

$\text{cov}(x, y) = \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y})$, where \bar{x} and \bar{y} are means of the variables x and y respectively

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

$$\text{var}(x) = \frac{1}{n} \sum (x_i - \bar{x})^2 = \frac{1}{n} \sum x_i^2 - \bar{x}^2$$

and similarly $\text{var}(y) = \frac{1}{n} \sum y_i^2 - \bar{y}^2$

So, we can write

$$r_{xy} = \frac{\frac{1}{n} \sum x_i y_i - \bar{x} \bar{y}}{\sqrt{\frac{1}{n} \sum x_i^2 - \bar{x}^2} \sqrt{\frac{1}{n} \sum y_i^2 - \bar{y}^2}} = \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

7.4 Properties of Correlation coefficients :-

(a) Correlation coefficient of any two variables is a pure number, i.e., is independent of the units of measurement. This follows from the definition of $r (= r_{xy})$

(b) $r_{xy} = r_{yx}$. It also follows from the definition

(c) The numerical value of Correlation Coefficient is independent of the change of origins and scales of the variables.

Suppose we are given n pair of values (x_i, y_i) , $i=1, 2, \dots, n$, of the variables x and y . Let us

introduce $u = \frac{x-a}{c}$, $v = \frac{y-b}{d}$, where

a, b, c and d are arbitrary constants and $c \neq 0, d \neq 0$.

Then corresponding to each pair of values (x_i, y_i) of

x and y , we have a pair of values

(u_i, v_i) of u and v where $u_i = \frac{x_i - a}{c}$, $v_i = \frac{y_i - b}{d}$

$$\text{So, } x_i = a + cu_i, \quad y_i = b + d v_i$$

$$\therefore x_i - \bar{x} = c(u_i - \bar{u})$$

$$\text{Similarly, } y_i - \bar{y} = d(v_i - \bar{v})$$

$$\text{Var}(x) = c^2 \text{var}(u)$$

$$\text{var}(y) = d^2 \text{var}(v)$$

$$\text{Cov}(x, y) = cd \text{Cov}(u, v)$$

$$\text{Hence } r_{xy} = \frac{cd \text{Cov}(u, v)}{\sqrt{c^2 \text{var}(u)} \sqrt{d^2 \text{var}(v)}} = \frac{cd}{|c||d|} r_{uv}$$

$$\text{So, } r_{xy} = \begin{cases} r_{uv} & \text{when } c, d \text{ are of same sign} \\ -r_{uv} & \text{when } c, d \text{ are of opposite sign.} \end{cases}$$

$$(d) \quad -1 \leq r \leq 1$$

Suppose we are given a pair of values (x_i, y_i) , $i=1, 2, \dots, n$, of

variable x and y . Then $r = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{s_x s_y}$

$$= \frac{1}{n} \sum_i \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right) = \frac{1}{n} \sum_i k_i q_i$$

$$\text{where } k_i = \frac{x_i - \bar{x}}{s_x}, \quad q_i = \frac{y_i - \bar{y}}{s_y}$$

$$\text{So, } \sum k_i q_i = nr$$

$$\sum k_i^2 = \sum_i \left(\frac{x_i - \bar{x}}{s_x} \right)^2 = \frac{1}{s_x^2} \sum (x_i - \bar{x})^2 = \frac{n s_x^2}{s_x^2} = n$$

$$\text{and similarly } \sum q_i^2 = n$$

Now since $\sum_{i=1}^n (k_i + q_i)^2 \geq 0$ (since squares of real quantities are non-negative)

$$\text{or, } \sum k_i^2 + \sum q_i^2 + 2 \sum k_i q_i \geq 0$$

$$\text{or, } n + n + 2nr \geq 0$$

$$\text{or, } r \geq -1 \quad \dots (i)$$

Again $\sum_{i=1}^n (b_i a - a_i)^2 \geq 0$ or $\sum b_i^2 + \sum a_i^2 - 2 \sum b_i a_i \geq 0$

or, $n \times n - 2nr \geq 0$ or $r \leq 1$... (ii)

From (i) and (ii), we have $-1 \leq r \leq 1$

Remark: When $r = -1$ when, for each i , $(b_i + a_i)^2 = 0$

or, $a_i = -b_i$ or, $\frac{y_i - \bar{y}}{s_y} = - \frac{x_i - \bar{x}}{s_x}$

or, $y_i = \bar{y} - \frac{s_y}{s_x} (x_i - \bar{x})$. Hence all pair of values

of x, y lies on the straight line $y = \bar{y} - \frac{s_y}{s_x} (x - \bar{x})$, slope of the line being negative

Similarly when $r = 1$, for each i , $(b_i - a_i)^2 = 0$

or, $y_i = \bar{y} + \frac{s_y}{s_x} (x_i - \bar{x})$

So, all pair of values of x, y lies on the straight line $y = \bar{y} + \frac{s_y}{s_x} (x - \bar{x})$, the slope of the line is positive.

Example 7.1 $n = 10, \sum x = 140, \sum y = 150, \sum (x-10)^2 = 180$

$\sum (y-15)^2 = 215, \sum (x-10)(y-15) = 60$. Find the

correlation coefficient between x and y .

Solution: Here $\bar{x} = \frac{\sum x}{n} = \frac{140}{10} = 14, \bar{y} = \frac{\sum y}{n} = \frac{150}{10} = 15$

now $\sum (x-10)^2 = 180$

or, $\sum (x-14+4)^2 = 180$

or, $\sum (x-14)^2 + 8 \sum (x-14) + 16 \times 10 = 180$ [since $n = 10$]

or, $\sum (x-14)^2 + 160 = 180$ [$\because \sum (x-14) = \sum (x-\bar{x}) = 0$]

or, $\sum (x-14)^2 = 20$ i.e. $\sum (x-\bar{x})^2 = 20$

or, $\sum (x-14)^2 = 20$ i.e. $\sum (x-\bar{x})^2 = 20$

$$\sum (y-15)^2 = 215 \quad \text{or,} \quad \sum (y-\bar{y})^2 = 215$$

Again $\sum (x-10)(y-15) = 60$

$$\text{or,} \quad \sum (x-14+4)(y-15) = 60$$

$$\text{or,} \quad \sum (x-14)(y-15) + 4 \sum (y-15) = 60$$

$$\text{or,} \quad \sum (x-14)(y-15) = 60 \quad \left[\text{Since } \sum (y-15) = \sum (y-\bar{y}) = 0 \right]$$

$$\text{i.e.,} \quad \sum (x-\bar{x})(y-\bar{y}) = 60$$

$$\text{Hence } r = \frac{\sum (x-\bar{x})(y-\bar{y})}{\sqrt{\sum (x-\bar{x})^2 \sum (y-\bar{y})^2}} = \frac{60}{\sqrt{20 \times 215}} = \frac{60}{\sqrt{4300}} = 0.92$$

Example 7.2 Let x and y be uncorrelated variables

with standard deviations s_1 and s_2 . Show that the correlation coefficient between x and $x+y$ is $\frac{s_1}{\sqrt{s_1^2 + s_2^2}}$

Solution: $r_{xy} = 0$ so that $\text{cov}(x, y) = 0$

Let $u = x+y$ or $\bar{u} = \bar{x} + \bar{y}$ or $u - \bar{u} = (x - \bar{x}) + (y - \bar{y})$

$$\text{Now, } \text{cov}(x, x+y) = \text{cov}(x, u) = \frac{1}{n} \sum (x - \bar{x})(u - \bar{u})$$

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

$$= \frac{1}{n} \sum (x - \bar{x})^2 + \frac{1}{n} \sum (x - \bar{x})(y - \bar{y})$$

$$= \text{var}(x) + \text{cov}(x, y) = s_1^2$$

$$\text{Var}(u) = \text{var}(x+y) = \text{var}(x) + \text{var}(y) + 2\text{cov}(x, y) = s_1^2 + s_2^2$$

$$\text{So, } r_{xu} = \frac{\text{cov}(x, u)}{\sqrt{\text{var}(x) \text{var}(u)}}$$

$$= \frac{s_1^2}{s_1 \sqrt{s_1^2 + s_2^2}} = \frac{s_1}{\sqrt{s_1^2 + s_2^2}}$$