

Correlation (r)

(In a bivariate distribution, If the change in one variable affects a change in the other variable, the variables are said to be correlated.)

* (The relationship between two variables such that a change in one variable gives in a positive or negative change in the other variable is known as correlation.)

2m (If the two variables deviate in the same direction, correlation is said to be direct or positive.) That is if the increase (or) decrease in one variable gives a corresponding increase (or) decrease in the other variable.

If the two variables deviate in the opposite direction, correlation is said to be inverse (or) negative. That is if a increase (or) decrease in one variable gives a corresponding decrease (or) increase in the other variable.

Example for positive correlation is Height and weight of a group of persons.

Example for negative correlation is The Volume and pressure of a perfect gas.

Correlation is said to be perfect if the deviation in one variable is followed by a corresponding and proportional deviation in the other.) 15M

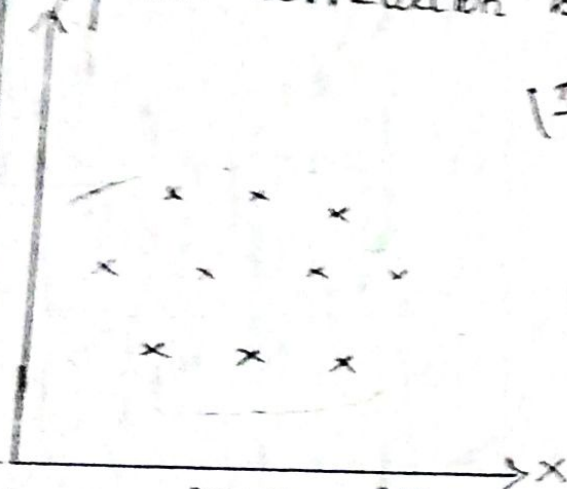
Scatter or Dot Diagram.

The values of the variables X and Y of a bivariate distribution can be plotted along the x axis and y axis. The diagram of dots so obtained is called a Scatter Diagram (or) dot diagram;

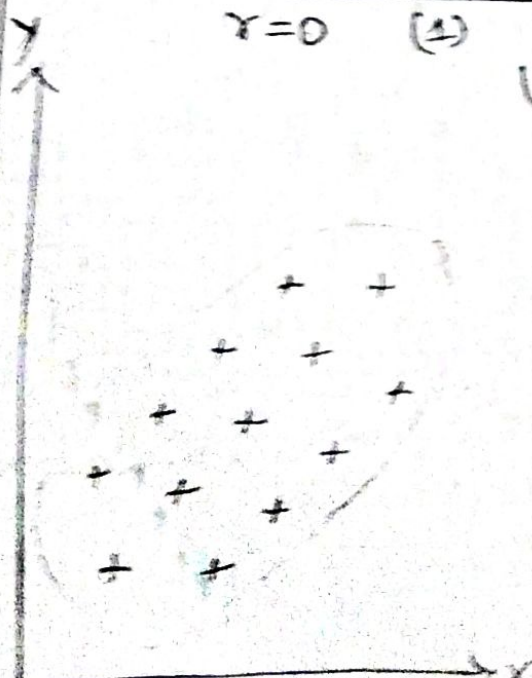
We can get a fairly good idea, whether the variables are correlated or not.

If the points are very close to other, we expect a fairly good amount of correlation between the variables.

If the points are widely scattered, a poor correlation is expected.

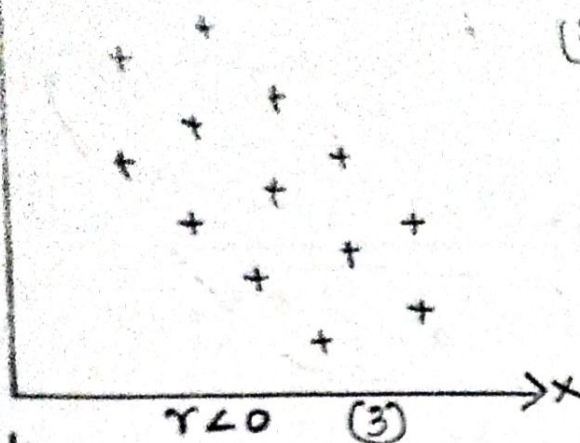


(1) This figure shows that there is no relationship between the variables.
i.e. $r=0$



(2) This figure shows that the variables are positively correlated. Here x increases, y also increases, when x is small y is also small.
i.e. $r > 0$

$r > 0$ (2)



(3) This figure shows that the variables are negatively correlated. When x is large y is small. When x is small y is large. i.e. $r < 0$.

Karl-Pearson's Coefficient of Correlation

Correlation coefficient between two random variables X and Y is denoted by $r(X, Y)$ or r_{xy} and defined as,

$$r_{xy} = \frac{\overset{\text{Covariance}}{\text{Cov}(X, Y)}}{\sigma_x \sigma_y}$$

$$\begin{aligned} \text{Cov}(X, Y) &= \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \sum [x_i y_i - \bar{x} y_i - x_i \bar{y} + \bar{x} \bar{y}] \\ &= \frac{1}{n} \sum x_i y_i - \bar{x} \left(\frac{1}{n} \sum y_i \right) - \left(\frac{1}{n} \sum x_i \right) \bar{y} + \bar{x} \bar{y} \\ &= \frac{1}{n} \sum x_i y_i - \bar{x} \bar{y} - \bar{x} \bar{y} + \bar{x} \bar{y} \end{aligned}$$

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

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

Similarly $\sigma_y^2 = \frac{1}{n} \sum y_i^2 - (\bar{y})^2$

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

Correlation coefficient cannot exceed unity. Correlation coeff limit is always lies between -1 and +1.

If $r = +1$ Correlation is perfect and positive
 If $r = -1$ Correlation is perfect and negative
 If the variables are independent $r = 0$; but the converse is not true.

Theorem:

Correlation Coefficient is independent of change of origin and scale.

Proof:

Let $U = \frac{X-a}{h}; V = \frac{Y-b}{k}$

a, b, k, h are constants. We have to

prove that $r_{xy} = r_{uv}$.

$$U = \frac{X-a}{h}$$

$$Uh = X - a$$

$$X = a + Uh$$

$$E(X) = a + h \cdot E(U) \quad X - E(X) = a + Uh - a - h E(U)$$

$$X - E(X) = h[U - E(U)] \quad X - E(X) = h[U - E(U)]$$

$$V = \frac{y-b}{k}$$

$$Vk = y - b$$

$$y = b + k \cdot v$$

$$E(y) = b + k \cdot E(v)$$

$$y - E(y) = k[v - E(v)]$$

$$\begin{aligned} \text{Cov}(x, y) &= E\left[\{x - E(x)\}\{y - E(y)\}\right] \\ &= E\left[h\{u - E(u)\} \cdot k\{v - E(v)\}\right] \\ &= E\left[hk\{u - E(u)\}\{v - E(v)\}\right] \\ &= hk \cdot E\left[\{u - E(u)\}\{v - E(v)\}\right] \\ &= hk \cdot \text{Cov}(u, v). \end{aligned}$$

$$\begin{aligned} \sigma_x^2 &= E\left[x - E(x)\right]^2 \\ &= E\left[h\{u - E(u)\}\right]^2 \\ &= E\left[h^2\{u - E(u)\}^2\right] \\ &= h^2 \cdot E\left[u - E(u)\right]^2 \end{aligned}$$

$$\sigma_x^2 = h^2 \cdot \sigma_u^2$$

$$\sigma_x = h \cdot \sigma_u$$

$$\begin{aligned} \sigma_y^2 &= E\left[y - E(y)\right]^2 \\ &= E\left[k\{v - E(v)\}\right]^2 \\ &= E\left[k^2\{v - E(v)\}^2\right] \\ &= k^2 E\left[v - E(v)\right]^2 \\ &= k^2 \cdot \sigma_v^2 \end{aligned}$$

$$\sigma_y = k \cdot \sigma_v$$

$$\begin{aligned} \sigma_y^2 &= E\left[y - E(y)\right]^2 \\ &= k \end{aligned}$$

Correlation Coefficient of X and Y is,

$$\begin{aligned} \gamma_{xy} &= \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} \\ &= \frac{k \cdot \text{Cov}(U, V)}{k \cdot \sigma_u \cdot k \cdot \sigma_v} \\ &= \frac{\text{Cov}(U, V)}{\sigma_u \sigma_v} \\ &= \gamma_{uv}. \end{aligned}$$

Hence $\gamma_{xy} = \gamma_{uv}$

Theorem.

Two independent Variables are uncorrelated

Proof
 $E = \text{expectation}$

$$\begin{aligned} \text{Cov}(X, Y) &= E \left[\{X - E(X)\} \{Y - E(Y)\} \right] \\ &= E \left[XY - X \cdot E(Y) - E(X) \cdot Y + E(X)E(Y) \right] \\ &= E(XY) - E(X) \cdot E(Y) - E(X) \cdot E(Y) + E(X) \cdot E(Y) \\ &= E(XY) - E(X) \cdot E(Y) \end{aligned}$$

If X and Y are independent $E(XY) = E(X)E(Y)$

$$\therefore \text{Cov}(X, Y) = E(X) \cdot E(Y) - E(X) \cdot E(Y) = 0.$$

$$\text{Hence } \gamma_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y} = 0.$$

i.e. The two independent Variables are uncorrelated.

Note Two independent Variables are uncorrelated, but the converse is not true, i.e. Two uncorrelated variables may not be independent.

5m

Probable Error of Correlation Coefficient.

If ' r ' is the correlation coefficient then the standard error $\left(S.E = \frac{1-r^2}{\sqrt{n}} \right)$

$$\begin{aligned} \text{Probable error (P.E)} &= 0.6745 \times S.E \\ &= 0.6745 \times \left\{ \frac{1-r^2}{\sqrt{n}} \right\} \end{aligned}$$

Probable error is helpful in testing the reliability of an observed correlation coefficient.

If $r < P.E$ correlation is not at all significant.

If $r > 6 P.E$ it is definitely significant.

Probable error also helps us to find the limits within which the correlation coefficient can be expected to vary.

The limits are $r \pm P.E.$

Limits for Correlation Coefficient = 10 + 1

5m

$$r_{xy} = \frac{Cov(x,y)}{\sigma_x \sigma_y}$$

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

Squaring both sides.

$$r_{xy}^2 = \frac{\frac{1}{n} \left[\sum (x_i - \bar{x})(y_i - \bar{y}) \right]^2}{\frac{1}{n} \sum (x_i - \bar{x})^2 \cdot \frac{1}{n} \sum (y_i - \bar{y})^2}$$

Let $a_i = (x_i - \bar{x})$ and $b_i = (y_i - \bar{y})$

$$r_{xy}^2 = \frac{(\sum a_i b_i)^2}{(\sum a_i^2)(\sum b_i^2)}$$

➤ Schwartz Inequality states that, if a_i, b_i ($i=1, 2, \dots, n$) are real quantities then,

$$(\sum a_i b_i)^2 \leq (\sum a_i^2)(\sum b_i^2)$$

↓
equality holds
iff

The sign of equality holding iff,

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

Using this we get,

$$r_{xy}^2 \leq \frac{(\sum a_i^2)(\sum b_i^2)}{(\sum a_i^2)(\sum b_i^2)}$$

≤ 1 .

$$|r_{xy}| \leq 1$$

$$\Rightarrow -1 \leq r_{xy} \leq 1$$

i.e. Correlation Coefficient cannot exceed unity

Hence, Correlation Coefficient lies in between -1 and $+1$.

Rank Correlation

(The ranks $1, 2, 3, \dots, n$ are given to the n Students for their marks in Mathematics and English.)

Let X denote the ranks for Mathematics and Y denote the ranks for English.

Let $(x_i, y_i) \ i=1, 2, \dots, n$ be their ranks in Mathematics and English. Pearsonian Coefficient of correlation between the ranks x_i 's and y_i 's is called the Rank Correlation.

Assume that no two students get the same rank in either Mathematics or in English.

X and Y takes the values $1, 2, \dots, n$.

$$\begin{array}{l} X: 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad \dots \quad n \\ Y: \quad 2 \quad 5 \quad 4 \quad 3 \quad n \quad \dots \quad 1 \end{array}$$

$$\begin{aligned} \therefore \bar{X} = \bar{Y} &= \frac{1}{n} (1+2+3+\dots+n) \\ &= \frac{1}{n} \left[\frac{n(n+1)}{2} \right] \end{aligned}$$

$$\bar{X} = \bar{Y} = \frac{n+1}{2}$$

$$\begin{aligned} \sigma_x^2 &= \frac{1}{n} \sum x_i^2 - (\bar{x})^2 \\ &= \frac{1}{n} [1^2 + 2^2 + \dots + n^2] - \left(\frac{n+1}{2} \right)^2 \\ &= \frac{1}{n} \left[\frac{n(n+1)(2n+1)}{6} \right] - \frac{(n+1)^2}{4} \\ &= \frac{(n+1)(2n+1)}{6} - \frac{(n+1)^2}{4} \Rightarrow \frac{2(2n+1) - 3(n+1)}{12} \\ &= \frac{(n+1)}{12} [2(2n+1) - 3(n+1)] \end{aligned}$$

$$\begin{aligned}\sigma_x^2 &= \frac{(n+1)}{12} [4n+2-3n-3] \\ &= \frac{(n+1)}{12} [n-1] \\ &= \frac{(n^2-1)}{12}\end{aligned}$$

$$x_i \neq y_i; \quad \therefore \sigma_x^2 = \sigma_y^2 = \frac{(n^2-1)}{12}$$

$z = x - y$
 $d_i = (x_i - y_i)$

Add and Subtract \bar{x}

$$d_i = (x_i - \bar{x} - y_i + \bar{x})$$

$$d_i = x_i - y_i + \bar{x} - \bar{x}$$

$$= (x_i - \bar{x}) - (y_i - \bar{y}) \text{ taking } \bar{x} = \bar{y}$$

$$[\bar{x} = \bar{y}]$$

\leq and square both side

$$\sum d_i^2 = \sum [(x_i - \bar{x}) - (y_i - \bar{y})]^2$$

$$= \sum [(x_i - \bar{x})^2 + (y_i - \bar{y})^2 - 2(x_i - \bar{x})(y_i - \bar{y})]$$

Dividing by 'n'

$$\frac{1}{n} \sum d_i^2 = \frac{1}{n} \sum [(x_i - \bar{x})^2 + (y_i - \bar{y})^2 - 2(x_i - \bar{x})(y_i - \bar{y})]$$

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

$$= \sigma_x^2 + \sigma_y^2 - 2 \cdot \text{Cov}(x, y)$$

$$r = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} \Rightarrow r \cdot \sigma_x \sigma_y = \text{Cov}(x, y)$$

$$\therefore \frac{1}{n} \sum d_i^2 = \sigma_x^2 + \sigma_y^2 - 2r \cdot \sigma_x \sigma_y$$

$$\begin{aligned}\sigma_x = \sigma_y \quad \therefore \frac{1}{n} \sum d_i^2 &= \sigma_x^2 + \sigma_x^2 - 2r \cdot \sigma_x \sigma_x \\ &= 2\sigma_x^2 - 2r \cdot \sigma_x^2\end{aligned}$$

$$\frac{1}{n} \sum d_i^2 = 2\sigma_x^2(1-r)$$

$$\frac{\sum d_i^2}{n \cdot 2\sigma_x^2} = 1-r$$

$$r = 1 - \left[\frac{\sum d_i^2}{2 \cdot n \cdot \sigma_x^2} \right]$$

$$= 1 - \left[\frac{\sum d_i^2}{2 \cdot n \cdot \left(\frac{n^2-1}{12} \right)} \right]$$

$$= 1 - \left[\frac{\sum d_i^2}{\frac{n(n^2-1)}{6}} \right]$$

$$= 1 - \left[\frac{6 \sum d_i^2}{n(n^2-1)} \right]$$

This is Spearman's formula for Rank Correlation Coefficient.

We use the symbol ' ρ ' for Rank Correlation Coefficient.

$$\therefore \rho = 1 - \left[\frac{6 \sum d_i^2}{n(n^2-1)} \right]$$

Repeated Ranks.

2nd [If any two or more students get the same rank, then the Spearman's formula for calculating the rank correlation coefficient breaks down.]

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 and the next item will get the rank next to the ranks already given.

As a result of this we add the correction factor $\frac{m(m^2-1)}{12}$ to $\sum d^2$, where 'm' is the number of items having the same rank. This correction factor is to be added for each repeated value.

$$r = 1 - \left\{ \frac{6(\sum d_i^2 + \text{correction factor})}{n(n^2-1)} \right\}$$

Limits for Rank Correlation Coefficient.

(Spearman's Rank Correlation coefficient is,

$$r = 1 - \left\{ \frac{6 \sum d_i^2}{n(n^2-1)} \right\}$$

[If $x_i = y_i$, then $d_i = 0$. $\therefore r = +1$
i.e. If the ranks are equal then the maximum value of r is $+1$]

If the ranks are in the opposite direction, then $r_s = -1$.

$$x: 1 \quad 2 \quad 3 \quad \dots \quad n$$

$$y: n \quad (n-1) \quad (n-2) \quad \dots \quad 1$$

$$\sum d^2 = (1-n)^2 + \{2-(n-1)\}^2 + \{3-(n-2)\}^2 + \dots + (n-1)^2 \quad d: x-y$$

$$= (1-n)^2 + (3-n)^2 + (5-n)^2 + \dots + (n-1)^2$$

$$= (n-1)^2 + (n-3)^2 + (n-5)^2 \dots \text{upto 'n' terms.}$$

$$\sum d^2 = \sum_{r=1}^n \left[n - (2r-1) \right]^2 \quad \text{sum of the natural } 1^{\text{st}} \text{ number}$$

$$= \sum_{r=1}^n \left[(n+1) - 2r \right]^2$$

$$= \sum_{r=1}^n \left[(n+1)^2 - 4r(n+1) + 4r^2 \right]$$

$$= n(n+1)^2 - 4 \sum_{r=1}^n r(n+1) + 4 \sum_{r=1}^n r^2$$

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

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

$$= n(n+1) \left[(n+1) - 2(n+1) + \frac{2}{3}(2n+1) \right]$$

$$= n(n+1) \left[n+1 - 2n - 2 + \frac{2}{3}(2n+1) \right]$$

$$= n(n+1) \left[-n - 1 + \frac{2}{3}(2n+1) \right]$$

$$= n(n+1) \left[\frac{-3n - 3 + 4n + 2}{3} \right]$$

$$\sum d^2 = n(n+1) \left[\frac{(n-1)}{3} \right] = \frac{n(n^2-1)}{3}$$

$$\begin{aligned}
 \therefore \rho &= 1 - \left\{ \frac{6 \sum d^2}{n(n^2-1)} \right\} \\
 &= 1 - \left\{ \frac{6 n(n^2-1)}{3 n(n^2-1)} \right\} \\
 &= 1 - (2/3) \\
 &= 1 - 2 \\
 &= -1.
 \end{aligned}$$

∴ The limits for rank correlation coefficient is $-1 \leq \rho \leq 1$.

Qom