 0$

[Sardar Patel Univ. M.Sc., 1992]

Ans. (a) Hint. See Example 5-25, page 5-55, (b)  $\frac{x+3}{2x+3}$ .

27. A ball is drawn at random from an urn containing three white balls numbered 0, 1, 2; two red balls numbered 0, 1 and one black ball numbered 0. If the colours white, red and black are again numbered 0, 1 and 2 respectively, find the correlation coefficient between the variates  $X$ , the colour number and  $Y$  the number of the ball. Write down the equation of regression line of  $Y$  on  $X$ .

[Calcutta Univ. B.Sc. (Maths. Hons.), 1986]

### OBJECTIVE TYPE QUESTIONS

I. State, giving reasons, whether each of the following statements is true or false.

- (i) Both regression lines of  $Y$  on  $X$  and of  $X$  on  $Y$  do not intersect at all.
- (ii) In a bivariate regression,  $b_{YX} = \frac{1}{5}$ ,  $b_{XY} = 10$
- (iii) The regression coefficient of  $Y$  on  $X$  is 3.2 and that of  $X$  on  $Y$  is 0.8.
- (iv) There is no relationship between correlation coefficient and regression coefficient.
- (v) Both the regression coefficients cannot exceed unity.

- (vi) The greater the value of ' $r$ ', the better are the estimates obtained through regression analysis.
- (vii) If  $X$  and  $Y$  are negatively correlated variables, and  $(0, 0)$  is on the least squares line of  $Y$  on  $X$ , and if  $X = 1$  is the observed value then predicted value of  $Y$  must be negative.
- (viii) Let the correlation between  $X$  and  $Y$  be perfect and positive. Suppose the points  $(3, 5)$  and  $(1, 4)$  are on the regression lines. With this knowledge it is possible to determine the least squares line exactly.
- (ix) If the lines of regression are  $Y = \frac{1}{4}X$  and  $X = \frac{1}{9}Y + 1$ , then  $\rho = \frac{1}{6}$  and  $E(X | Y = 0) = 1$ .
- (x) In a bivariate distribution,  $b_{YX} = 2.8$  and  $b_{XY} = 0.3$ .

**II. Fill in the blanks :**

- (i) The regression analysis measures ... between  $X$  and  $Y$ .
- (ii) Lines of regression are ... if  $r_{XY} = 0$  and they are ... if  $r_{XY} = \pm 1$ .
- (iii) If the regression coefficients of  $X$  on  $Y$  and  $Y$  on  $X$  are  $-0.4$  and  $-0.9$  respectively then the correlation coefficient between  $X$  and  $Y$  is ...
- (iv) If the two regression lines are  $X + 3Y - 5 = 0$  and  $4X + 3Y - 8 = 0$ , then the correlation coefficient between  $X$  and  $Y$  is ...
- (v) If one of the regression coefficients is ... unity, the other must be ... unity.
- (vi) The farther the two regression lines cut each other, the ... will be the degree of correlation.
- (vii) When one regression coefficient is positive, the other would be ...
- (viii) The sign of regression coefficient is ... as that of correlation coefficient.
- (ix) Correlation coefficient is the... between regression coefficients.
- (x) Arithmetic mean of regression coefficients is ... correlation coefficient.
- (xi) When the correlation coefficient is zero, the two regression lines are ... and when it is  $\pm 1$ , then the regression lines are ...

**III. Indicate the correct answer :**

- (i) The regression line of  $Y$  on  $X$  (a) minimises total of the squares of horizontal deviations, (b) total of the squares of the vertical deviations, (c) both vertical and horizontal deviations, (d) none of these.
- (ii) The regression coefficients are  $b_2$  and  $b_1$ . Then the correlation coefficient  $r$  is (a)  $b_1/b_2$ , (b)  $b_2/b_1$ , (c)  $b_1b_2$  (d)  $\pm \sqrt{b_1 b_2}$ .
- (iii) The farther the two regression lines cut each other (a) the greater will be the degree of correlation, (b) the lesser will be the degree of correlation, (c) does not matter.

- (iv) If one regression coefficient is greater than unity, then the other must be (a) greater than the first one, (b) equal to unity, (c) less than unity, (d) equal to zero.
- (v) When the correlation coefficient  $r = \pm 1$ , then the two regression lines (a) are perpendicular to each other; (b) coincide, (c) are parallel to each other, (d) do not exist.
- (vi) The two lines of regression are given as  $\bar{X} + 2\bar{Y} - 5 = 0$  and  $2\bar{X} + 3\bar{Y} = 8$ . Then the mean values of  $X$  and  $Y$  respectively are (a) 2, 1, (b) 1, 2, (c) 2, 5, (d) 2, 3.
- (vii) The tangent of the angle between two regression lines is given as 0.6 and the s.d. of  $Y$  is known to be twice that of  $X$ . Then the value of correlation coefficient between  $X$  and  $Y$  is (a)  $-\frac{1}{2}$ , (b)  $\frac{1}{2}$ , (c) 0.7, (d) 0.3.

IV.  $\sigma_x$  and  $\sigma_y$  are the standard deviations of two correlated variables  $X$  and  $Y$  respectively in a large sample, and  $r$  is the sample correlation coefficient.

- (i) State the "Standard Error of Estimate" for linear regression of  $Y$  on  $X$ .
- (ii) What is the standard error in estimating  $Y$  from  $X$  if  $r = 0$ ?
- (iii) By how much is this error reduced if  $r$  is increased to 0.30?
- (iv) How large must  $r$  be in order to reduce this standard error to one-half its value for  $r = 0$ ?
- (v) Give your interpretations for the cases  $r = 0$  and  $r = 1$ .

V. Explain why we have two lines of regression.

**10-8. Correlation Ratio.** As discussed earlier, when variables are linearly related, we have the regression lines of one variable on another variable and correlation coefficient can be computed to tell us about the extent of association between them. However, if the variables are not linearly related but some sort of curvilinear relationship exists between them, the use of  $r$  which is a measure of the degree to which the relation approaches a straight line "law" will be misleading. We might come across bivariate distributions where  $r$  may be very low or even zero but the regression may be strong, or even perfect. *Correlation ratio* ' $\eta$ ' is the appropriate measure of curvilinear relationship between the two variables. Just as  $r$  measures the concentration of points about the straight line of best fit,  $\eta$  measures the concentration of points about the curve of best fit. If regression is linear  $\eta = r$ , otherwise  $\eta > r$  (c.f. Remark 2, § 10-8-1).

**10-8-1. Measure of Correlation Ratio.** In the previous articles we have assumed that there is a single observed value  $Y$  corresponding to the given value  $x_i$  of  $X$  but sometimes there are more than one such value of  $Y$ .

Suppose corresponding to the values  $x_i$ , ( $i = 1, 2, \dots, m$ ) of the variable  $X$ , the variable  $Y$  takes the values  $y_{ij}$  with respective frequencies  $f_{ij}$ ,  $j = 1, 2, \dots, n$ .

Though all the  $x$ 's in the  $i$ th vertical array have the same value, the  $y$ 's are different. A typical pair of values in the  $i$ th array is  $(x_i, y_{ij})$ , with frequency  $f_{ij}$ .

Thus the first suffix  $i$  indicates the vertical array while the second suffix  $j$  indicates the positions of  $y$  in that array. Let

$$\sum_{j=1}^n f_{ij} = n_i \quad \text{and} \quad \sum_{i=1}^m \sum_{j=1}^n f_{ij} = \sum_{i=1}^m \left( \sum_{j=1}^n f_{ij} \right) = \sum_{i=1}^m n_i = N, \quad (\text{say}).$$

If  $\bar{y}_i$  and  $\bar{y}$  denote the means of the  $i$ th array and the overall mean respectively, then

$$\bar{y}_i = \frac{\sum_{j=1}^n f_{ij} y_{ij}}{\sum_{j=1}^n f_{ij}} = \frac{\sum_{j=1}^n f_{ij} y_{ij}}{n_i} = \frac{T_i}{n_i} \quad \text{and} \quad \bar{y} = \frac{\sum_{i=1}^m \sum_{j=1}^n f_{ij} y_{ij}}{\sum_{i=1}^m \sum_{j=1}^n f_{ij}} = \frac{\sum_{i=1}^m n_i \bar{y}_i}{\sum_{i=1}^m n_i} = \frac{T}{N}$$

In other words  $\bar{y}$  is the weighted mean of all the array means, the weights being the array frequencies.

**Def.** The correlation ratio of  $Y$  on  $X$ , usually denoted by  $\eta_{yx}$  is given by

$$\eta_{yx}^2 = 1 - \frac{\sigma_{ey}^2}{\sigma_y^2} \quad \dots (10.21)$$

where  $\sigma_{ey}^2$  and  $\sigma_y^2$  are given by

$$\sigma_{ey}^2 = \frac{1}{N} \sum_i \sum_j f_{ij} (y_{ij} - \bar{y}_i)^2 \quad \text{and} \quad \sigma_y^2 = \frac{1}{N} \sum_i \sum_j f_{ij} (y_{ij} - \bar{y})^2$$

A convenient expression for  $\eta_{yx}$  can be obtained in terms of standard deviation  $\sigma_{my}$  of the means of the vertical arrays, each mean being weighted by the array frequency.

We have

$$\begin{aligned} N\sigma_y^2 &= \sum_i \sum_j f_{ij} (y_{ij} - \bar{y})^2 = \sum_i \sum_j f_{ij} \{ (y_{ij} - \bar{y}_i) + (\bar{y}_i - \bar{y}) \}^2 \\ &= \sum_i \sum_j f_{ij} (y_{ij} - \bar{y}_i)^2 + \sum_i \sum_j f_{ij} (\bar{y}_i - \bar{y})^2 + 2 \sum_i \sum_j f_{ij} (y_{ij} - \bar{y}_i)(\bar{y}_i - \bar{y}) \end{aligned}$$

The term  $2[\sum_i (\bar{y}_i - \bar{y}) \{\sum_j f_{ij} (y_{ij} - \bar{y}_i)\}]$  vanishes since  $\sum_j f_{ij} (y_{ij} - \bar{y}_i) = 0$ ,

being the algebraic sum of the deviations from mean.

$$\begin{aligned} \therefore N\sigma_y^2 &= \sum_i \sum_j f_{ij} (y_{ij} - \bar{y}_i)^2 + \sum_i n_i (\bar{y}_i - \bar{y})^2 \\ \Rightarrow N\sigma_y^2 &= N\sigma_{ey}^2 + N\sigma_{my}^2 \Rightarrow \sigma_y^2 = \sigma_{ey}^2 + \sigma_{my}^2 \\ \Rightarrow 1 - \frac{\sigma_{ey}^2}{\sigma_y^2} &= \frac{\sigma_{my}^2}{\sigma_y^2} \end{aligned}$$

which on comparison with (10.21) gives

$$\eta_{yx}^2 = \frac{\sigma_{my}^2}{\sigma_y^2} = \frac{\sum_i n_i (\bar{y}_i - \bar{y})^2}{\sum_i \sum_j f_{ij} (y_{ij} - \bar{y})^2} \quad \dots (10.22)$$

We have

$$\begin{aligned} N\sigma_{mY}^2 &= \sum_i n_i (\bar{y}_i - \bar{y})^2 = \sum_i n_i \bar{y}_i^2 - N \bar{y}^2 = \sum_i \frac{T_i^2}{n_i} - \frac{T^2}{N} \\ \therefore \quad \eta_{YX}^2 &= \left[ \sum_i \left( \frac{T_i^2}{n_i} \right) - \frac{T^2}{N} \right] / N\sigma_Y^2, \end{aligned} \quad \dots(10-23)$$

a formula, much more convenient for computational purposes.

**Remarks 1.** (10-21) implies that

$$\sigma_{eY}^2 = \sigma_Y^2 (1 - \eta_{YX}^2)$$

Since  $\sigma_{eY}^2$  and  $\sigma_Y^2$  are non-negative, we have

$$1 - \eta_{YX}^2 \geq 0 \Rightarrow \eta_{YX}^2 \leq 1 \Rightarrow |\eta_{YX}| \leq 1$$

**2.** Since the sum of squares of deviations in any array is minimum when measured from its mean, we have

$$\sum_i \sum_j f_{ij} (y_{ij} - \bar{y}_i)^2 \leq \sum_i \sum_j f_{ij} (y_{ij} - \hat{y}_{ij})^2 \quad \dots(*)$$

where  $\hat{y}_{ij}$  is the estimate of  $y_{ij}$  for given value of  $X = x_i$ , say, as given by the line of regression of  $Y$  on  $X$  i.e.,  $\hat{y}_{ij} = a + bx_i$ , ( $j = 1, 2, \dots, n$ ).

$$\text{But } \sum_i \sum_j f_{ij} (y_{ij} - \bar{y}_i)^2 = N\sigma_{eY}^2 = N\sigma_Y^2 (1 - \eta_{YX}^2)$$

$$\text{and } \sum_i \sum_j f_{ij} (y_{ij} - a - bx_i)^2 = N\sigma_Y^2 (1 - r^2) \quad (\text{c.f. } \S \text{ 10-7-6})$$

$$\therefore (*) \Rightarrow 1 - \eta_{YX}^2 \leq 1 - r^2$$

$$\text{i.e., } \eta_{YX}^2 \geq r^2 \Rightarrow |\eta_{YX}| \geq |r|$$

Thus the absolute value of the correlation ratio can never be less than the absolute of  $r$ , the correlation coefficient.

When the regression of  $Y$  on  $X$  is linear, straight line of means of arrays coincides with the line of regression and  $\eta_{YX}^2 = r^2$ . Thus  $\eta_{YX}^2 - r^2$  is the departure of regression from linearity. It is also clear (from Remark 1) that the more nearly  $\eta_{YX}^2$  approaches unity, the smaller is  $\sigma_{eY}^2$  and, therefore, closer are the points to the curve of means of vertical arrays.

$$\text{When } \eta_{YX}^2 = 1, \sigma_{eY}^2 = 0 \Rightarrow \sum_i \sum_j f_{ij} (y_{ij} - \bar{y}_i)^2 = 0$$

$\Rightarrow y_{ij} = \bar{y}_i$ ,  $\forall j = 1, 2, \dots, n$ , i.e., all the points lie on the curve of means. This implies that there is a functional relationship between  $X$  and  $Y$ .  $\eta_{YX}$  is, therefore, the measure of the degree to which the association between the variables approaches a functional relationship of the form  $Y = F(X)$ , where  $F(X)$  is a single valued function of  $X$ ,  $[F(X) = a + bX]$ .

**3.** It is worth noting that the value of  $\eta_{YX}$  is not independent of the classification of the data. As the class intervals become narrower  $\eta_{YX}$  approaches unity, since in that case  $\sigma_{mY}^2$  gets nearer to  $\sigma_Y^2$ . If the grouping is so fine that only one item appears in each row (related to each  $x$ -class), that item will constitute the mean of that column and thus in this case  $\sigma_{mY}^2$  and  $\sigma_Y^2$  become equal so that  $\eta_{YX}^2 = 1$ . On the other hand, a very coarse grouping tends to make the value of  $\eta_{YX}$  approach  $r$ . "Student" has given a formula for 'the correction'

to be made in the correlation ratio 'for grouping' in Biometrika (Vol IX page 316-320.)

4. It can be easily proved that  $\eta_{YX}^2$  is independent of change of origin and scale of measurements.

5.  $\eta_{XY}^2$ , the second correlation ratio of  $X$  on  $Y$  depends upon the scatter of observations about the line of column means.

6.  $r_{XY}$  and  $r_{YX}$  are same but  $\eta_{YX}$  is, in general, different from  $\eta_{XY}$ .

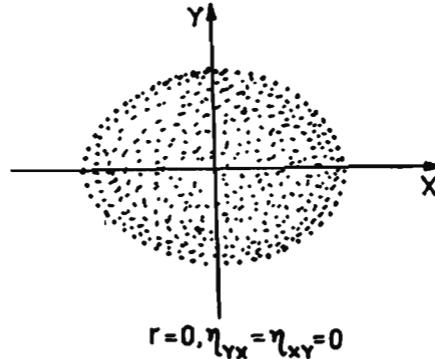
7. In terms of expectation, correlation ratio is defined as follows :

$$\eta_{YX}^2 = \frac{E_X [E(Y|X) - E(Y)]^2}{E[Y - E(Y)]^2} = \frac{E[E(Y|X) - E(Y)]^2}{\sigma_Y^2}$$

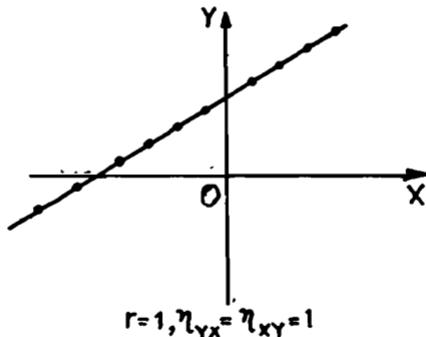
and  $\eta_{XY}^2 = \frac{E_Y [E(X|Y) - E(X)]^2}{E[X - E(X)]^2} = \frac{E[E(X|Y) - E(X)]^2}{\sigma_X^2}$

8. We give below some diagrams, exhibiting the relationship between  $r$  and  $\eta_{YX}$ .

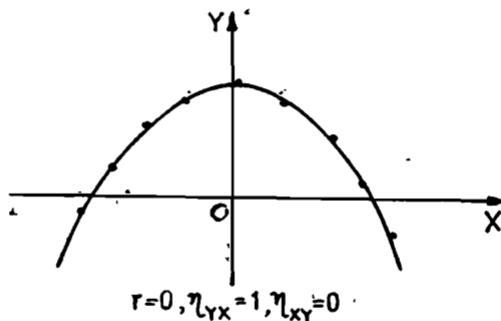
(i) For completely random scattering of the dots with no trend, both  $r$  and  $\eta$  are zero.



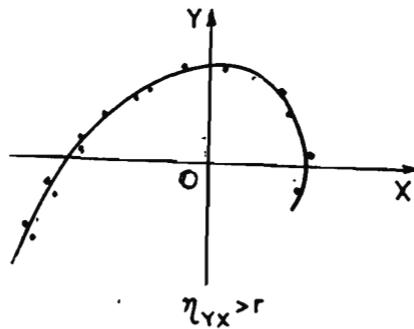
(ii) If dots lie precisely on a line,  $r = 1$  and  $\eta = 1$ .



(iii) If dots lie on a curve, such that no ordinate cuts it more than once,  $\eta_{YX} = 1$  and if furthermore, the dots are symmetrically placed about Y-axis, then  $\eta_{XY} = 0, r = 0$ .



(iv) If  $\eta_{YX} > r$ , the dots are scattered around a definitely curved trend line.



### EXERCISE 10(e)

1. (a) Define correlation coefficient and correlation ratio. When is the latter a more suitable measure of correlation than the former ? Show that the correlation ratio is never less than the correlation coefficient. What do you infer if the two are equal ? Further, show that none of these can exceed one.

[*Delhi Univ. B.Sc. (Stat. Hons.), 1988*]

(b) Show that  $1 \geq \eta_{YX}^2 \geq r_{YX}^2 \geq 0$

Interpret each of the following statements.

(i)  $r = 0$ , (ii)  $r^2 = 1$ , (iii)  $\eta^2 = 1$ , (iv)  $\eta^2 = r^2$  and (v)  $\eta = 0$

(c) When the correlation coefficient is equal to unity, show that the two correlation ratios are also equal to unity. Is the converse true ?

(d) Define correlation ratio  $\eta_{XY}$  and prove that

$$1 \geq \eta^2_{XY} \geq r^2,$$

where  $r$  is the coefficient of correlation between  $X$  and  $Y$ . Show further that  $(\eta^2_{XY} - r^2)$  is a measure of non-linearity of regression.

**2. For the joint p.d.f.**

$$f(x, y) = \frac{1}{2}x^3 \exp[-x(y+1)], \quad y > 0, x > 0 \\ = 0 \quad \text{otherwise},$$

find :

- (i) Two lines of regression.
- (ii) The regression curves for the means.
- (iii)  $r(X, Y)$ .
- (iv)  $\eta^2_{YX}$  and  $\eta^2_{XY}$ .

[Delhi Univ. B.A. (Stat. Hons. Spl. Course), 1987]

Ans. (i)  $y = -\frac{1}{6}x + 1$  ;  $x = -\frac{2}{3}y + \frac{10}{3}$

(ii)  $y = E(Y|x) = \frac{1}{x}$  ;  $x = E(X|y) = \frac{4}{1+y}$

(iii)  $r(X, Y) = -\frac{1}{3}$  (iv)  $\eta^2_{YX} = \frac{1}{3}$ ,  $\eta^2_{XY} = \frac{1}{5}$

**3. Compute  $r(X, Y)$  and  $\eta^2_{YX}$  for the following data :**

$X :$	0.5	— 1.5	1.5	— 2.5	2.5	— 3.5	3.5	— 4.5	4.5	— 5.5
$f :$	20		30		35		25		15	
$\bar{y} :$	11.3		12.7		14.7		16.5		19.1	

$\text{Var}(Y) = 9.61$

Ans.  $\eta^2_{YX} = 0.77$ ,  $r = 0.85$

**4. Compute  $\eta^2_{XY}$  for the following table :**

		$X$					
			47	52	57	62	67
$Y \downarrow$	57	4	4	2	...	...	
	62	4	8	8	1	...	
	67	...	7	12	1	4	
	72	...	3	1	8	5	
	77	...	...	3	5	6	

**10-9. Intra-class Correlation.** Intra-class correlation means within class correlation. It is distinguishable from product moment correlation in as much as here both the variables measure the same characteristics. Sometimes specially in biological and agricultural study, it is of interest to know how the members of a family or group are correlated among themselves with respect to some one of their common characteristic. For example, we may require the correlation between the heights of brothers of a family or between yields of plots of an experimental block. In such cases both the variables measure the same characteristic, e.g., height and height or weight and weight. There is

nothing to distinguish one from the other so that one may be treated as  $X$ -variable and the other as the  $Y$ -variable.

Suppose we have  $A_1, A_2, \dots, A_n$  families with  $k_1, k_2, \dots, k_n$  members, each of which may be represented as

$x_{11}$	$x_{21} \dots$	$x_{i1} \dots$	$x_{n1}$
$x_{12}$	$x_{22}$	$x_{i2}$	$x_{n2}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_{1j}$	$x_{2j} \dots$	$x_{ij} \dots$	$x_{nj}$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
$x_{1k_1}$	$x_{2k_2} \dots$	$x_{ik_i} \dots$	$x_{nk_n}$

and let  $x_{ij}$  ( $i = 1, 2, \dots, n; j = 1, 2, \dots, k_i$ ) denote the measurement on the  $j$ th member in the  $i$ th family.

We shall have  $k_i(k_i - 1)$  pairs for the  $i$ th family or group like  $(x_{ij}, x_{il})$ ,  $j \neq l$ . There will be  $\sum_{i=1}^n k_i(k_i - 1) = N$  pairs for all the  $n$  families or groups. If we prepare a correlation table there will be  $k_i(k_i - 1)$  entries for the  $i$ th group or family and  $\sum_i k_i(k_i - 1) = N$  entries for all the  $n$  families or groups. The table is symmetrical about the principal diagonal. Such a table is called an *intra-class correlation table* and the correlation is called *intra-class correlation*.

In the bivariate table  $x_{il}$  occurs  $(k_i - 1)$  times,  $x_{i2}$  occurs  $(k_i - 1)$  times, ...,  $x_{ik_i}$  occurs  $(k_i - 1)$  times, i.e., from the  $i$ th family we have  $(k_i - 1) \sum_j x_{ij}$  and hence for all the  $n$  families we have  $\sum_i (k_i - 1) \sum_j x_{ij}$  as the marginal frequency, the table being symmetrical about principal diagonal.

$$\therefore \bar{x} = \bar{y} = \frac{1}{N} \left[ \sum_i (k_i - 1) \sum_j x_{ij} \right]$$

Similarly,

$$\sigma_x^2 = \sigma_y^2 = \frac{1}{N} \left[ \sum_i (k_i - 1) \sum_j (x_{ij} - \bar{x})^2 \right]$$

Further

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{1}{N} \sum_i \left[ \sum_{j \neq l} (x_{ij} - \bar{x})(x_{il} - \bar{x}) \right], j \neq l \\ &= \frac{1}{N} \sum_i \left[ \sum_{j=1}^{k_i} \sum_{l=1}^{k_i} (x_{ij} - \bar{x})(x_{il} - \bar{x}) - \sum_{j=1}^{k_i} (x_{ij} - \bar{x})^2 \right] \end{aligned}$$

**Correlation and Regression**

If we write  $\bar{x}_i = \sum_j x_{ij} / k_i$ , then

$$\begin{aligned}\sum_i \left[ \sum_{j=1}^k \sum_{l=1}^k (x_{ij} - \bar{x}) (x_{il} - \bar{x}) \right] &= \sum_i \left[ \sum_j (x_{ij} - \bar{x}) \sum_l (x_{il} - \bar{x}) \right] \\ &= \sum_i [k_i (\bar{x}_i - \bar{x}) k_i (\bar{x}_i - \bar{x})] \\ &= \sum_i k_i^2 (\bar{x}_i - \bar{x})^2\end{aligned}$$

Therefore intra-class correlation coefficient is given by

$$r(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{V(X) V(Y)}} = \frac{\sum_i k_i^2 (\bar{x}_i - \bar{x})^2 - \sum_i \sum_j (x_{ij} - \bar{x})^2}{\sum_i \sum_j (k_i - 1) (x_{ij} - \bar{x})^2} \quad \dots(10-24)$$

If we put  $k_i = k$ , i.e., if all families have equal members then

$$\begin{aligned}r &= \frac{k^2 \sum_i (\bar{x}_i - \bar{x})^2 - \sum_i \sum_j (x_{ij} - \bar{x})^2}{(k-1) \sum_i \sum_j (x_{ij} - \bar{x})^2} \\ &= \frac{nk^2 \sigma_m^2 - nk\sigma^2}{(k-1) nk\sigma^2} = \frac{1}{(k-1)} \left\{ \frac{k \sigma_m^2}{\sigma^2} - 1 \right\} \quad \dots(10-24a)\end{aligned}$$

where  $\sigma^2$  denotes the variance of  $X$  and  $\sigma_m^2$  the variance of means of families.

**Limits.** We have from (10-24a),

$$1 + (k-1)r = \frac{k\sigma_m^2}{\sigma^2} \geq 0 \Rightarrow r \geq -\frac{1}{(k-1)}$$

Also  $1 + (k-1)r \leq k$ , as the ratio  $\frac{\sigma_m^2}{\sigma^2} \leq 1 \Rightarrow r \leq 1$

so that  $-\frac{1}{(k-1)} \leq r \leq 1$

**Interpretation.** Intraclass correlation cannot be less than  $-1/(k-1)$ , though it may attain the value +1 on the positive side, so that it is a skew coefficient and a negative value has not the same significance as a departure from independence as an equivalent positive value.

**10.11. Multiple and Partial Correlation.** When the values of one variable are associated with or influenced by other variable, e.g., the age of husband and wife, the height of father and son, the supply and demand of a commodity and so on, Karl Pearson's coefficient of correlation can be used as a measure of linear relationship between them. But sometimes there is interrelation between many variables and the value of one variable may be influenced by many others, e.g., the yield of crop per acre say ( $X_1$ ) depends upon quality of seed ( $X_2$ ), fertility of soil ( $X_3$ ), fertilizer used ( $X_4$ ), irrigation facilities ( $X_5$ ), weather conditions ( $X_6$ ) and so on. Whenever we are interested in studying the joint effect of a group of variables upon a variable not included in that group, our study is that of *multiple correlation and multiple regression*.

Suppose in a trivariate or multi-variate distribution we are interested in the relationship between two variables only. There are two alternatives, viz., (i) we

consider only those two members of the observed data in which the other members have specified values or (ii) we may eliminate mathematically the effect of other variates on two variates. The first method has the disadvantage that it limits the size of the data and also it will be applicable to only the data in which the other variates have assigned values. In the second method it may not be possible to eliminate the entire influence of the variates but the linear effect can be easily eliminated. The correlation and regression between only two variates eliminating the linear effect of other variates in them is called the *partial correlation and partial regression*.

**10-11-1. Yule's Notation.** Let us consider a distribution involving three random variables  $X_1, X_2$  and  $X_3$ . Then the equation of the plane of regression of  $X_1$  on  $X_2$  and  $X_3$  is

$$X_1 = a + b_{12.3}X_2 + b_{13.2}X_3 \quad \dots(10-28)$$

Without loss of generality we can assume that the variables  $X_1, X_2$  and  $X_3$  have been measured from their respective means, so that

$$E(X_1) = E(X_2) = E(X_3) = 0$$

Hence on taking expectation of both sides in (10-28), we get  $a = 0$ .

Thus the plane of regression of  $X_1$  on  $X_2$  and  $X_3$  becomes

$$X_1 = b_{12.3}X_2 + b_{13.2}X_3 \quad \dots(10-28a)$$

The coefficients  $b_{12.3}$  and  $b_{13.2}$  are known as the *partial regression coefficients* of  $X_1$  on  $X_2$  and of  $X_1$  on  $X_3$  respectively.

$$e_{1.23} = b_{12.3}X_2 + b_{13.2}X_3$$

is called the estimate of  $X_1$  as given by the plane of regression (10-28a) and the quantity

$$X_{1.23} = X_1 - b_{12.3}X_2 - b_{13.2}X_3,$$

is called the *error of estimate or residual*.

In the general case of  $n$  variables  $X_1, X_2, \dots, X_n$ , the equation of the plane of regression of  $X_1$  on  $X_2, X_3, \dots, X_n$  becomes

$$X_1 = b_{12.34\dots n}X_2 + b_{13.24\dots n}X_3 + \dots + b_{1n.23\dots(n-1)}X_n$$

The error of estimate or residual is given by

$$X_{1.23\dots n} = X_1 - (b_{12.34\dots n}X_2 + b_{13.24\dots n}X_3 + \dots + b_{1n.23\dots(n-1)}X_n)$$

The notations used here are due to Yule. The subscripts before the dot(.) are known as *primary subscripts* and those after the dot are called *secondary subscripts*. The order of a regression coefficient is determined by the number of secondary subscripts, e.g.,

$$b_{12.3}, b_{12.34}, \dots, b_{12.34\dots n}$$

are the regression coefficients of order 1, 2, ...,  $(n - 2)$  respectively. Thus in general, a regression coefficient with  $p$ -secondary subscripts will be called a regression coefficient of order ' $p$ '. It may be pointed out that the order in which the secondary subscripts are written is immaterial but the order of the primary subscripts is important, e.g., in  $b_{12.34\dots n}$ ,  $X_2$  is independent while  $X_1$  is dependent variable but in  $b_{21.34\dots n}$ ,  $X_1$  is independent while  $X_2$  is dependent

variable. Thus of the two primary subscripts, former refers to dependent variable and the latter to independent variable.

The order of a residual is also determined by the number of secondary subscripts in it, e.g.,  $X_{1 \cdot 23}, X_{1 \cdot 234}, \dots, X_{1 \cdot 23 \dots n}$  are the residuals of order 2, 3, ...,  $(n - 1)$  respectively.

**Remark.** In the following sequences we shall assume that the variables under consideration have been measured from their respective means.

**10-12. Plane of Regression.** The equation of the plane of regression of  $X_1$  on  $X_2$  and  $X_3$  is

$$X_1 = b_{12 \cdot 3} X_2 + b_{13 \cdot 2} X_3 \quad \dots(10-29)$$

The constants  $b$ 's in (10-29) are determined by the principle of least squares, i.e., by minimising the sum of the squares of the residuals, viz.,

$$S = \sum X_{1 \cdot 23}^2 = \sum (X_1 - b_{12 \cdot 3} X_2 - b_{13 \cdot 2} X_3)^2,$$

the summation being extended to the given values ( $N$  in number) of the variables.

The normal equations for estimating  $b_{12 \cdot 3}$  and  $b_{13 \cdot 2}$  are

$$\left. \begin{aligned} \frac{\partial S}{\partial b_{12 \cdot 3}} &= 0 = -2 \sum X_2 (X_1 - b_{12 \cdot 3} X_2 - b_{13 \cdot 2} X_3) \\ \frac{\partial S}{\partial b_{13 \cdot 2}} &= 0 = -2 \sum X_3 (X_1 - b_{12 \cdot 3} X_2 - b_{13 \cdot 2} X_3) \end{aligned} \right\} \quad \dots(10-30)$$

$$\text{i.e.,} \quad \sum X_2 X_{1 \cdot 23} = 0 \quad \text{and} \quad \sum X_3 X_{1 \cdot 23} = 0 \quad \dots(10-30a)$$

$$\Rightarrow \left. \begin{aligned} \sum X_1 X_2 - b_{12 \cdot 3} \sum X_2^2 - b_{13 \cdot 2} \sum X_2 X_3 &= 0 \\ \sum X_1 X_3 - b_{12 \cdot 3} \sum X_2 X_3 - b_{13 \cdot 2} \sum X_3^2 &= 0 \end{aligned} \right\} \quad \dots(10-30b)$$

Since  $X_i$ 's are measured from their respective means, we have

$$\left. \begin{aligned} \sigma_i^2 &= \frac{1}{N} \sum X_i^2, \quad \text{Cov}(X_i, X_j) = \frac{1}{N} \sum X_i X_j \\ \text{and} \quad r_{ij} &\doteq \frac{\text{Cov}(X_i, X_j)}{\sigma_i \sigma_j} = \frac{\sum X_i X_j}{N \sigma_i \sigma_j} \end{aligned} \right\} \quad \dots(10-30c)$$

Hence from (10-30b), we get

$$\left. \begin{aligned} r_{12} \sigma_1 \sigma_2 - b_{12 \cdot 3} \sigma_2^2 - b_{13 \cdot 2} r_{23} \sigma_2 \sigma_3 &= 0 \\ r_{13} \sigma_1 \sigma_3 - b_{12 \cdot 3} r_{23} \sigma_2 \sigma_3 - b_{13 \cdot 2} \sigma_3^2 &= 0 \end{aligned} \right\} \quad \dots(10-30d)$$

Solving the equations (10-30d) for  $b_{12 \cdot 3}$  and  $b_{13 \cdot 2}$ , we get

$$b_{12 \cdot 3} = \frac{\begin{vmatrix} r_{12} \sigma_1 & r_{23} \sigma_3 \\ r_{13} \sigma_1 & \sigma_3 \end{vmatrix}}{\begin{vmatrix} \sigma_2 & r_{23} \sigma_3 \\ r_{23} \sigma_2 & \sigma_3 \end{vmatrix}} = \frac{\sigma_1}{\sigma_2} \cdot \frac{\begin{vmatrix} r_{12} & r_{23} \\ r_{13} & 1 \end{vmatrix}}{\begin{vmatrix} 1 & r_{23} \\ r_{23} & 1 \end{vmatrix}} \quad \dots(10-31)$$

Similarly, we will get

$$b_{13 \cdot 2} = \frac{\sigma_1}{\sigma_3} \cdot \frac{\begin{vmatrix} 1 & r_{12} \\ r_{23} & r_{13} \end{vmatrix}}{\begin{vmatrix} 1 & r_{23} \\ r_{23} & 1 \end{vmatrix}} \quad \dots(10-31a)$$

If we write

$$\omega = \begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{vmatrix} \quad \dots(10-32)$$

and  $\omega_{ij}$  is the cofactor of the element in the  $i$ th row and  $j$ th column of  $\omega$ , we have from (10-31) and (10-31a)

$$b_{12 \cdot 3} = -\frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}} \quad \text{and} \quad b_{13 \cdot 2} = -\frac{\sigma_1}{\sigma_3} \cdot \frac{\omega_{13}}{\omega_{11}} \quad \dots(10-33)$$

Substituting these values in (10-29), we get the required equation of the plane of regression of  $X_1$  on  $X_2$  and  $X_3$  as

$$\begin{aligned} X_1 &= -\frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}} \cdot X_2 - \frac{\sigma_1}{\sigma_3} \cdot \frac{\omega_{13}}{\omega_{11}} \cdot X_3 \\ \Rightarrow \quad \frac{X_1}{\sigma_1} \cdot \omega_{11} + \frac{X_2}{\sigma_2} \cdot \omega_{12} + \frac{X_3}{\sigma_3} \cdot \omega_{13} &= 0 \end{aligned} \quad \dots(10-34)$$

Aliter. Eliminating the coefficient  $b_{12 \cdot 3}$  and  $b_{13 \cdot 2}$  in (10-29) and (10-30a), the required equation of the plane of regression of  $X_1$  on  $X_2$  and  $X_3$  becomes

$$\begin{vmatrix} X_1 & X_2 & X_3 \\ r_{12}\sigma_1\sigma_2 & \sigma_2^2 & r_{23}\sigma_2\sigma_3 \\ r_{13}\sigma_1\sigma_3 & r_{23}\sigma_2\sigma_3 & \sigma_3^2 \end{vmatrix} = 0$$

Dividing  $C_1$ ,  $C_2$  and  $C_3$  by  $\sigma_1$ ,  $\sigma_2$  and  $\sigma_3$  respectively and also  $R_2$  and  $R_3$  by  $\sigma_2$  and  $\sigma_3$  respectively, we get

$$\begin{aligned} \begin{vmatrix} \frac{X_1}{\sigma_1} & \frac{X_2}{\sigma_2} & \frac{X_3}{\sigma_3} \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{vmatrix} &= 0 \\ \Rightarrow \quad \frac{X_1}{\sigma_1} \omega_{11} + \frac{X_2}{\sigma_2} \omega_{12} + \frac{X_3}{\sigma_3} \omega_{13} &= 0 \end{aligned}$$

where  $\omega_{ij}$  is defined in (10-32).

**10-12-1. Generalisation.** In general, the equation of the plane of regression of  $X_1$  on  $X_2, X_3, \dots, X_n$  is

$$X_1 = b_{12 \cdot 34 \dots n} X_2 + b_{13 \cdot 24 \dots n} X_3 + \dots + b_{1n \cdot 23 \dots (n-1)} X_n \quad \dots(10-35)$$

The sum of the squares of residuals is given by

$$S = \sum X_{1 \cdot 23 \dots n}^2$$

$$= \sum (X_1 - b_{12\cdot34\ldots n} X_2 - b_{13\cdot24\ldots n} X_3 - \dots - b_{1n\cdot23\ldots(n-1)} X_n)^2$$

Using the principle of least squares, the normal equations for estimating the  $(n-1)$ ,  $b$ 's are

$$\frac{\partial S}{\partial b_{12\cdot34\ldots n}} = 0 = -2 \sum X_2 (X_1 - b_{12\cdot34\ldots n} X_2 - b_{13\cdot24\ldots n} X_3 - \dots - b_{1n\cdot23\ldots(n-1)} X_n)$$

$$\frac{\partial S}{\partial b_{13\cdot24\ldots n}} = 0 = -2 \sum X_3 (X_1 - b_{12\cdot34\ldots n} X_2 - b_{13\cdot24\ldots n} X_3 - \dots - b_{1n\cdot23\ldots(n-1)} X_n)$$

$$\frac{\partial S}{\partial b_{1n\cdot23\ldots(n-1)}} = 0 = -2 \sum X_n (X_1 - b_{12\cdot34\ldots n} X_2 - b_{13\cdot24\ldots n} X_3 - \dots - b_{1n\cdot23\ldots(n-1)} X_n) \quad \boxed{\dots(10.36)}$$

$$\text{i.e., } \sum X_i X_{1\cdot23\ldots n} = 0, \quad (i = 2, 3, \dots, n) \quad \boxed{\dots(10.36a)}$$

which on simplification after using (10.30c), give

$$r_{12}\sigma_1\sigma_2 = b_{12\cdot34\ldots n}\sigma_2^2 + b_{13\cdot24\ldots n}r_{23}\sigma_2\sigma_3 + \dots + b_{1n\cdot23\ldots(n-1)}r_{2n}\sigma_2\sigma_n$$

$$r_{13}\sigma_1\sigma_3 = b_{12\cdot34\ldots n}r_{23}\sigma_2\sigma_3 + b_{13\cdot24\ldots n}\sigma_3^2 + \dots + b_{1n\cdot23\ldots(n-1)}r_{3n}\sigma_3\sigma_n$$

$$r_{1n}\sigma_1\sigma_n = b_{12\cdot34\ldots n}r_{2n}\sigma_2\sigma_n + b_{13\cdot24\ldots n}r_{3n}\sigma_3\sigma_n + \dots + b_{1n\cdot23\ldots(n-1)}\sigma_n^2 \quad \boxed{\dots(10.36b)}$$

Hence the eliminant of  $b$ 's between (10.35) and (10.36b) is

$$\left| \begin{array}{ccccc} X_1 & X_2 & X_3 & \dots & X_n \\ r_{12}\sigma_1\sigma_2 & \sigma_2^2 & r_{23}\sigma_2\sigma_3 & \dots & r_{2n}\sigma_2\sigma_n \\ r_{13}\sigma_1\sigma_3 & r_{23}\sigma_2\sigma_3 & \sigma_3^2 & \dots & r_{3n}\sigma_3\sigma_n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_{1n}\sigma_1\sigma_n & r_{2n}\sigma_2\sigma_n & r_{3n}\sigma_3\sigma_n & \dots & \sigma_n^2 \end{array} \right| = 0$$

Dividing  $C_1, C_2, \dots, C_n$  by  $\sigma_1, \sigma_2, \dots, \sigma_n$  respectively and  $R_2, R_3, \dots, R_n$  by  $\sigma_2, \sigma_3, \dots, \sigma_n$  respectively, we get

$$\left| \begin{array}{ccccc} \frac{X_1}{\sigma_1} & \frac{X_2}{\sigma_2} & \frac{X_3}{\sigma_3} & \dots & \frac{X_n}{\sigma_n} \\ r_{12} & 1 & r_{32} & \dots & r_{2n} \\ r_{13} & r_{23} & 1 & \dots & r_{3n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ r_{1n} & r_{2n} & r_{3n} & \dots & 1 \end{array} \right| = 0 \quad \boxed{\dots(10.37)}$$

If we write

$$\omega = \begin{vmatrix} 1 & r_{12} & r_{13} & \dots & r_{1n} \\ r_{21} & 1 & r_{23} & \dots & r_{2n} \\ r_{31} & r_{32} & 1 & \dots & r_{3n} \\ \vdots & \vdots & \vdots & & \vdots \\ r_{n1} & r_{n2} & r_{n3} & \dots & 1 \end{vmatrix} \quad \dots(10.38)$$

and  $\omega_{ij}$  is the cofactor of the element in the  $i$ th row and  $j$ th column of  $\omega$ , we get from (10.37)

$$\frac{X_1}{\sigma_1} \cdot \omega_{11} + \frac{X_2}{\sigma_2} \omega_{12} + \frac{X_3}{\sigma_3} \omega_{13} + \dots + \frac{X_n}{\sigma_n} \omega_{1n} = 0 \quad \dots(10.39)$$

as the required equation of the plane of regression of  $X_1$  on  $X_2, X_3, \dots, X_n$ .

Equation (10.39) can be re-written as

$$X_1 = -\frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}} X_2 - \frac{\sigma_1}{\sigma_3} \cdot \frac{\omega_{13}}{\omega_{11}} X_3 - \dots - \frac{\sigma_1}{\sigma_n} \cdot \frac{\omega_{1n}}{\omega_{11}} X_n \quad \dots(10.39a)$$

Comparing (10.39a) with (10.35), we get

$$\left. \begin{aligned} b_{12.34\dots n} &= -\frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}}, \\ b_{13.24\dots n} &= -\frac{\sigma_1}{\sigma_3} \cdot \frac{\omega_{13}}{\omega_{11}}, \\ &\vdots &&\vdots \\ b_{1n.23\dots(n-1)} &= -\frac{\sigma_1}{\sigma_n} \cdot \frac{\omega_{1n}}{\omega_{11}} \end{aligned} \right\} \quad \dots(10.40)$$

**Remarks 1.** From the symmetry of the result obtained in (10.40), the equation of the plane of regression of  $X_i$ , (say), on the remaining variables  $X_j$  ( $j \neq i = 1, 2, \dots, n$ ), is given by

$$\frac{X_1}{\sigma_1} \omega_{i1} + \frac{X_2}{\sigma_2} \omega_{i2} + \dots + \frac{X_i}{\sigma_i} \omega_{ii} + \dots + \frac{X_n}{\sigma_n} \omega_{in} = 0 ; i = 1, 2, \dots, n \quad \dots(10.41)$$

**2. We have**

$$b_{12.34\dots n} = -\frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}}$$

$$\text{and} \quad b_{21.34\dots n} = -\frac{\sigma_2}{\sigma_1} \cdot \frac{\omega_{21}}{\omega_{22}}$$

Since each of  $\sigma_1, \sigma_2, \omega_{11}$  and  $\omega_{22}$  is non-negative and  $\omega_{12} = \omega_{21}$ , [c.f. Remarks 3 and 4 to §10.14, page 10.113], the sign of each regression coefficient  $b_{12.34\dots n}$  and  $b_{21.34\dots n}$  depends on  $\omega_{12}$ .

### 10-13. Properties of residuals

**Property 1.** *The sum of the product of any residual of order zero with any other residual of higher order is zero, provided the subscript of the former occurs among the secondary subscripts of the latter.*

The normal equations for estimating  $b$ 's in trivariate and  $n$ -variate distributions, as obtained in equations (10-30a) and (10-36a), are

$$\sum X_2 X_{1 \cdot 23} = 0, \sum X_3 X_{1 \cdot 23} = 0$$

and

$$\sum X_i X_{1 \cdot 23 \dots n} = 0; i = 2, 3, \dots, n$$

respectively. Here  $X_i$ , ( $i = 1, 2, 3, \dots, n$ ) can be regarded as a residual of order zero. Hence the result.

**Property 2.** *The sum of the product of any two residuals in which all the secondary subscripts of the first occur among the secondary subscripts of the second is unaltered if we omit any or all of the secondary subscripts of the first. Conversely, the product sum of any residual of order ' $p$ ' with a residual of order  $p + q$ , the ' $p$ ' subscripts being the same in each case is unaltered by adding to the secondary subscripts of the former any or all the ' $q$ ' additional subscripts of the latter.*

Let us consider

$$\begin{aligned}\sum X_{1 \cdot 2} X_{1 \cdot 23} &= \sum (X_1 - b_{12} X_2) X_{1 \cdot 23} = \sum X_1 X_{1 \cdot 23} - b_{12} \sum X_2 X_{1 \cdot 23} \\ &= \sum X_1 X_{1 \cdot 23} \quad (\text{c.f. Property 1})\end{aligned}$$

$$\begin{aligned}\text{Also } \sum X_{1 \cdot 23}^2 &= \sum X_{1 \cdot 23} X_{1 \cdot 23} = \sum (X_1 - b_{12 \cdot 3} X_2 - b_{13 \cdot 2} X_3) X_{1 \cdot 23} \\ &= \sum X_1 X_{1 \cdot 23} - b_{12 \cdot 3} \sum X_2 X_{1 \cdot 23} - b_{13 \cdot 2} \sum X_3 X_{1 \cdot 23} \\ &= \sum X_1 X_{1 \cdot 23} \quad (\text{c.f. Property 1})\end{aligned}$$

$$\therefore \sum X_{1 \cdot 23}^2 = \sum X_{1 \cdot 2} X_{1 \cdot 23} = \sum X_1 X_{1 \cdot 23}$$

Again  $\sum X_{1 \cdot 34 \dots n} X_{2 \cdot 34 \dots n}$

$$\begin{aligned}&= \sum [(X_1 - b_{13 \cdot 4 \dots n} X_3 - b_{14 \cdot 35 \dots n} X_4 - \dots - b_{1n \cdot 34 \dots (n-1)} X_n) X_{2 \cdot 34 \dots n}] \\ &= \sum X_1 X_{2 \cdot 34 \dots n} \quad (\text{c.f. Property 1})\end{aligned}$$

Hence the property ?

**Property 3.** *The sum of the product of two residuals is zero if all the subscripts (primary as well as secondary) of the one occur among the secondary subscripts of the other, e.g.,*

$$\sum X_{1 \cdot 2} X_{3 \cdot 12} = \sum (X_1 - b_{12} X_2) X_{3 \cdot 12} = \sum X_1 X_{3 \cdot 12} - b_{12} \sum X_2 X_{3 \cdot 12} = 0 \quad (\text{c.f. Property 1})$$

$$\sum X_{2 \cdot 34 \dots n} X_{1 \cdot 23 \dots n}$$

$$\begin{aligned}&= \sum [(X_2 - b_{23 \cdot 4 \dots n} X_3 - b_{24 \cdot 35 \dots n} X_4 - \dots - b_{2n \cdot 34 \dots (n-1)} X_n) X_{1 \cdot 23 \dots n}] \\ &= \sum X_2 X_{1 \cdot 23 \dots n} - b_{23 \cdot 4 \dots n} \sum X_3 X_{1 \cdot 23 \dots n} - b_{24 \cdot 35 \dots n} \sum X_4 X_{1 \cdot 23 \dots n} \\ &\quad \dots - b_{2n \cdot 34 \dots (n-1)} \sum X_n X_{1 \cdot 23 \dots n} \\ &= 0 \quad (\text{c.f. Property 1})\end{aligned}$$

Hence the property 3.

**10.13.1. Variance of the Residual.** Let us consider the plane of regression of  $X_1$  on  $X_2, X_3, \dots, X_n$  viz.,

$$X_1 = b_{123\dots n} X_2 + b_{132\dots n} X_3 + \dots + b_{1n23\dots(n-1)} X_n$$

Since all the  $X_i$ 's are measured from their respective means, we have

$$E(X_i) = 0; i = 1, 2, \dots, n \Rightarrow E(X_{123\dots n}) = 0$$

Hence the variance of the residual is given by

$$\begin{aligned}\sigma^2_{123\dots n} &= \frac{1}{N} \sum [X_{123\dots n} - E(X_{123\dots n})]^2 = \frac{1}{N} \sum X_{123\dots n}^2 \\ &= \frac{1}{N} \sum X_{123\dots n} X_{123\dots n} = \frac{1}{N} \sum X_1 X_{123\dots n},\end{aligned}$$

(c.f. Property 2 § 10.13)

$$\begin{aligned}&= \frac{1}{N} \sum X_1 (X_1 - b_{123\dots n} X_2 - b_{132\dots n} X_3 - \dots - b_{1n23\dots(n-1)} X_n) \\ &= \sigma_1^2 - b_{123\dots n} r_{12} \sigma_1 \sigma_2 - b_{132\dots n} r_{13} \sigma_1 \sigma_3 - \dots - b_{1n23\dots(n-1)} r_{1n} \sigma_1 \sigma_n \\ \Rightarrow \quad \sigma_1^2 - \sigma^2_{123\dots n} &= b_{123\dots n} r_{12} \sigma_1 \sigma_2 - b_{132\dots n} r_{13} \sigma_1 \sigma_3 - \dots \\ &\quad - b_{1n23\dots(n-1)} r_{1n} \sigma_1 \sigma_n \dots \quad (10.42)\end{aligned}$$

Eliminating the  $b$ 's in equations (10.42) and (10.36b), we get

$$\left| \begin{array}{cccc} \sigma_1^2 - \sigma^2_{123\dots n} & r_{12} \sigma_1 \sigma_2 & \dots & r_{1n} \sigma_1 \sigma_n \\ r_{12} \sigma_1 \sigma_2 & \sigma_2^2 & \dots & r_{2n} \sigma_2 \sigma_n \\ \vdots & \vdots & \ddots & \vdots \\ r_{1n} \sigma_1 \sigma_n & r_{2n} \sigma_2 \sigma_n & \dots & \sigma_n^2 \end{array} \right| = 0$$

Dividing  $R_1, R_2, \dots, R_n$ , by  $\sigma_1, \sigma_2, \dots, \sigma_n$  respectively and also  $C_1, C_2, \dots, C_n$  by  $\sigma_1, \sigma_2, \dots, \sigma_n$  respectively, we get

$$\begin{aligned}&\left| \begin{array}{cccc} 1 - \frac{\sigma^2_{123\dots n}}{\sigma_1^2} & r_{12} & \dots & r_{1n} \\ r_{12} & 1 & \dots & r_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1n} & r_{2n} & \dots & 1 \end{array} \right| = 0 \\ \Rightarrow \quad &\left| \begin{array}{cccc} 1 - \frac{\sigma^2_{123\dots n}}{\sigma_1^2} & r_{12} & \dots & r_{1n} \\ r_{12} + 0 & 1 & \dots & r_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1n} + 0 & r_{2n} & \dots & 1 \end{array} \right| = 0\end{aligned}$$

$$\left| \begin{array}{cccc} 1 & r_{12} & \dots & r_{1n} \\ r_{12} & 1 & \dots & r_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1n} & r_{2n} & \dots & 1 \end{array} \right| - \left| \begin{array}{cccc} \frac{\sigma^2_{1,23,\dots,n}}{\sigma_1^2} & r_{12} & \dots & r_{1n} \\ 0 & 1 & \dots & r_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ 0 & r_{2n} & \dots & 1 \end{array} \right| = 0$$

$$\Rightarrow \omega - \frac{\sigma^2_{1,23,\dots,n}}{\sigma_1^2} \omega_{11} = 0$$

$$\therefore \sigma^2_{1,23,\dots,n} = \sigma_1^2 \frac{\omega}{\omega_{11}} \quad \dots(10.43)$$

**Remark.** In a tri-variate distribution,

$$\sigma_{1,23}^2 = \sigma_1^2 \frac{\omega}{\omega_{11}} \quad \dots(10.43a)$$

where  $\omega$  and  $\omega_{11}$  are defined in (10.32).

**10.14. Coefficient of Multiple Correlation.** In a tri-variate distribution in which each of the variables  $X_1$ ,  $X_2$ , and  $X_3$  has  $N$  observations, the multiple correlation coefficient of  $X_1$  on  $X_2$  and  $X_3$ , usually denoted by  $R_{1,23}$ , is the simple correlation coefficient between  $X_1$  and the joint effect of  $X_2$  and  $X_3$  on  $X_1$ . In other words  $R_{1,23}$  is the correlation coefficient between  $X_1$  and its estimated value as given by the plane of regression of  $X_1$  on  $X_2$  and  $X_3$  viz.,

$$e_{1,23} = b_{12,3} X_2 + b_{13,2} X_3$$

$$\text{We have } X_{1,23} = X_1 - b_{12,3} X_2 - b_{13,2} X_3 = X_1 - e_{1,23}$$

$$\Rightarrow e_{1,23} = X_1 - \bar{X}_{1,23}$$

Since  $X_i$ 's are measured from their respective means, we have

$$E(X_{1,23}) = 0 \text{ and } E(e_{1,23}) = 0 \quad (\because E(X_i) = 0; i = 1, 2, 3)$$

By def.,

$$R_{1,23} = \frac{\text{Cov}(X_1, e_{1,23})}{\sqrt{V(X_1) V(e_{1,23})}} \quad \dots(10.44)$$

$$\begin{aligned} \text{Cov}(X_1, e_{1,23}) &= E[(X_1 - E(X_1))(e_{1,23} - E(e_{1,23}))] = E(X_1 e_{1,23}) \\ &= \frac{1}{N} \sum X_1 e_{1,23} = \frac{1}{N} \sum X_1 (X_1 - \bar{X}_{1,23}) \\ &= \frac{1}{N} \sum X_1^2 - \frac{1}{N} \sum X_1 X_{1,23} = \frac{1}{N} \sum X_1^2 - \frac{1}{N} \sum X_{1,23}^2 \\ &= \sigma_1^2 - \sigma_{1,23}^2 \quad (\text{c.f. Property 2, § 10.13}) \end{aligned}$$

$$\begin{aligned} \text{Also } V(e_{1,23}) &= E(e_{1,23}^2) = \frac{1}{N} \sum e_{1,23}^2 = \frac{1}{N} \sum (X_1 - \bar{X}_{1,23})^2 \\ &= \frac{1}{N} \sum (X_1^2 + X_{1,23}^2 - 2 X_1 X_{1,23}) \\ &= \frac{1}{N} \sum X_1^2 + \frac{1}{N} \sum X_{1,23}^2 - \frac{2}{N} \sum X_1 X_{1,23} \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{N} \sum X_1^2 + \frac{1}{N} \sum X_{1-23}^2 - \frac{2}{N} \sum X_{1-23}^2 \\
 &= \sigma_1^2 - \sigma_{1-23}^2 \quad (\text{cf. Property 2, § 10.13}) \\
 \therefore R_{1-23} &= \frac{\sigma_1^2 - \sigma_{1-23}^2}{\sqrt{\sigma_1^2(\sigma_1^2 - \sigma_{1-23}^2)}} \\
 \Rightarrow R_{1-23}^2 &= \frac{\sigma_1^2 - \sigma_{1-23}^2}{\sigma_1^2} = 1 - \frac{\sigma_{1-23}^2}{\sigma_1^2} \\
 \Rightarrow 1 - R_{1-23}^2 &= \frac{\sigma_{1-23}^2}{\sigma_1^2}
 \end{aligned}$$

Using (10.43a), we get

$$1 - R_{1-23}^2 = \frac{\omega}{\omega_{11}} \quad \dots(10.45)$$

where

$$\omega = \begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{vmatrix} = 1 - r_{12}^2 - r_{13}^2 - r_{23}^2 + 2r_{12}r_{13}r_{23} \quad (\text{On simplification}).$$

$$\text{and } \omega_{11} = \begin{vmatrix} 1 & r_{23} \\ r_{32} & 1 \end{vmatrix} = 1 - r_{23}^2$$

Hence from (10.45), we get

$$R_{1-23}^2 = 1 - \frac{\omega}{\omega_{11}} = \frac{r_{12}^2 + r_{13}^2 - 2r_{12}r_{13}r_{23}}{1 - r_{23}^2} \quad \dots(10.45a)$$

This formula expresses the multiple correlation coefficient in terms of the total correlation coefficients between the pairs of variables.

**Generalisation.** In case of  $n$ -variate distribution, the multiple correlation coefficient of  $X_1$  on  $X_2, X_3, \dots, X_n$ , usually denoted by  $R_{1-23\dots n}$ , is the correlation coefficient between  $X_1$  and

$$\begin{aligned}
 e_{1-23\dots n} &= X_1 - X_{1-23\dots n} \\
 \therefore R_{1-23\dots n} &= \frac{\text{Cov}(X_1, e_{1-23\dots n})}{\sqrt{V(X_1) V(e_{1-23\dots n})}} \\
 \text{Cov}(X_1, e_{1-23\dots n}) &= \frac{1}{N} \sum X_1 e_{1-23\dots n} = \frac{1}{N} \sum X_1 (X_1 - X_{1-23\dots n}) \\
 &= \frac{1}{N} \sum X_1^2 - \frac{1}{N} \sum X_1 X_{1-23\dots n} \\
 &= \frac{1}{N} \sum X_1^2 - \frac{1}{N} \sum X_{1-23\dots n}^2 = \sigma_1^2 - \sigma_{1-23\dots n}^2 \quad \dots(*) \\
 V(e_{1-23\dots n}) &= \frac{1}{N} \sum e_{1-23\dots n}^2 = \frac{1}{N} \sum (X_1 - X_{1-23\dots n})^2
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{N} \sum (X_1^2 + X_{1.23...n}^2 - 2X_1 X_{1.23...n}) \\
 &= \frac{1}{N} \sum X_1^2 + \frac{1}{N} \sum X_{1.23...n}^2 - 2 \frac{1}{N} \sum X_1 X_{1.23...n} \\
 &= \frac{1}{N} \sum X_1^2 + \frac{1}{N} \sum X_{1.23...n}^2 - \frac{2}{N} \sum X_{1.23...n}^2 \\
 &= \sigma_1^2 - \sigma_{1.23...n}^2 \\
 \therefore R_{1.23...n} &= \frac{\sigma_1^2 - \sigma_{1.23...n}^2}{\sqrt{\sigma_1^2(\sigma_1^2 - \sigma_{1.23...n}^2)}} = \left( \frac{\sigma_1^2 - \sigma_{1.23...n}^2}{\sigma_1^2} \right)^{1/2} \\
 R_{1.23...n}^2 &= 1 - \frac{\sigma_{1.23...n}^2}{\sigma_1^2} = 1 - \frac{\omega}{\omega_{11}} \quad \dots(10.45c)
 \end{aligned}$$

where  $\omega$  and  $\omega_{11}$  are defined in (10.38).

**Remarks 1.** It may be pointed out here that multiple correlation coefficient can never be negative, because from (\*) and (\*\*), we get

$$\text{Cov}(X_1, e_{1.23...n}) = \sigma_1^2 - \sigma_{1.23...n}^2 = \text{Var}(e_{1.23...n}) \geq 0$$

Since the sign of  $R_{1.23...n}$  depends upon the covariance term  $\text{Cov}(X_1, e_{1.23...n})$ , we conclude that  $R_{1.23...n} \geq 0$ .

2. Since  $R_{1.23...n}^2 \geq 0$ , we have :

$$1 - \frac{\omega}{\omega_{11}} \geq 0 \Rightarrow \omega \leq \omega_{11} \quad \dots(10.45d)$$

$$3. \text{ Also, } R_{1.23...n}^2 \leq 1 \Rightarrow 1 - \frac{\omega}{\omega_{11}} \leq 1$$

$$\Rightarrow 0 \leq \frac{\omega}{\omega_{11}} \Rightarrow \frac{\omega}{\omega_{11}} \geq 0 \Rightarrow \omega \geq 0 \quad \dots(10.45e)$$

From the above results, we get

$$\omega_{11} \geq \omega \geq 0 \quad \dots(10.45f)$$

In general, we have

$$\omega_{ii} \geq 0; i = 1, 2, \dots, n$$

4. Since  $\omega$  is symmetric in  $r_{ij}$ 's, we have

$$\omega_{ij} = \omega_{ji}; i \neq j = 1, 2, \dots, n \quad \dots(10.45g)$$

#### 10.14.1. Properties of Multiple Correlation Coefficient

1. Multiple correlation co-efficient measures the closeness of the association between the observed values and the expected values of a variable obtained from the multiple linear regression of that variable on other variables.

2. Multiple correlation coefficient between observed values and expected values, when the expected values are calculated from a linear relation of the variables determined by the method of least squares, is always greater than that where expected values are calculated from any other linear combination of the variables.

3. Since  $R_{1.23}$  is the simple correlation between  $X_1$  and  $e_{1.23}$ , it must lie between  $-1$  and  $+1$ . But as seen in Remark 1 above,  $R_{1.23}$  is a non-negative quantity and we conclude that  $0 \leq R_{1.23} \leq 1$ .

4. If  $R_{1.23} = 1$ , then association is perfect and all the regression residuals are zero, and as such  $\sigma^2_{1.23} = 0$ . In this case, since  $X_1 = e_{1.23}$ , the predicted value of  $X_1$ , the multiple linear regression equation of  $X_1$  on  $X_2$  and  $X_3$  may be said to be a perfect prediction formula.

5. If  $R_{1.23} = 0$ , then all total and partial correlations involving  $X_1$  are zero. [See Example 10.37]. So  $X_1$  is completely uncorrelated with all the other variables in this case and the multiple regression equation fails to throw any light on the value of  $X_1$  when  $X_2$  and  $X_3$  are known.

6.  $R_{1.23}$  is not less than any total correlation coefficient, i.e.,

$$R_{1.23} \geq r_{12}, r_{13}, r_{23}$$

**10.15. Coefficient of Partial Correlation.** Sometimes the correlation between two variables  $X_1$  and  $X_2$  may be partly due to the correlation of a third variable,  $X_3$  with both  $X_1$  and  $X_2$ . In such a situation, one may want to know what the correlation between  $X_1$  and  $X_2$  would be if the effect of  $X_3$  on each of  $X_1$  and  $X_2$  were eliminated. This correlation is called the *partial correlation* and the correlation coefficient between  $X_1$  and  $X_2$  after the linear effect of  $X_3$  on each of them has been eliminated is called the *partial correlation coefficient*.

The residual  $X_{1.3} = X_1 - b_{13}X_3$ , may be regarded as that part of the variable  $X_1$  which remains after the linear effect of  $X_3$  has been eliminated. Similarly, the residual  $X_{2.3}$  may be interpreted as the part of the variable  $X_2$  obtained after eliminating the linear effect of  $X_3$ . Thus the partial correlation coefficient between  $X_1$  and  $X_2$ , usually denoted by  $r_{12.3}$ , is given by

$$r_{12.3} = \frac{\text{Cov}(X_{1.3}, X_{2.3})}{\sqrt{\text{Var}(X_{1.3}) \text{Var}(X_{2.3})}} \quad \dots(10.46)$$

We have

$$\begin{aligned} \text{Cov}(X_{1.3}, X_{2.3}) &= \frac{1}{N} \sum X_{1.3} X_{2.3} = \frac{1}{N} \sum X_1 X_{2.3} \\ &= \frac{1}{N} \sum X_1 (X_2 - b_{23} X_3) = \frac{1}{N} \sum X_1 X_2 - b_{23} \frac{1}{N} \sum X_1 X_3 \end{aligned}$$

$$\begin{aligned} &= r_{12} \sigma_1 \sigma_2 - r_{23} \frac{\sigma_2}{\sigma_3} \cdot (r_{13} \sigma_1 \sigma_3) \\ &= \sigma_1 \sigma_2 (r_{12} - r_{13} r_{23}) \end{aligned}$$

$$\text{Also } V(X_{1.3}) = \frac{1}{N} \sum X_{1.3}^2 = \frac{1}{N} \sum X_{1.3} X_{1.3}$$

$$= \frac{1}{N} \sum X_1 X_{1.3} = \frac{1}{N} \sum X_1 (X_1 - b_{13} X_3)$$

$$= \frac{1}{N} \sum X_1^2 - b_{13} \cdot \frac{1}{N} \sum X_1 X_3$$

$$= \sigma_1^2 - r_{13} \frac{\sigma_1}{\sigma_3} r_{13} \sigma_1 \sigma_3$$

$$= \sigma_1^2 (1 - r_{13}^2)$$

Similarly, we shall get

$$V(X_{2,3}) = \sigma_2^2 (1 - r_{23}^2)$$

Hence

$$r_{12,3} = \frac{\sigma_1 \sigma_2 (r_{12} - r_{13} r_{23})}{\sqrt{\sigma_1^2 (1 - r_{13}^2) \sigma_2^2 (1 - r_{23}^2)}} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{(1 - r_{13}^2) (1 - r_{23}^2)}} \quad \dots(10-46a)$$

Aliter. We have

$$0 = \sum X_{2,3} X_{1,3}$$

$$= \sum X_{2,3} (X_1 - b_{12,3} X_2 - b_{13,2} X_3)$$

$$= \sum X_1 X_{2,3} - b_{12,3} \sum X_{2,3} X_2 - b_{13,2} \sum X_{2,3} X_3$$

$$= \sum X_{1,3} X_{2,3} - b_{12,3} \sum X_{2,3} X_{2,3}$$

$$\therefore b_{12,3} = \frac{\sum X_{1,3} X_{2,3}}{\sum X_{2,3}^2}$$

From this it follows that  $b_{12,3}$  is coefficient of regression of  $X_{1,3}$  on  $X_{2,3}$ .

Similarly,  $b_{21,3}$  is the coefficient of regression of  $X_{2,3}$  on  $X_{1,3}$ .

Since correlation coefficient is the geometric mean between regression coefficients, we have

$$r^2_{12,3} = b_{12,3} \times b_{21,3}$$

But by def.,

$$b_{12,3} = - \frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}} \quad \text{and} \quad b_{21,3} = - \frac{\sigma_2}{\sigma_1} \cdot \frac{\omega_{21}}{\omega_{22}}$$

$$\therefore r^2_{12,3} = \left( - \frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}} \right) \left( - \frac{\sigma_2}{\sigma_1} \cdot \frac{\omega_{21}}{\omega_{22}} \right) = \frac{\omega_{12}^2}{\omega_{11} \omega_{22}}$$

$$(\because \omega_{12} = \omega_{21})$$

$$\Rightarrow r_{12,3} = - \frac{\omega_{12}}{\sqrt{\omega_{11} \omega_{22}}} ,$$

the negative sign being taken since the sign of regression coefficients is the same as that of  $(-\omega_{12})$ .

Substituting the values of  $\omega_{12}$ ,  $\omega_{11}$  and  $\omega_{22}$  from (10-32), we get

$$r_{12,3} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}}$$

**Remarks 1.** The expressions for  $r_{13,2}$  and  $r_{23,1}$  can be similarly obtained, to give

$$r_{13 \cdot 2} = \frac{r_{13} - r_{12} r_{32}}{\sqrt{(1 - r_{12}^2)(1 - r_{32}^2)}} \quad \text{and} \quad r_{23 \cdot 1} = \frac{r_{23} - r_{21} r_{31}}{\sqrt{(1 - r_{21}^2)(1 - r_{31}^2)}}$$

2. If  $r_{12 \cdot 3} = 0$ , we have then  $r_{12} = r_{13} r_{23}$ , it means that  $r_{12}$  will not be zero if  $X_3$  is correlated with both  $X_1$  and  $X_2$ . Thus, although  $X_1$  and  $X_2$  may be uncorrelated when effect of  $X_3$  is eliminated, yet  $X_1$  and  $X_2$  may appear to be correlated because they carry the effect of  $X_3$  on them.

3. Partial correlation coefficient helps in deciding whether to include or not an additional independent variable in regression analysis.

4. We know that  $\sigma_1^2(1 - r_{12}^2)$  and  $\sigma_1^2(1 - r_{13}^2)$  are the residual variances if  $X_1$  is estimated from  $X_2$  and  $X_3$  individually, while  $\sigma_1^2(1 - R_{1 \cdot 23}^2)$  is the residual variance if  $X_1$  is estimated from  $X_2$  and  $X_3$  taken together. So from the above remark and  $R_{1 \cdot 23}^2 \geq r_{12}^2$  and  $r_{13}^2$ , it follows that inclusion of an additional variable can only reduce the residual variance. Now inclusion of  $X_3$  when  $X_2$  has already been taken for predicting  $X_1$ , is worthwhile only when the resultant reduction in the residual variance is substantial. This will be the case when  $r_{13 \cdot 2}$  is sufficiently large. Thus in this respect partial correlation coefficient has its significance in regression analysis.

**10-15-1. Generalisation.** In the case of  $n$  variables  $X_1, X_2, \dots, X_n$  the partial correlation coefficient  $r_{12 \cdot 34 \dots n}$  between  $X_1$  and  $X_2$  (after the linear effect of  $X_3, X_4, \dots, X_n$  on them has been eliminated), is given by

$$r_{12 \cdot 34 \dots n}^2 = b_{12 \cdot 34 \dots n} \times b_{21 \cdot 34 \dots n}$$

But, we have

$$\begin{aligned} \text{and } b_{12 \cdot 34 \dots n} &= -\frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}} \\ b_{21 \cdot 34 \dots n} &= -\frac{\sigma_2}{\sigma_1} \cdot \frac{\omega_{21}}{\omega_{22}} \end{aligned} \left. \right\} [\text{cf. Equation (10-40)}]$$

$$\therefore r_{12 \cdot 34 \dots n}^2 = \left( -\frac{\sigma_1}{\sigma_2} \cdot \frac{\omega_{12}}{\omega_{11}} \right) \left( -\frac{\sigma_2}{\sigma_1} \cdot \frac{\omega_{21}}{\omega_{22}} \right) = \frac{\omega_{12}^2}{\omega_{11} \omega_{22}}$$

$$\Rightarrow r_{12 \cdot 34 \dots n} = -\frac{\omega_{12}}{\sqrt{\omega_{11} \omega_{22}}} \quad (10-46b)$$

negative sign being taken since the sign of the regression coefficient is same as that of  $(-\omega_{12})$ .

#### 10-16. Multiple Correlation in Terms of Total and Partial Correlations.

$$1 - R_{1 \cdot 23}^2 = (1 - r_{12}^2)(1 - r_{13 \cdot 2}^2) \quad \dots (10-46c)$$

**Proof.** We have

$$\begin{aligned} 1 - R_{1 \cdot 23}^2 &= 1 - \frac{r_{12}^2 + r_{13}^2 - 2r_{12} r_{13} r_{23}}{1 - r_{23}^2} \\ &= \frac{1 - r_{23}^2 - r_{12}^2 - r_{13}^2 + 2r_{12} r_{13} r_{23}}{1 - r_{23}^2} \end{aligned}$$

Also

$$1 - r_{13.2}^2 = 1 - \frac{(r_{13} - r_{12} r_{23})^2}{(1 - r_{12}^2)(1 - r_{23}^2)} = \frac{1 - r_{12}^2 - r_{23}^2 - r_{13}^2 + 2r_{12}r_{13}}{(1 - r_{12}^2)(1 - r_{23}^2)}$$

Hence the result.

**Theorem.** Any standard deviation of order 'p' may be expressed in terms of a standard deviation of order  $(p - 1)$  and a partial correlation coefficient of order  $(p - 1)$ .

**Proof.** Let us consider the sum :

$$\begin{aligned} \sum X_{1.23...n}^2 &= \sum X_{1.23...n} X_{1.23...n} \\ &= \sum [X_{1.23...(n-1)} X_{1.23...n} \\ &\quad (c.f. \text{ Property 2, § 10.13}) \\ &= \sum [X_{1.23...(n-1)} (X_1 - b_{1.23...n} X_2 - \dots - b_{1(n-1).23...n} X_{n-1} \\ &\quad - b_{1n.23...(n-1)} X_n)] \\ &= \sum X_{1.23.(n-1)} X_1 - b_{1n.23...(n-1)} \sum X_{1.23...(n-1)} X_n \\ &\quad (c.f. \text{ Property 2 § 10.13}) \\ &= \sum X_{1.23...(n-1)}^2 - b_{1n.23...(n-1)} \sum X_{1.23...(n-1)} X_{n.23...(n-1)} \end{aligned}$$

Dividing both sides by  $N$  (total number of observations), we get  $\sigma_{1.23...n}^2 = \sigma_{1.23...(n-1)}^2 - b_{1n.23...(n-1)} \text{Cov}(X_{1.23...(n-1)}, X_{n.23...(n-1)})$

The regression coefficient of  $X_{n.23...(n-1)}$  on  $X_{1.23...(n-1)}$  is given by

$$\begin{aligned} b_{n1.23...(n-1)} &= \frac{\text{Cov}(X_{1.23...(n-1)}, X_{n.23...(n-1)})}{\sigma_{1.23...(n-1)}^2}, \\ \therefore \sigma_{1.23...n}^2 &= \sigma_{1.23...(n-1)}^2 [1 - b_{1n.23...(n-1)} b_{n1.23...(n-1)}] \\ &= \sigma_{1.23...(n-1)}^2 [1 - r_{1n.23...(n-1)}^2], \quad \dots(10.47) \end{aligned}$$

a formula which expresses the standard deviation of order  $(n - 1)$  in terms of standard deviation of order  $(n - 2)$  and partial correlation coefficient of order  $(n - 2)$ . If we take  $p = (n - 1)$ , the theorem is established.

**Cor. 1.** From (10.47), we have

$$\sigma_{1.23...(n-1)}^2 = \sigma_{1.23...(n-2)}^2 (1 - r_{1(n-1).23...(n-2)}^2) \quad \dots(10.47a)$$

and so on. Thus the repeated application of (10.47) gives

$$\sigma_{1.23...n}^2 = \sigma_1^2 (1 - r_{12}^2) (1 - r_{13.2}^2) (1 - r_{14.32}^2) \dots (1 - r_{1n.23...(n-1)}^2) \quad \dots(10.47b)$$

Since partial correlation coefficients cannot exceed unity numerically, we get from (10.47), (10.47a), and so on,

$$\left. \begin{array}{l} \sigma_{1 \cdot 23 \dots n}^2 \leq \sigma_{1 \cdot 23 \dots (n-1)}^2 \\ \sigma_{1 \cdot 23 \dots (n-1)}^2 \leq \sigma_{1 \cdot 23 \dots (n-2)}^2 \\ | \qquad \qquad | \\ \sigma_{1 \cdot 23}^2 \leq \sigma_{1 \cdot 2}^2 \\ \sigma_{1 \cdot 2}^2 \leq \sigma_1^2 \end{array} \right\} \quad \sigma_1 \geq \sigma_{1 \cdot 2} \geq \sigma_{1 \cdot 23} \geq \dots \geq \sigma_{1 \cdot 23 \dots n} \quad \dots(10-47c)$$

**Cor. 2.** Also, we have

**Cor. 2.** Also, we have

$$\sigma^2_{1,23\dots n} = \sigma_1^2(1 - R^2_{1,23\dots n})$$

On using (10-47b), we get

$$1 - R_{1,2,3,\dots,n}^2 = (1 - r_{12}^2)(1 - r_{13,2}^2)\dots(1 - r_{1,n-3,\dots,(n-1)}^2) \quad \dots(10.47d)$$

This is the generalisation of the result obtained in (10-46c).

Since  $|r_{ij,(s)}| \leq 1$ ;  $s = 0, 1, 2, \dots, (n - 1)$ ,

where  $r_{ij,(s)}$  is a partial correlation coefficient of order  $s$ . we get from (10-47d)

$$1 - R^2_{1,23\dots n} \leq 1 - r_{12}^2$$

$$1 - R^2_{1,23,\dots,n} \leq 1 - r^2_{13,2},$$

and so on.

$$\text{i.e.,} \quad R^2_{1,2,3,\dots,n} \geq r_{12}^2, r_{13,2}^2, \dots, r_{1,n,2,3,\dots,(n-1)}^2 \quad \dots(10.47e)$$

Since  $R_{1,23\dots n}$  is symmetric in its secondary subscripts, we have

$$\left. \begin{array}{l} R^2_{1 \cdot 2 \cdot 3 \dots n} \geq r_{1i}^2, \quad (i = 2, 3, \dots, n) \\ R^2_{1 \cdot 2 \cdot 3 \dots n} \geq r_{1ij} \quad (i \neq j = 2, 3, \dots, n) \end{array} \right\} \quad \dots(10 \cdot 47f)$$

and so on

**Example 10.33.** From the data relating to the yield of dry bark ( $X_1$ ), height ( $X_2$ ) and girth  $X_3$  for 18 cinchona plants the following correlation coefficients were obtained :

$$r_{12} = 0.77, r_{13} = 0.72 \text{ and } r_{23} = 0.52$$

Find the the partial correlation coefficient  $r_{12.3}$  and multiple correlation coefficient  $R_{1.23}$ .

**Solution.**

$$r_{12.3} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}} = \frac{0.77 - 0.72 \times 0.52}{\sqrt{[1 - (0.72)^2][1 - (0.52)^2]}} = 0.62$$

$$\begin{aligned} R_{1.23}^2 &= \frac{r_{12}^2 + r_{13}^2 - 2r_{12}r_{13}r_{23}}{1 - r_{23}^2} \\ &= \frac{(0.77)^2 + (0.72)^2 - 2 \times 0.77 \times 0.72 \times 0.52}{1 - (0.52)^2} = 0.7334 \end{aligned}$$

$$\therefore R_{1.23} = \pm 0.8564$$

(since multiple correlation coefficient is non-negative).

**Example 10.34.** In a trivariate distribution :

$$\sigma_1 = 2, \sigma_2 = \sigma_3 = 3, r_{12} = 0.7, r_{23} = r_{31} = 0.5.$$

Find (i)  $r_{23.1}$ , (ii)  $R_{1.23}$ , (iii)  $b_{12.3}$ ,  $b_{13.2}$  and (iv)  $\sigma_{1.23}$ .

**Solution.** We have

$$(i) r_{23.1} = \frac{r_{23} - r_{21}r_{31}}{\sqrt{(1 - r_{21}^2)(1 - r_{31}^2)}} = \frac{0.5 - (0.7)(0.5)}{\sqrt{(1 - 0.49)(1 - 0.25)}} = 0.2425$$

$$\begin{aligned} (ii) R_{1.23}^2 &= \frac{r_{12}^2 + r_{13}^2 - 2r_{12}r_{13}r_{23}}{1 - r_{23}^2} \\ &= \frac{0.49 + 0.25 - 2(0.7)(0.5)(0.5)}{1 - 0.25} = 0.52 \end{aligned}$$

$$\therefore R_{1.23} = + 0.7211$$

$$(iii) b_{12.3} = r_{12.3} \frac{\sigma_{1.3}}{\sigma_{2.3}} \text{ and } b_{13.2} = r_{13.2} \frac{\sigma_{1.2}}{\sigma_{3.2}}$$

$$r_{12.3} = \frac{r_{12} - r_{13}r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}} = 0.6, r_{13.2} = \frac{r_{13} - r_{12}r_{32}}{\sqrt{(1 - r_{12}^2)(1 - r_{32}^2)}} = 0.2425$$

$$\sigma_{1:3} = \sigma_1 \sqrt{(1 - r_{13}^2)} = 2 \sqrt{(1 - 0.25)} = 1.7320$$

$$\sigma_{2:3} = \sigma_2 \sqrt{(1 - r_{23}^2)} = 3 \sqrt{(1 - 0.25)} = 2.5980$$

$$\sigma_{1:2} = \sigma_1 \sqrt{(1 - r_{12}^2)} = 2 \sqrt{(1 - 0.49)} = 1.4282$$

$$\sigma_{3:2} = \sigma_3 \sqrt{(1 - r_{32}^2)} = 3 \sqrt{(1 - 0.25)} = 2.5980$$

Hence  $b_{12:3} = 0.4$  and  $b_{13:2} = 0.1333$

$$(iv) \quad \sigma_{1:23} = \sigma_1 \sqrt{\frac{\omega}{\omega_{11}}}$$

$$\text{where } \omega = \begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{vmatrix} = 1 - r_{12}^2 - r_{13}^2 - r_{23}^2 + 2r_{12}r_{13}r_{23} = 0.36$$

$$\text{and } \omega_{11} = \begin{vmatrix} 1 & r_{23} \\ r_{32} & 1 \end{vmatrix} = 1 - r_{23}^2 = 1 - 0.25 = 0.75$$

$$\therefore \sigma_{1:23} = 2 \times \sqrt{0.48} = 2 \times 0.6928 = 1.3856$$

**Example 10.35.** Find the regression equation of  $X_1$  on  $X_2$  and  $X_3$  given the following results :—

Trait	Mean	Standard deviation	$r_{12}$	$r_{23}$	$r_{31}$
$X_1$	28.02	4.42	+ 0.80	—	—
$X_2$	4.91	1.10	—	-0.56	—
$X_3$	594	85	—	—	-0.40

where  $X_1$  = Seed per acre;  $X_2$  = Rainfall in inches

$X_3$  = Accumulated temperature above 42°F.

**Solution.** Regression equation of  $X_1$  on  $X_2$  and  $X_3$  is given by

$$(X_1 - \bar{X}_1) \frac{\omega_{11}}{\sigma_1} + (X_2 - \bar{X}_2) \frac{\omega_{12}}{\sigma_2} + (X_3 - \bar{X}_3) \frac{\omega_{13}}{\sigma_3} = 0$$

$$\text{where } \omega = \begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{vmatrix}$$

$$\omega_{11} = \begin{vmatrix} 1 & r_{23} \\ r_{32} & 1 \end{vmatrix} = 1 - r_{23}^2 = 1 - (-0.56)^2 = 0.686$$

$$\omega_{12} = \begin{vmatrix} r_{21} & r_{23} \\ r_{31} & 1 \end{vmatrix} = r_{13}r_{23} - r_{21} = -0.576$$

$$\omega_{13} = r_{23}r_{12} - r_{13} = (-0.56)(0.80) - (-0.40) = -0.048$$

∴ Required equation of plane of regression of  $X_1$  on  $X_2$  and  $X_3$  is given by

$$\frac{0.686}{4.42} (X_1 - 28.02) + \frac{(-0.576)}{1.10} (X_2 - 4.91) + \frac{(-0.048)}{85.00} (X_3 - 594) = 0$$

**Example 10-36.** Five hundred students were examined in three subjects I, II and III, each subject carrying 100 marks. A student getting 120 or more but less than 150 marks was put in pass class. A student getting 150 or more but less than 180 marks was put in second class and a student getting 180 or more marks was put in the first class. The following marks were obtained :

	I	II	III
Mean :	35.8	52.4	48.8
S.D. :	4.2	5.3	6.1
Correlation :	$r_{12} = 0.6$ ,	$r_{13} = 0.7$	$r_{23} = 0.8$

- (i) Find the number of students in each of the three classes.  
(ii) Find the total number of students with total marks lying between 120 and 190.

(iii) Find the probability that a student gets more than 240 marks.

(iv) What should be the correlation between marks in subjects I and II among students who scored equal marks in subject III ?

(v) If  $r_{23}$  was not known, obtain the limits within which it may lie from the values of  $r_{12}$  and  $r_{13}$  (ignoring sampling errors).

**Solution.** If  $Z$  denotes the total marks of the students in the three subjects and  $X_1, X_2, X_3$  the total marks of the students in subjects I, II and III respectively, then

$$\begin{aligned} Z &= X_1 + X_2 + X_3 \\ \therefore E(Z) &= E(X_1) + E(X_2) + E(X_3) = 35.8 + 52.4 + 48.8 = 137 \\ V(Z) &= V(X_1) + V(X_2) + V(X_3) \\ &\quad + 2[\text{Cov}(X_1, X_2) + \text{Cov}(X_2, X_3) + \text{Cov}(X_3, X_1)] \\ &= 17.64 + 28.09 + 37.21 + 26.712 + 35.868 + 51.728 \\ &= 197.248 \quad [\text{Using } \text{Cov}(X_i, X_j) = r_{ij} \sigma_i \sigma_j] \\ \Rightarrow \sigma_Z^2 &= 197.248 \text{ or } \sigma_Z = 14.045 \\ \text{Now } \xi &= \frac{Z - E(Z)}{\sigma_Z} \sim N(0, 1) \end{aligned}$$

Z	$\xi = \frac{Z - 137}{14.045}$	$p = \int_{-\infty}^{\xi} p(\xi) d\xi$	Class	Area under the curve in this class (A)	Frequency $500 \times (A)$
120 -	1.21050	0.11314	120 - 150	0.70937	354.685
150	0.92567	0.82251	150 - 180	0.17639	88.195
180	3.06180	0.99890	180 -	0.00102	0.510
190	3.77400	0.99992	120 - 190	0.88678	443.390
240	7.33410	1.00000	240 -	0.00000	0.000

- (i) The number of students in first, second and third class respectively are 355, 88 and 0 (approx.).  
(ii) Total number of students with total marks between 120 and 190 is 443.  
(iii) Probability that a student gets more than 240 marks is zero.  
(iv) The correlation coefficient between marks in subjects I and II of the students who secured equal marks in subject III is  $r_{123}$  and is given by

$$r_{12 \cdot 3} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}} = \frac{0.04}{\sqrt{(1 - 0.49)(1 - 0.64)}} = 0.0934$$

(v) We have .

$$\begin{aligned} r_{12 \cdot 3}^2 &= \frac{(r_{12} - r_{13} r_{23})^2}{(1 - r_{13}^2)(1 - r_{23}^2)} \leq 1 \\ \therefore \quad &\frac{(0.6 - 0.7a)^2}{(1 - 0.49)(1 - a^2)} \leq 1, \text{ where } a = r_{23}. \\ \Rightarrow \quad &0.36 + 0.49a^2 - 0.84a \leq 0.51(1 - a^2) \\ \Rightarrow \quad &a^2 - 0.84a - 0.15 \leq 0 \end{aligned}$$

Thus 'a' lies between the roots of the equation :

$$a^2 - 0.84a - 0.15 = 0,$$

which are 0.99 and -0.15.

Hence  $r_{23}$  should lie between -0.15 and 0.99.

**Example 10-37.** Show that

$$1 - R_{1 \cdot 23}^2 = (1 - r_{12}^2)(1 - r_{13 \cdot 2}^2)$$

Deduce that

$$(i) R_{1 \cdot 23} \geq r_{12}, \quad (ii) R_{1 \cdot 23}^2 = r_{12}^2 + r_{13}^2, \text{ if } r_{23} = 0$$

$$(iii) 1 - R_{1 \cdot 23}^2 = \frac{(1 - \rho)(1 + 2\rho)}{(1 + \rho)}, \text{ provided all the coefficients of zero order are equal to } \rho.$$

(iv) If  $R_{1 \cdot 23} = 0$ ,  $X_1$  is uncorrelated with any of other variables, i.e.,  $r_{12} = r_{13} = 0$ . [Delhi Univ. B.Sc. (Stat. Hons.), 1989]

**Solution.** (i) Since  $|r_{13 \cdot 2}| \leq 1$ , we have from (10-46c)

$$1 - R_{1 \cdot 23}^2 \leq 1 - r_{12}^2 \Rightarrow R_{1 \cdot 23} \geq r_{12}$$

(ii) We have

$$r_{13 \cdot 2} = \frac{r_{13} - r_{12} r_{32}}{\sqrt{(1 - r_{12}^2)(1 - r_{32}^2)}} = \frac{r_{13}}{\sqrt{1 - r_{12}^2}}. \quad (\text{if } r_{23} = 0)$$

∴ From (10-46c), we get

$$1 - R_{1 \cdot 23}^2 = (1 - r_{12}^2) \left[ 1 - \frac{r_{13}^2}{1 - r_{12}^2} \right] = 1 - r_{12}^2 - r_{13}^2$$

Hence  $R_{1 \cdot 23}^2 = r_{12}^2 + r_{13}^2$ , if  $r_{23} = 0$ .

(iii) Here, we are given that  $r_{12} = r_{13} = r_{23} = \rho$

$$\therefore r_{13 \cdot 2} = \frac{\rho - \rho^2}{\sqrt{(1 - \rho^2)(1 - \rho^2)}} = \frac{\rho(1 - \rho)}{(1 - \rho^2)} = \frac{\rho}{1 + \rho}$$

Hence from (10-46c), we have

$$1 - R_{1 \cdot 23}^2 = (1 - \rho^2) \left[ 1 - \frac{\rho^2}{(1 + \rho)^2} \right] = \frac{(1 - \rho)(1 + 2\rho)}{(1 + \rho)}$$

(iv) If  $R_{1 \cdot 23} = 0$ , (10-46c) gives

$$1 = (1 - r_{12}^2)(1 - r_{13 \cdot 2}^2)$$

...(\*)

Since  $0 \leq r_{12}^2 \leq 1$  and  $0 \leq r_{13-2}^2 \leq 1$ , (\*) will hold if and only if

$$r_{12} = 0 \quad \text{and} \quad r_{13-2} = 0$$

$$\text{Now } r_{13-2} = 0 \Rightarrow \frac{r_{13}-r_{12}r_{32}}{\sqrt{(1-r_{12}^2)(1-r_{32}^2)}} = 0$$

$$\Rightarrow \frac{r_{13}}{\sqrt{1-r_{32}^2}} = 0 \quad (\because r_{12} = 0)$$

$$\Rightarrow r_{13} = 0$$

Thus if  $R_{1-23} = 0$ , then  $r_{13} = r_{12} = 0$ , i.e.,  $X_1$  is uncorrelated with  $X_2$  and  $X_3$ .

**Example 10-38.** Show that the correlation coefficient between the residuals  $X_{1-23}$  and  $X_{2-13}$  is equal and opposite to that between  $X_{1-3}$  and  $X_{2-3}$ .

[Poona Univ. B.Sc., 1991]

**Solution.** The correlation coefficient between  $X_{1-23}$  and  $X_{2-13}$  is given by

$$\begin{aligned}\frac{\text{Cov}(X_{1-23}, X_{2-13})}{\sigma_{1-23} \sigma_{2-13}} &= \frac{\sum X_{1-23} X_{2-13}}{N \sigma_{1-23} \sigma_{2-13}} = \frac{\frac{1}{N} \sum X_{2-13} (X_1 - b_{12-3} X_2 - b_{13-2} X_3)}{\sigma_{1-23} \sigma_{2-13}} \\ &= -b_{12-3} \frac{\sum X_{2-13} X_2}{N \sigma_{1-23} \sigma_{2-13}} \quad (\text{c.f. Property 1, § 10-13}) \\ &= -b_{12-3} \frac{\sum X_{2-13}^2}{N \sigma_{1-23} \sigma_{2-13}} \quad (\text{c.f. Property 2, § 10-13}) \\ &= -b_{12-3} \frac{\sigma_{2-13}}{\sigma_{1-23}} = -b_{12-3} \frac{(\sigma_2 \sqrt{\omega/\omega_{22}})}{(\sigma_1 \sqrt{\omega/\omega_{11}})}\end{aligned}$$

$$\text{where } \omega = \begin{vmatrix} 1 & r_{12} & r_{13} \\ r_{21} & 1 & r_{23} \\ r_{31} & r_{32} & 1 \end{vmatrix}$$

$$\omega_{11} = \begin{vmatrix} 1 & r_{23} \\ r_{32} & 1 \end{vmatrix} = 1 - r_{23}^2 \quad \text{and} \quad \omega_{22} = \begin{vmatrix} 1 & r_{13} \\ r_{31} & 1 \end{vmatrix} = 1 - r_{13}^2$$

$$\therefore r(X_{1-23}, X_{2-13}) = -b_{12-3} \frac{\sigma_2}{\sigma_1} \cdot \sqrt{\frac{1-r_{23}^2}{1-r_{13}^2}} = -b_{12-3} \frac{\sigma_{2-3}}{\sigma_{1-3}}$$

[since  $\sigma_{2-3}^2 = \sigma_2^2(1-r_{23}^2)$  and  $\sigma_{1-3}^2 = \sigma_1^2(1-r_{13}^2)$ ]

$$\therefore r(X_{1-23}, X_{2-13}) = -\frac{\text{Cov}(X_{1-3}, X_{2-3})}{\sigma_{2-3}^2} \cdot \frac{\sigma_{2-3}}{\sigma_{1-3}}$$

$$= -\frac{\text{Cov}(X_{1-3}, X_{2-3})}{\sigma_{2-3} \sigma_{1-3}} = -r(X_{1-3}, X_{2-3})$$

Hence the result.

**Example 10·42.** If  $r_{12}$  and  $r_{13}$  are given, show that  $r_{23}$  must lie in the range :  $r_{12}r_{13} \pm (1 - r_{12}^2 - r_{13}^2 + r_{12}^2r_{13}^2)^{1/2}$

If  $r_{12} = k$  and  $r_{13} = -k$ , show that  $r_{23}$  will lie between  $-1$  and  $1 - 2k^2$ .

[*Sardar Patel Univ. B.Sc. Oct., 1992; Madras Univ. B.Sc. (Stat. Main) 1991*]

**Solution.** We have

$$r_{12}^2 = \left[ \frac{r_{12} - r_{13}r_{23}}{\sqrt{(1 - r_{13}^2)(1 - r_{23}^2)}} \right]^2 \leq 1$$

$$\therefore (r_{12} - r_{13}r_{23})^2 \leq (1 - r_{13}^2)(1 - r_{23}^2)$$

$$\Rightarrow r_{12}^2 + r_{13}^2r_{23}^2 - 2r_{12}r_{13}r_{23} \leq 1 - r_{13}^2 - r_{23}^2 + r_{13}^2r_{23}^2$$

$$\Rightarrow r_{12}^2 + r_{13}^2 + r_{23}^2 - 2r_{12}r_{13}r_{23} \leq 1 \quad \dots (*)$$

This condition holds for consistent values of  $r_{12}$ ,  $r_{13}$  and  $r_{23}$ .  $(*)$  may be rewritten as :

$$r_{23}^2 - (2r_{12}r_{13})r_{23} + (r_{12}^2 + r_{13}^2 - 1) \leq 0.$$

Hence, for given values of  $r_{12}$  and  $r_{13}$ ,  $r_{23}$  must lie between the roots of the quadratic (in  $r_{23}$ ) equation

$$r_{23}^2 - (2r_{12}r_{13})r_{23} + (r_{12}^2 + r_{13}^2 - 1) = 0,$$

which are given by :

$$r_{23} = r_{12}r_{13} \pm \sqrt{r_{12}^2r_{13}^2 - (r_{12}^2 + r_{13}^2 - 1)}$$

Hence

$$r_{12}r_{13} - \sqrt{1 - r_{12}^2 - r_{13}^2 + r_{12}^2r_{13}^2} \leq r_{23} \leq r_{12}r_{13} + \sqrt{1 - r_{12}^2 - r_{13}^2 + r_{12}^2r_{13}^2} \quad \dots (**)$$

In other words,  $r_{23}$  must lie in the range

$$r_{12}r_{13} \pm \sqrt{1 - r_{12}^2 - r_{13}^2 + r_{12}^2r_{13}^2}$$

In particular, if  $r_{12} = k$  and  $r_{13} = -k$ , we get from  $(**)$

$$-k^2 - \sqrt{1 - k^2 - k^2 + k^4} \leq r_{23} \leq -k^2 + \sqrt{1 - k^2 - k^2 + k^4}$$

$$\Rightarrow -k^2 - (1 - k^2) \leq r_{23} \leq -k^2 + (1 - k^2)$$

$$\therefore -1 \leq r_{23} \leq 1 - 2k^2$$

**EXERCISE 10(g)**

1. (a) Explain partial correlation and multiple correlation.

(b) Explain the concepts of multiple and partial correlation coefficients.

Show that the multiple correlation coefficient  $R_{1.23}$  is, in the usual notations given by :

$$R_{1.23}^2 = 1 - \frac{\omega}{\omega_{11}}$$

2 (a) In the usual notations, prove that

$$R_{1.23}^2 = \frac{r_{12}^2 + r_{13}^2 - 2r_{12}r_{23}r_{31}}{1 - r_{23}^2} \leq r_{12}^2$$

(b) If  $R_{1.23} = 1$ , prove that  $r_{2.13}$  is also equal to 1. If  $R_{1.23} = 0$ , does it necessarily mean that  $R_{2.13}$  is also zero?

3. (a) Obtain an expression for the variance of the residual  $X_{1.23}$  in terms of the correlations  $r_{12}$ ,  $r_{23}$  and  $r_{31}$  and deduce that  $R_{1(23)} \geq r_{12}$  and  $r_{13}$ .

(b) Show that the standard deviation of order  $p$  may be expressed in terms of standard deviation of order  $(p - 1)$  and a correlation coefficient of order  $(p - 1)$ . Hence deduce that :

$$(i) \sigma_1 \geq \sigma_{1.2} \geq \sigma_{1.23} \geq \dots \geq \sigma_{1.23\dots n}$$

$$(ii) 1 - R_{1.23\dots n}^2 = (1 - r_{12}^2)(1 - r_{13}^2)\dots(1 - r_{1n-23\dots(n-1)}^2)$$

[Delhi Univ. M.Sc. (Stat.) 1987]

4. (a) In a  $p$ -variate distribution all the total (zero order) correlation coefficients are equal to  $\rho_0 \neq 0$ . If  $\rho_k$  denotes the partial correlation coefficient of order  $k$ , find  $\rho_k$ . Hence deduce that :

$$(i) \rho_k - \rho_{k-1} = -\rho_k \rho_{k-1}$$

$$(ii) \rho_0 \geq -1/(p-1).$$

[Delhi Univ. M.Sc. (Stat.), 1989]

(b) Show that the multiple correlation coefficient  $R_{1.23\dots j}$  between  $X_1$  and  $(X_2, X_3, \dots, X_j)$ ,  $j = 2, 3, \dots, p$  satisfies the inequalities :

$$R_{1.2} \leq R_{1.23} \leq \dots \leq R_{1.23\dots p}$$

[Delhi Univ. M.Sc. (Maths.), 1989]

5. (a)  $X_0, X_1, \dots, X_n$  are  $(n+1)$  variates. Obtain a linear function of  $X_1, X_2, \dots, X_n$  which will have a maximum correlation with  $X_0$ . Show that the correlation  $R$  of  $X_0$  with the linear function is given by

$$R = \left(1 - \frac{\omega}{\omega_{00}^2}\right)^{\frac{1}{2}}$$

where  $\omega = \begin{vmatrix} 1 & r_{01} & r_{02} & \dots & r_{0n} \\ r_{10} & 1 & r_{12} & \dots & r_{1n} \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ r_{n0} & r_{n1} & r_{n2} & \dots & 1 \end{vmatrix}$

and  $\omega_{00}$  is the determinant obtained by deleting the first row and the first column of  $\omega$ .

(b) With the usual notations, prove that

$$\sigma^2_{1,234\dots n} = \frac{\omega}{\omega_{11}} \sigma_1^2 = \sigma_1^2 (1 - r_{12}^2)(1 - r_{13}^2)\dots(1 - r_{1n}^2)$$

(c) For a trivariate distribution, prove that

$$r_{12,3} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{(1 - r_{12}^2)(1 - r_{23}^2)}}$$

6. (a) The simple correlation coefficients between temperature ( $X_1$ ), corn yield ( $X_2$ ) and rainfall ( $X_3$ ) are,  $r_{12} = 0.59$ ,  $r_{13} = 0.46$  and  $r_{23} = 0.77$ .

Calculate the partial correlation coefficients  $r_{12,3}$ ,  $r_{23,1}$  and  $r_{31,2}$ . Also calculate  $R_{1,23}$ .

(b) If  $r_{12} = 0.80$ ,  $r_{13} = -0.40$  and  $r_{23} = -0.56$ , find the values of  $r_{12,3}$ ,  $r_{13,2}$  and  $r_{23,1}$ . Calculate further  $R_{1(23)}$ ,  $R_{2(13)}$  and  $R_{3(12)}$ .

7. (a) In certain investigation, the following values were obtained :

$$r_{12} = 0.6, r_{13} = -0.4 \text{ and } r_{23} = 0.7$$

Are the values consistent ?

(b) Comment on the consistency of

$$r_{12} = \frac{3}{5}, r_{23} = \frac{4}{5}, r_{31} = -\frac{1}{2}.$$

(c) Suppose a computer has found, for a given set of values of  $X_1$ ,  $X_2$  and  $X_3$ ,

$$r_{12} = 0.91, \quad r_{13} = 0.33 \text{ and } r_{23} = 0.81$$

Examine whether the computations may be said to be free from error.

8. (a) Show that if  $r_{12} = r_{13} = 0$ , then  $R_{1(23)} = 0$ . What is the significance of this result in regard to the multiple regression equation of  $X_1$  on  $X_2$  and  $X_3$  ?

(b) For what value of  $R_{1,23}$  will  $X_2$  and  $X_3$  be uncorrelated ?

(c) Given the data :  $r_{12} = 0.6$ ,  $r_{13} = 0.4$ , find the value of  $r_{23}$  so that  $R_{1,23}$ , the multiple correlation coefficient of  $X_1$  on  $X_2$  and  $X_3$  should be unity.

9. From the heights ( $X_1$ ), weights ( $X_2$ ) and ages ( $X_3$ ) of a group of students the following standard deviations and correlation coefficients were obtained :  $\sigma_1 = 2.8$  inches,  $\sigma_2 = 12$  lbs, and  $\sigma_3 = 1.5$  years,  $r_{12} = 0.75$ ,  $r_{23} = 0.54$ , and  $r_{31} = 0.43$ . Calculate (i) partial regression coefficients and (ii) partial correlation coefficients.

10. For a trivariate distribution :

$\bar{X}_1 = 40$	$\bar{X}_2 = 70$	$\bar{X}_3 = 90$
$\sigma_1 = 3$	$\sigma_2 = 6$	$\sigma_3 = 7$
$r_{12} = 0.4$	$r_{23} = 0.5$	$r_{13} = 0.6$

Find

(i)  $R_{1.23}$ , (ii)  $r_{23.1}$ , (iii) the value of  $X_3$  when  $X_1 = 30$  and  $X_2 = 45$ .

11. (a) In a study of a random sample of 120 students, the following results are obtained :

$$\begin{aligned}\bar{X}_1 &= 68, & \bar{X}_2 &= 70, & \bar{X}_3 &= 74 \\ S_1^2 &= 100, & S_2^2 &= 25, & S_3^2 &= 81, \\ r_{12} &= 0.60, & r_{13} &= 0.70, & r_{23} &= 0.65\end{aligned}$$

[ $S_i^2 = \text{Var}(X_i)$ ], where  $X_1$ ,  $X_2$ ,  $X_3$  denote percentage of marks obtained by a student in I test, II test and the final examination respectively.

(i) Obtain the least square regression equation of  $X_3$  on  $X_1$  and  $X_2$ .

(ii) Compute  $r_{12.3}$  and  $R_{3.12}$ .

(iii) Estimate the percentage marks of a student in the final examination if he gets 60% and 67% in I and II tests respectively.

(b)  $X_1$  is the consumption of milk per head,  $X_2$  the mean price of milk, and  $X_3$ , the per capita income. Time series of the three variables are rendered trend free and the standard deviations and correlation coefficients calculated :

$$s_1 = 7.22, \quad s_2 = 5.47, \quad s_3 = 6.87$$

$$r_{12} = -0.83, \quad r_{13} = 0.92, \quad r_{23} = -0.61$$

Calculate the regression equation of  $X_1$  on  $X_2$  and  $X_3$  and interpret the regression as a demand equation.

12. (a) Five thousand candidates were examined in the subjects (a), (b); (c); each of these subjects carrying 100 marks. The following constants relate to these data :

	<i>Subjects</i>		
	(a)	(b)	(c)
Mean	39.46	52.31	45.26
Standard deviation	6.2	9.4	8.7
$r_{bc} = 0.47$	$r_{ca} = 0.38$	$r_{ab} = 0.29$	

Assuming normally correlated population, find the number of candidates who will pass if minimum pass marks are an aggregate of 150 marks for the three subjects together.

(b) Establish the equation of plane of regression for variates  $X_1$ ,  $X_2$ ,  $X_3$  in the determinant form

$$\begin{vmatrix} X_1/\sigma_1 & X_2/\sigma_2 & X_3/\sigma_3 \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{vmatrix} = 0$$

[Delhi Univ. B.Sc. (Maths. Hons.), 1986]

13. (a) Prove the identity

$$b_{12.3} b_{23.1} b_{31.2} = r_{12.3} r_{23.1} r_{31.2}$$

[Gujarat Univ. B.Sc., 1992]

(b) Prove that

$$R_{1 \cdot 2 \cdot 3}^2 = b_{12 \cdot 3} r_{12} \frac{\sigma_2}{\sigma_1} + b_{13 \cdot 2} r_{13} \frac{\sigma_3}{\sigma_1}$$

[Sardar Patel Univ. B.Sc., 1991]

14. (a) If  $X_3 = aX_1 + bX_2$  for all sets of values of  $X_1$ ,  $X_2$ , and  $X_3$ , find the value of  $r_{23 \cdot 1}$ .

(b) If the relation  $aX_1 + bX_2 + cX_3 = 0$  holds for all sets of values  $X_1$ ,  $X_2$  and  $X_3$ , what must be the partial correlation coefficients?

15. (a) If  $r_{12} = r_{23} = r_{31} = \rho \neq 1$ , then

$$r_{12 \cdot 3} = r_{23 \cdot 1} = r_{31 \cdot 2} = \frac{\rho}{1 + \rho} \text{ and } R_{1(23)} = R_{2(13)} = R_{3(12)} = \frac{\rho \sqrt{2}}{\sqrt{(1 + \rho)^2}}$$

(b)  $Y_1$ ,  $Y_2$ ,  $Y_3$  are uncorrelated standard variates.  $X_1 = Y_2 + Y_3$ ,  $X_2 = Y_3 + Y_1$ , and  $X_3 = Y_1 + Y_2$ . Find the multiple correlation coefficient between  $X_3$  and  $(X_1 \text{ and } X_2)$ .

16.  $X$ ,  $Y$ ,  $Z$  are independent random variables with the same variance. If

$$X_1 = \frac{1}{\sqrt{2}}(X - Z), X_2 = \frac{1}{\sqrt{3}}(X + Y + Z), X_3 = \frac{1}{\sqrt{6}}(X + 2Y + Z),$$

show that  $X_1$ ,  $X_2$ ,  $X_3$  have equal variances. Calculate  $r_{12 \cdot 3}$  and  $R_{1(23)}$ .

17. (a) If  $X_1$ ,  $X_2$  and  $X_3$  are three variables measured from their respective means as origin and if  $e_1$  is the expected value of  $X_1$  for given values of  $X_2$  and  $X_3$  from the linear regression of  $X_1$  on  $X_2$  and  $X_3$ , prove that

$$\text{Cov}(X_1, e_1) = \text{Var}(e_1) = \text{Var}(X_1) - \text{Var}(X_1 - e_1)$$

(b) If  $r_{12} = k$  and  $r_{23} = -k$ , show that  $r_{13}$  will lie between  $-1$  and  $1 - 2k^2$ .

18. (a) For three variables  $X$ ,  $Y$  and  $Z$ , prove that

$$r_{XY} + r_{YZ} + r_{ZX} \geq -\frac{3}{2} \quad \dots (*)$$

**Hint.** Let us transform  $X$ ,  $Y$ ,  $Z$  to their standard variables  $U$ ,  $V$  and  $W$ , (say), respectively, where

$$U = \frac{X - E(X)}{\sigma_X}, V = \frac{Y - E(Y)}{\sigma_Y}, W = \frac{Z - E(Z)}{\sigma_Z}$$

so that

$$E(U) = E(V) = E(W) = 0 \quad \left. \begin{array}{l} \sigma_U^2 = \sigma_V^2 = \sigma_W^2 = 1 \Rightarrow E(U^2) = E(V^2) = E(W^2) = 1 \\ \text{and} \end{array} \right\} \dots (**)$$

$$r_{UV} = \frac{\text{Cov}(U, V)}{\sigma_U \sigma_V} = \frac{E(UV) - E(U)E(V)}{\sigma_U \sigma_V} = E(UV) \quad \left. \begin{array}{l} \\ \text{and} \end{array} \right\} \dots (***)$$

Since correlation coefficient is independent of change of origin and scale, proving (\*) is equivalent to proving  
 $r_{UV} + r_{VW} + r_{UW} \geq -3/2 \quad \dots (****)$

To establish (\*\*\*\*) let us consider the  $E(U + V + W)^2$ , which is always non-negative i.e.,  $E(U + V + W)^2 \geq 0$ , and use (\*\*) and (\*\*\*).

(b)  $X, Y, Z$  are three reduced (standard) variates and  $E(YZ) = E(ZX) = -1/2$ , find the limits between which the coefficient of correlation  $r(X, Y)$  is necessarily placed.

**Hint.** Consider  $E(X + Y + Z)^2 \geq 0 \Rightarrow r \geq -\frac{1}{2}$ .

(c) If  $r_{12}$ ,  $r_{23}$  and  $r_{31}$  are correlation coefficients of any three random variables  $X_1$ ,  $X_2$  and  $X_3$  taken in pairs  $(X_1, X_2)$ ,  $(X_2, X_3)$  and  $(X_3, X_1)$  respectively, show that

$$1 + 2r_{12}r_{23}r_{31} \geq r_{12}^2 + r_{13}^2 + r_{23}^2$$

19. (a) If the relation  $aX_1 + bX_2 + cX_3 = 0$ , holds for all sets of values of  $X_1$ ,  $X_2$  and  $X_3$ , where  $X_1$ ,  $X_2$  and  $X_3$  are three standardised variables, find the three total correlation coefficients  $r_{12}$ ,  $r_{23}$  and  $r_{13}$  in terms of  $a$ ,  $b$  and  $c$ . What are the values of partial correlation coefficients if  $a$ ,  $b$  and  $c$  are positive?

(b) Suppose  $X_1$ ,  $X_2$  and  $X_3$  satisfy the relation  $a_1X_1 + a_2X_2 + a_3X_3 = k$ .

(i) Determine the three total correlation coefficients in terms of standard deviations and the constants  $a_1$ ,  $a_2$  and  $a_3$ .

(ii) State what the partial correlation coefficients would be.

20. (a) Show that the multiple correlation between  $Y$  and  $X_1, X_2, \dots, X_p$  is the maximum correlation between  $Y$  and any linear function of  $X_1, X_2, \dots, X_p$ .

(b) Show that for  $p$  variates there are  ${}^pC_2$  correlation coefficients of order zero and  ${}^{p-2}C_s \cdot {}^pC_2$  of order  $s$ . Show further that there are  ${}^pC_2 \cdot 2^{p-2}$  correlation coefficients altogether and  ${}^pC_2 \cdot 2^{p-1}$  regression coefficients.

### ADDITIONAL EXERCISES ON CHAPTER X

1. Find the correlation coefficient between

(i)  $aX + b$  and  $Y$ , (ii)  $lx + mY$  and  $X + Y$ , when correlation coefficient between  $X$  and  $Y$  is  $\rho$ .

2. If  $X_1$  and  $X_2$  are independent normal variates and  $U$  and  $V$  are defined by

$$U = X_1 \cos \alpha + X_2 \sin \alpha, \quad V = X_2 \cos \alpha - X_1 \sin \alpha,$$

show that the correlation coefficient  $\rho$  between  $U$  and  $V$  is given by

$$\rho^2 = 1 - \frac{4\sigma_1^2\sigma_2^2}{4\sigma_1^2\sigma_2^2 + (\sigma_1^2 - \sigma_2^2)\sin^2 2\alpha},$$

where  $\sigma_1^2$  and  $\sigma_2^2$  are variances of  $X_1$  and  $X_2$  respectively.

3. The variables  $X$  and  $Y$  are normally correlated, and  $\xi$ ,  $\eta$  are defined by

$$\xi = X \cos \theta + Y \sin \theta, \quad \eta = Y \cos \theta - X \sin \theta$$

Obtain  $\theta$  so that the distributions of  $\xi$  and  $\eta$  are independent.

4. A set of  $n$  observations of simultaneous values of  $X$  and  $Y$  are made by an observer and the standard deviations and product moment coefficient about the mean are found to be  $\sigma_X$ ,  $\sigma_Y$  and  $\rho_{XY}$ . A second observer repeating the same observations made a constant error  $e$  in observing each  $X$  and a constant error  $E$  in observing each  $Y$ . The two sets of observations are combined into a single set and coefficient of correlation calculated from it. Show that its value is

$$\sqrt{(\rho_{XY} + \frac{1}{4}eE) + \sqrt{(\sigma_X^2 + \frac{1}{4}e^2)(\sigma_Y^2 + \frac{1}{4}E^2)}}$$

**Hint.** here we have two sets of observations :

**1st Set :**  $(x_i, y_i)$ ,  $i = 1, 2, \dots, n$ ; Mean =  $\bar{x}$ , s.d. =  $\sigma_x$ .

Product moment coefficient  $\rho_{xy} = r_{xy} \sigma_x \sigma_y$

**2nd Set :**  $(x_i + e, y_i + E)$ ,  $i = 1, 2, \dots, n$

$$\text{Mean } (\bar{x}') = \frac{1}{N} \sum (x_i + e) = \bar{x} + e$$

$$\text{Variance} = \sigma_x'^2 = \frac{1}{n} \sum [(x_i + e) - (\bar{x} + e)]^2 = \frac{1}{n} \sum (x_i - \bar{x})^2 = \sigma_x^2$$

$$\text{Mean } (\bar{y}') = \bar{y} + E, \sigma_y'^2 = \sigma_y^2.$$

Product moment coefficient :

$$\rho_{xy}' = \frac{1}{n} \sum [(x_i + e) - (\bar{x} + e)][(y_i + E) - (\bar{y} + E)] = \rho_{xy}$$

To obtain the correlation coefficient for the combined set of  $2n$  observations use Formula (10-5); Example 10-11(a) page 10-15.

5. Each of  $n$  independent trials can materialise in exactly one of the results

$A_1, A_2, \dots, A_k$ . If the probability of  $A_i$  is  $p_i$  in every trial  $\left( \sum_{i=1}^k p_i = 1 \right)$ , find the probability of obtaining the frequencies  $r_1, r_2, \dots, r_k$  for  $A_1, A_2, \dots, A_k$  respectively in these trials. Also find  $E(r_j)$ ,  $\text{Var}(r_j)$  and show that the correlation coefficient between  $r_i$  and  $r_j$  is independent of  $n$ .

6. In a sample of size  $n$  from a multinomial population  $n_1, n_2, \dots, n_k$  are of type  $1, 2, \dots, k$  with  $\sum p_i = 1$ , where  $p_i$  is the probability of type  $i$  ( $i = 1, 2, \dots, k$ ). Show that the expected value of  $n_2$  when  $n_1$  is given is  $(n - n_1)p_2(1 - p_1)$  and hence or otherwise show that the coefficient of correlation between  $n_i$  and  $n_j$  is

$$= \left[ \frac{p_i p_j}{(1 - p_i)(1 - p_j)} \right]^{\frac{1}{2}}$$

7. A ball is drawn at random from an urn containing 3 white balls numbered 0, 1, 2; 2 red balls numbered 0, 1 and 1 black ball numbered 0. If the colours white, red and black are again numbered 0, 1 and 2 respectively, show that the correlation coefficient between the variables :  $X$ , the colour number and  $Y$ , the number of the ball is  $-\frac{1}{2}$ .

8. If  $X_1$  and  $X_2$  are two independent normal variates with a common mean zero and variances  $\sigma_1^2$  and  $\sigma_2^2$  respectively, show that the variates defined by

$$U_1 = X_1 + X_2 \quad \text{and} \quad U_2 = -\frac{\sigma_2}{\sigma_1} X_1 + \frac{\sigma_1}{\sigma_2} X_2$$

are independent and that each is normally distributed with mean zero and common variance  $(\sigma_1^2 + \sigma_2^2)$ .

9. If  $X_1, X_2$  and  $X_3$  are uncorrelated variables with equal mean  $M$  and variances  $V_1^2, V_2^2$  and  $V_3^2$  respectively, prove that correlation coefficient  $\rho$

between  $Z_1 = \frac{X_1}{X_3}$  and  $Z_2 = \frac{X_2}{X_3}$  is given by

$$\rho = \frac{V_3^2}{\sqrt{[(V_1^2 + V_3^2)(V_2^2 + V_3^2)]}}$$

**Hint.** Neglecting the cubes and higher powers of  $\frac{x_i}{M}$ ,  $x_i$  being the deviation of  $X_i$  from  $M$  and letting the means and s.d.'s of  $Z_1$  and  $Z_2$  to be  $I_1, I_2$  and  $s_1, s_2$  respectively, we get

$$\begin{aligned} I_1 &= \frac{1}{N} \sum \frac{X_1}{X_3} = \frac{1}{N} \sum_i (x_{1i} + M)(x_{3i} + M)^{-1} \\ &= \frac{1}{N} \sum \left( 1 + \frac{x_{1i}}{M} \right) \left( 1 + \frac{x_{3i}}{M} \right)^{-1} \\ &= \frac{1}{N} \sum \left[ \left( 1 - \frac{x_{3i}}{M} + \frac{x_{3i}^2}{M^2} - \dots \right) + \frac{x_{1i}}{M} - \frac{x_{1i}x_{3i}}{M^2} + \dots \right] \\ &= 1 + \frac{V_3^2}{M^2} \end{aligned}$$

$$\text{Similarly } I_2 = 1 + \frac{V_2^2}{M^2}$$

$$\therefore I_1 = I_2$$

$$\text{Now } s_1^2 = \frac{1}{N} \sum \left( \frac{X_1}{X_3} \right)^2 - I_1^2$$

$$\text{or } s_1^2 + I_1^2 = 1 + \frac{3V_3^2}{M^2} + \frac{V_1^2}{M^2}, \text{ and so we have } s_1^2 = \frac{V_3^2}{M^2} + \frac{V_1^2}{M^2}.$$

$$\text{Similarly } s_2^2 = \frac{V_2^2}{M^2} + \frac{V_3^2}{M^2}$$

$$\text{Now } N\rho s_1 s_2 = \sum \left( \frac{X_1}{X_3} - I_1 \right) \left( \frac{X_2}{X_3} - I_2 \right) = \frac{V_3^2}{M^2} \quad (\text{On simplification})$$

$$\text{Hence } \rho = \frac{N\rho s_1 s_2}{s_1 s_2} = \frac{V_3^2}{\sqrt{(V_3^2 + V_1^2)} \sqrt{(V_3^2 + V_2^2)}}$$

10. (*Weldon's Dice Problem*).  $n$  white dice and  $m$  red dice are shaken together and thrown on a table. The sum of the dots on the upper faces are noted. The red dice are then picked up and thrown again among the white dice left on the table. The sum of the dice on the upper faces is again noted. What is the correlation coefficient between the first and the second sums?

Ans.  $n/(n+m)$

11. Random variables  $X$  and  $Y$  have zero means and non-zero variances  $\sigma_X^2$  and  $\sigma_Y^2$ . If  $Z = Y - X$ , then find  $\sigma_Z^2$  and the correlation coefficient  $\rho(X, Z)$  of  $X$  and  $Z$  in terms of  $\sigma_X, \sigma_Y$  and the correlation coefficient  $r(X, Y)$  of  $X$  and  $Y$ .

For certain data  $Y = 1.2X$  and  $X = 0.6Y$ , are the regression lines. Compute  $r(X, Y)$  and  $\sigma_x/\sigma_y$ . Also compute  $\rho(X, Z)$ , if  $Z = Y - X$ .

[*Calcutta Univ. B.Sc. (Maths. Honors.), 1984*]

12. An item (say, a pen) from a production line can be acceptable, repairable or useless. Suppose a production is stable and let  $p, q, r$  ( $p + q + r = 1$ ), denote the probabilities for three possible conditions of an item. If the items are put into lots of 100 :

- (i) Derive an expression for the probability function of  $(X, Y)$  where  $X$  and  $Y$  are the number of items in the lots that are respectively in the first two conditions.
- (ii) Derive the moment generating function of  $X$  and  $Y$ .
- (iii) Find the marginal distribution  $X$ .
- (iv) Find the conditional distribution of  $Y$  given  $X = 90$ .
- (v) Obtain the regression function of  $Y$  on  $X$ .

[*Delhi Univ. M.A. (Eco.), 1985*]

13. If the regression of  $X_1$  on  $X_2, \dots, X_p$  is given by :

$$E(X_1 | X_2, \dots, X_p) = \alpha + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_p X_p$$

and

$$\begin{vmatrix} \sigma_{22} & \sigma_{23} & \dots & \sigma_{2p} \\ \sigma_{32} & \sigma_{33} & \dots & \sigma_{3p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p2} & \sigma_{p3} & \dots & \sigma_{pp} \end{vmatrix} > 0, \quad \begin{cases} \sigma_{ii} = \text{variances} \\ \sigma_{ij} = \text{covariances} \end{cases}$$

then the constants  $\alpha, \beta_2, \dots, \beta_p$  are given by

$$\alpha = \mu_1 + \frac{R_{12}}{R_{11}} \cdot \frac{\sigma_1}{\sigma_2} \cdot \mu_2 + \frac{R_{13}}{R_{11}} \cdot \frac{\sigma_1}{\sigma_3} \cdot \mu_3 + \dots + \frac{R_{1p}}{R_{11}} \cdot \frac{\sigma_1}{\sigma_p} \cdot \mu_p$$

and

$$\beta_j = -\frac{R_{1j}}{R_{11}} \cdot \frac{\sigma_1}{\sigma_j}, \quad (j = 1, 2, \dots, p)$$

where  $R_{ij}$  is the cofactor of  $\rho_{ij}$  in the determinant ( $R$ ) of the correlation matrix

$$R = \begin{vmatrix} \rho_{11} & \rho_{12} & \dots & \rho_{1p} \\ \rho_{21} & \rho_{22} & \dots & \rho_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{p1} & \rho_{p2} & \dots & \rho_{pp} \end{vmatrix}$$

[*Delhi Univ. M.Sc. (Stat.), 1988*]

14. Let  $X_1$  and  $X_2$  be random variables with means 0 and variances 1 and correlation coefficient  $\rho$ . Show that :

$$E[\max(X_1^2, X_2^2)] \leq 1 + \sqrt{1 - \rho^2}$$

Using the above inequality, show that for random variables  $X_1$  and  $X_2$  with means  $\mu_1$  and  $\mu_2$ , variances  $\sigma_1^2$  and  $\sigma_2^2$  and correlation coefficient  $\rho$  and for any  $k > 0$ ,

$$P [|X_1 - \mu_1| \geq k\sigma_1 \text{ or } |X_2 - \mu_2| \geq k\sigma_2] \leq \frac{1}{k^2} [1 + \sqrt{1 - \rho^2}]$$

15. Let the maximum correlation between  $X_0$  and any linear function of  $X_1, X_2, \dots, X_n$  be  $R$  and if  $r_{01} = r_{02} = \dots = r_{0n} = r$

and all other correlation coefficients are equal to  $s$ , then show that :

$$R = r \left[ \frac{n}{1 + (n - 1)s} \right]^{1/2}$$

16. If  $f = f(x, y)$  is the p.d.f. of  $BVN(0, 0, 1, 1, \rho)$  distribution, verify that :

$$\frac{\partial f}{\partial \rho} = \frac{\partial^2 f}{\partial x \partial y}$$

Further, if two new random variables  $U$  and  $V$  are defined by the relation

$$U = P(Z \leq x) \text{ and } V = P(Z \leq y) \text{ where } Z \sim N(0, 1),$$

prove the marginal distributions of both  $U$  and  $V$  are uniform in the interval  $(-\frac{1}{2}, \frac{1}{2})$  and their common variance is  $\frac{1}{12}$ .

Hence prove that  $R = \text{Corr.}(U, V)$ , satisfies the relation :  $\rho = 2 \sin(\pi R/6)$ .

[Delhi Univ. B.A. (Stat. Hons. Spl. Course), 1988]

17. If  $(X, Y) \sim BVN(\mu_x, \mu_y, \sigma_x^2, \sigma_y^2, \rho)$ , then prove that  $a + bX + cY$ , ( $b \neq 0, c \neq 0$ ) is distributed as  $N(a + b\mu_x + c\mu_y, b^2\sigma_x^2 + c^2\sigma_y^2 + 2bc\rho\sigma_x\sigma_y)$ .

[Delhi Univ. M.Sc. (Stat.), 1989]

18. Let  $X_1, X_2, X_3$  be a random sample of size  $n = 3$  from  $N(0, 1)$  distribution.

(a) Show that  $Y_1 = X_1 + \delta X_3, Y_2 = X_2 + \delta X_3$  has a bivariate normal distribution.

(b) Find the value of  $\delta$  so that  $\rho(Y_1, Y_2) = \frac{1}{2}$ .

(c) What additional transformation involving  $Y_1$  and  $Y_2$  would produce a bivariate normal distribution with means  $\mu_1$  and  $\mu_2$ , variances  $\sigma_1^2$  and  $\sigma_2^2$ , and the same correlation coefficient  $\rho$  ?

Ans. (b) -1 or 1. (c)  $Z_1 = \sigma_1 Y_1 + \mu_1, Z_2 = \sigma_2 Y_2 + \mu_2$ .

19. If  $(X, Y) \sim BVN(0, 0, 1, 1, \rho)$ , prove that :

$$E[\max(X, Y)] = [(1 - \rho)/\pi]^{1/2} \text{ and } E[\min(X, Y)] = -[(1 - \rho)/\pi]^{1/2}$$

20. If  $(X, Y) \sim BVN(0, 0, \sigma_1^2, \sigma_2^2, \rho)$ , show that  $r$ th cumulant of  $XY$  is given by :

$$\kappa_r = \frac{1}{2}(r - 1)! \sigma_1^r \sigma_2^r [(\rho + 1)^r + (\rho - 1)^r].$$

$$\text{Deduce that } E(X^2 Y^2) = \sigma_1^2 \sigma_2^2 (1 + 2\rho^2).$$

21. Let  $f$  and  $g$  be the p.d.f.'s of  $X$  and  $Y$  with corresponding distribution functions  $F$  and  $G$ . Also let

$$h(x, y) = f(x) g(y) [1 + \alpha (2f(x) - 1)(2G(x) - 1)]; |\alpha| \leq 1,$$

Show that  $h(x, y)$  is a joint p.d.f. with marginal p.d.f.'s  $f$  and  $g$ . Further, let  $f$  and  $g$  be  $N(0, 1)$  p.d.f.'s. Show that  $Z = X + Y$ , is not normally distributed, except in the trivial case  $\alpha = 0$ .

**Hint.** Find  $M_Z(t) = E(e^{tZ})$  and use  $\text{Cov}(X, Y) = \alpha/\pi$ .

22. State p.d.f. of bivariate normal distribution. Let  $X$  and  $Y$  have joint p.d.f. of the form :

$$f(x, y) = ke^{-\frac{1}{2} [a_{11}(x - b_1)^2 + 2a_{12}(x - b_1)(y - b_2) + a_{22}(y - b_2)^2]}; \\ -\infty < (x, y) <$$

Find (i)  $k$ , (ii) the correlation coefficient between  $X$  and  $Y$ .

23. Write down, but do not derive, the moment generating function for a pair of random variables which have a bivariate normal distribution with both means equal to zero.

The independent random variables  $X, Y, Z$ , are each normally distributed with mean 0 and variance 1. If  $U = X + Y + Z$  and  $V = X - Y + 2Z$ , show that  $U$  and  $V$  have bivariate normal distribution. Find the correlation of  $U$  with  $V$  and the expectation of  $U$  when  $V$  is equal to 1.

24. Let  $X_1$  and  $X_2$  have a joint m.g.f.

$$M(t_1, t_2) = [a(e^{t_1} + t_2 + 1) + b(e^{t_1} + e^{t_2})]^2$$

in which  $a$  and  $b$  are positive constants such that  $2a + 2b = 1$ .

Find  $E(X_1)$ ,  $E(X_2)$ ,  $\text{Var}(X_1)$ ,  $\text{Var}(X_2)$ ,  $\text{Cov}(X_1, X_2)$ .

Ans. Means = 1, Variances =  $\frac{1}{2}$ , Covariance =  $2a - \frac{1}{2}$ .

25.  $X_1, X_2, X_3$  have joint distribution as a multinomial distribution with parameters  $N, p_1, p_2, p_3$ . If  $r_{ij}$  is the correlation coefficient between  $X_i$  and  $X_j$ , find the expression for  $r_{12}, r_{23}$  and  $r_{31}$  and hence deduce the expression for the partial correlation coefficient  $r_{123}$ .

26. (i) If all the inter-correlations between  $(p+1)$  variates  $X_0, X_1, X_2, \dots, X_p$  are equal to  $r$ , show that each of the partial correlation co-efficients of order  $p-1$  is equal to  $r/[1+(p-1)r]$  and that the multiple correlation of  $X_0$  on  $X_1, X_2, \dots, X_p$  is given by

$$1 - R_{0(12\dots p)}^2 = \frac{(1-r)(1-pr)}{1+(p-1)r}$$

$$(ii) \quad r_{12} = (r_{123} - r_{13}r_{23})/\sqrt{(1-r_{13}^2)^{1/2}(1-r_{23}^2)^{1/2}}$$

- 27. If  $R$  denotes the multiple correlation co-efficient of  $X_1$  on  $X_2, X_3, \dots, X_p$  in  $p$ -variate distribution, prove that

(i)  $R^2 \geq R_0^2$ , where  $R_0$  is the correlation of  $X_1$  with any arbitrary linear function of  $X_2, X_3, \dots, X_p$ .

(ii)  $R^2 \geq R_1^2$ , where  $R_1$  is the multiple correlation coefficient of  $X_1$  with  $X_2, X_3, \dots, X_k, k < p$

$$(iii) \quad 1 - R^2 = \prod_{j=2}^p (1 - r_{1j23\dots(j-1)}^2)$$