

UNIT-IV CORRELATION

Concept of Correlation

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 a positive or negative change in the other variable is known as correlation.

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

Ex. For positive correlation

(i) heights and weights of a group of persons.

(ii) The price and supply of a commodity.

Ex. For negative correlation

- (i) Price and demand of a commodity
- (ii) 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.

Methods of correlation:-

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 the diagram of dots obtained is scatter diagram. We can get a fairly good idea whether the variables are correlated and 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.

▶ 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

$$r_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} \rightarrow (1)$$

For

deviation
method:

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

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

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

$$\bar{x}\bar{y} = \frac{1}{n}$$

$$\bar{x} = \frac{\sum x_i}{n}$$

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

$$n\bar{x} = \sum x_i$$

$$n\bar{y} = \sum y_i$$

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

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

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

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

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

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

$$= \bar{x}$$

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

$$= \frac{\sum x_i^2}{n} - \frac{2 \sum x_i \bar{x}}{n} + \frac{\sum \bar{x}^2}{n} = \frac{1}{n} \sum [x_i^2 - 2x_i \bar{x} + \bar{x}^2]$$

$$= \frac{\sum x_i^2}{n} - 2\bar{x} \frac{\sum x_i}{n} + \bar{x}^2$$

$$= \frac{\sum x_i^2}{n} - 2\bar{x} \bar{x} + \bar{x}^2$$

$$= \frac{\sum x_i^2}{n} - 2\bar{x} \bar{x} + \bar{x}^2$$

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

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

Similarly,

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

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

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

Correlation coefficient cannot exceed unity. It 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:

$$\text{Prove: } \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{\text{Cov}(u, v)}{\sigma_u \sigma_v}$$

$$\text{Let } u = \frac{x-a}{h} \quad ; \quad v = \frac{y-b}{k}$$

$$\therefore x = a + hu \quad ; \quad y = b + kv$$

$$\sigma_{uv} = r_{uv}$$

We have to prove that $r_{xy} = r_{uv}$

$$x = a + hu$$

$$y = b + kv$$

Taking Expectation

$$E(x) = a + hE(u)$$

$$E(y) = b + kE(v)$$

$$x - E(x) = (a + hu) - (a + hE(u)) \quad | \quad y - E(y) = (b + kv) - (b + kE(v))$$

$$= u \cdot h - h \cdot E(u) \quad | \quad = kv - k \cdot E(v)$$

$$= h(u - E(u)) \quad | \quad = k(v - E(v))$$

$$x - E(x) = h \{u - E(u)\}; \quad y - E(y) = k \{v - E(v)\}$$

$$\text{COV}(x, y) = E \{ [h \{u - E(u)\}] [k \{v - E(v)\}] \}$$

$$= E [h k \{u - E(u)\} \{v - E(v)\}]$$

$$= h k \cdot E [\{u - E(u)\} \cdot \{v - E(v)\}]$$

$$= h k \cdot \text{COV}(u, v)$$

$$\sigma_x^2 = E [(x - E(x))^2]$$

$$\sigma_y^2 = E [(y - E(y))^2]$$

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

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

$$\sigma_x = h \sigma_u$$

$$\sigma_y^2 = E [(y - E(y))^2]$$

$$= E [(k \{v - E(v)\})^2]$$

$$= E [k^2 \{v - E(v)\}^2]$$

$$= k^2 \cdot E [(v - E(v))^2]$$

$$\sigma_y^2 = k^2 \cdot \sigma_v^2$$

$$\sigma_y = k \sigma_v$$

correlation coefficient between x and y is

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

$$= \frac{h_x \cdot \text{cov}(u, v)}{h_x \cdot \sigma_u \cdot h_y \cdot \sigma_v}$$

$$r_{xy} = r_{uv}$$

correlation.
Regression

$$= \frac{\text{cov}(u, v)}{\sigma_u \cdot \sigma_v}$$

$$\therefore r_{xy} = r_{uv}$$

Theorem

Two independent variables are uncorrelated

$$\begin{aligned} \text{cov}(x, y) &= E[\{x - E(x)\}\{y - E(y)\}] \\ &= E[xy - x \cdot E(y) - E(x) \cdot y + E(x)E(y)] \\ &= E(xy) - E(x) \cdot E(y) - E(x)E(y) + E(x)E(y) \\ &= E(xy) - E(x) \cdot E(y) \end{aligned}$$

If x and y are independent,

$$E(xy) = E(x) \cdot E(y)$$

$$\therefore \text{cov}(x, y) = E(x)E(y) - E(x) \cdot E(y) = 0$$

Hence

$$r_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = 0$$

(i.e.) the two independent variables are uncorrelated, but the converse is not true.

(i.e.) Two uncorrelated variables may not be independent.

Example :- 1

a) If $Z = ax + by$ and r is the correlation coefficient between x and y . Show that

$$\sigma_z^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab r \cdot \sigma_x \sigma_y$$

b) Show that $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x \sigma_y}$ where σ_x , σ_y and σ_{x-y} are the standard derivations of x , y and $x-y$.

Soln:-

$$Z = ax + by$$

change into variable $\rightarrow [Z - E(Z)]$

$$E(Z) = aE(x) + bE(y)$$

$$Z - E(Z) = a\{x - E(x)\} + b\{y - E(y)\}$$

$$\begin{aligned} \sigma_z^2 &= \text{Var}(Z) = E\{[Z - E(Z)]^2\} \\ &= E\{[a\{x - E(x)\} + b\{y - E(y)\}]^2\} \\ &= E\{a^2\{x - E(x)\}^2 + b^2\{y - E(y)\}^2 + 2ab\{x - E(x)\}\{y - E(y)\}\} \end{aligned}$$

$$\sigma_x^2 = E\{[x - E(x)]^2\}$$

$$\begin{aligned} \sigma_y^2 &= E\{[y - E(y)]^2\} = a^2 E\{[x - E(x)]^2\} + b^2 E\{[y - E(y)]^2\} \\ &\quad + 2ab E\{[x - E(x)]\{y - E(y)\}\} \end{aligned}$$

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

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

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

$$\sigma_z^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab r \sigma_x \sigma_y$$

Putting $a=1$, $b=1$ we get $Z = x - y$

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

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

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

$$\begin{aligned} \text{Var}(ax + by) &= a^2 \text{Var}(x) + b^2 \text{Var}(y) \\ &\quad + 2ab \text{Cov}(x, y) \end{aligned}$$

$$\text{Var}(ax - by) = a^2 v(x) + b^2 v(y) - 2ab \text{cov}(x, y)$$

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

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

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

If $r > \text{P.E.}$ correlation is not at all significant.

If $r > 6 \text{ P.E.}$ 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 $r \pm \text{P.E.}$

Limits for correlation coefficient:-

$$r_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \frac{1/n \sum (x_i - \bar{x})(y_i - \bar{y})}{\left[\frac{1}{n} \sum (x_i - \bar{x})^2 \right]^{1/2} \left[\frac{1}{n} \sum (y_i - \bar{y})^2 \right]^{1/2}}$$

Squaring on both sides,

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

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

where $a_i = (x_i - \bar{x})$, $b_i = (y_i - \bar{y})$
 Schwarz's inequality states that,
 if $a_i, b_i (i=1, 2, \dots, n)$ are real
 quantities then,

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

the sign of equality holding

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

using this we get $r_{xy}^2 \leq 1$

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

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

(ie) correlation coefficient cannot
 exceed unity. Hence, correlation
 coefficient lies between -1 and +1.

Rank correlation

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, y_i is
 called the rank correlation.

Assume that no two students
 get the same rank in either
 classification.

x & y takes the values $1, 2, \dots, n$.

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

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

$$= \frac{n+1}{2}$$

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

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

$$= \frac{(n+1)}{12} [2(n+1) - 3(n+1)]$$

$$= \frac{n+1}{12} [4n+2 - 3n-3]$$

$$= \frac{n+1}{12} [n-1]$$

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

$$\therefore \bar{x} = \bar{y}, \sigma_x^2 = \sigma_y^2$$

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

$$x_i \neq y_i$$

Let

$$d_i = x_i - y_i$$

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

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

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

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

Dividing by n^2

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

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

$$= \sigma_x^2 + \sigma_y^2 - 2\rho\sigma_x\sigma_y$$

$$\frac{1}{n} \sum d_i^2 = 2\sigma_x^2 - 2\rho\sigma_x\sigma_y$$

$$= 2\sigma_x^2 - 2\rho \cdot \sigma_x^2$$

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

$$\sum d_i^2 = n \cdot 2\sigma_x^2(1-\rho)$$

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

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

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

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

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

This is known as Spearman's formula for the rank correlation co-efficient.

Note:-

For rank correlation let us use the symbol ρ for r .

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

Repeated Rank :-

If any two or more students get the same rank then the Spearman's formula for calculating the rank correlation co-efficient 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 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 we add the correlation factor

$$\frac{M(n^2-1)}{12}$$

to $\sum d_i^2$, where m is the number of times an item has repeated value. This correlation factor is to be added for each repeated value.

$$\therefore P = 1 - \left[\frac{6 \left\{ \sum d_i^2 + \text{correlation factor} \right\}}{n(n^2-1)} \right]$$

Limits for rank correlation co-efficient
Spearman's rank correlation coefficient is

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

If

$x_i = y_i$ then $d_i = 0 \therefore P = +1$

(ie) If the ranks are equal then the maximum value of P is ± 1

If the ranks are in the opposite direction the $P = -1$

(ie) $x: 1 \ 2 \ 3 \ \dots \ (n-1) \ n$

$y: n \ n-1 \ n-2 \ \dots \ 2 \ 1$

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

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

$$= (1-n)^n + (n-3)^2 + \dots \text{ up to } 6n^3 \text{ terms.}$$

$$\Sigma d^2 = \sum_{r=1}^n [n - (2r-1)]^2$$

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

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

$$= n(n+1)^2 - 4(n+1)\Sigma r + 4\Sigma r^2$$

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

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

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

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

$$\Sigma d^2 = \frac{n(n+1)}{3} [3n+3 - 6n - 6 + 4n+2]$$

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

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

$$P = 1 - \left\{ \frac{6 \Sigma d^2}{n(n^2-1)} \right\}$$

$$= 1 - \frac{6n(n^2-1)/3}{n(n^2-1)}$$

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

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

REGRESSION

definition:-

Regression analysis is a mathematical measure of the average relationship between two (or) more variables in terms of the original best fits of the data.

equation:-

Regression means stepping back towards

average:-

The line of regression is the line of "best fit" and is obtained by the principle of least squares.

Let the ^{regression} line of regression of Y on x is

$Y = a + bx_i$. x is independent variable Y is dependent variable.

According to the principle of least squares the constants a and b are to be determined such that

$$E = \sum_{i=1}^n (y_i - a - bx_i)^2 \text{ is minimum}$$

$$\frac{\partial E}{\partial a} = 2 \sum (y_i - a - bx_i)(-1)$$

$$\frac{\partial E}{\partial a} = 0 \Rightarrow -2 \sum (y_i - a - bx_i) = 0$$

$$\sum y_i - na - b \sum x_i = 0$$

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

$$\frac{\partial E}{\partial b} = 0$$

$$-2 \sum x_i (y_i - a - bx_i) = 0$$

$$\sum x_i y_i - a \sum x_i - b \sum x_i^2 = 0$$

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

These two are the normal equations for finding a and b . dividing (1) by n we get $\bar{y} = a + b\bar{x} \rightarrow (3)$

(ie) the line of regression of y on x passes through (\bar{x}, \bar{y})

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

$$\frac{1}{n} \sum x_i y_i = \text{cov}(x, y) + (\bar{x})(\bar{y}) \rightarrow (4)$$

Also,

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

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

$$\text{Dividing (2) by } n \cdot \frac{\sum x_i y_i}{n} = a \frac{\sum x_i}{n} + b \frac{\sum x_i^2}{n}$$

Using (4) & 4.a we get

$$\text{cov}(x, y) + \bar{x} \bar{y} = a \bar{x} + b [\sigma_x^2 + (\bar{x})^2]$$

Multiplying (3) by \bar{x}

$$\bar{x} \bar{y} = a \bar{x} + b (\bar{x})^2$$

Subtracting

$$\text{cov}(x, y) = b \cdot \sigma_x^2$$

$$b_{yx} = \frac{\text{cov}(x, y)}{\sigma_x^2}$$

b is the slope, (\bar{x}, \bar{y}) is one point.

\therefore The equation is $y - \bar{y} = b(x - \bar{x})$

$$y - \bar{y} = \frac{\text{cov}(x, y)}{\sigma_x^2} (x - \bar{x})$$

$$y - \bar{y} = \frac{r \cdot \sigma_x \sigma_y}{\sigma_x^2} (x - \bar{x}) \quad r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

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

Starting with the equation $x = a + by$ and proceeding similarly we get the line of regression of x on y as:

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

Always we have two lines of regression one is y on x and the other is x on y .

Line of regression of y on x is used to find y for any given x regression line of x on y is used to find x for a given y .

where $r = \pm 1$, both the lines reduce to

$$\frac{y - \bar{y}}{\sigma_y} = \pm \frac{x - \bar{x}}{\sigma_x}$$

(ie) the two lines of regression coincide and we have only one line

Regression coefficients

The slope of the line 'b' is known as regression co-efficient.

Regression co-efficient of Y on X is

$$b_{yx} = \frac{\text{Cov}(X, Y)}{\sigma_y^2}$$
$$= r \cdot \frac{\sigma_x}{\sigma_y}$$

Regression coefficient of X on Y is

$$b_{xy} = \frac{\text{Cov}(X, Y)}{\sigma_x^2}$$
$$= r \cdot \frac{\sigma_y}{\sigma_x}$$

Properties of regression coefficients

1. Correlation co-efficient is the geometric mean between the regression co-efficients.

$$b_{xy} = r \cdot \frac{\sigma_y}{\sigma_x} \text{ and } b_{yx} = r \cdot \frac{\sigma_x}{\sigma_y}$$

$$b_{xy} \times b_{yx} = \left\{ r \cdot \frac{\sigma_y}{\sigma_x} \right\} \left\{ r \cdot \frac{\sigma_x}{\sigma_y} \right\} = r^2$$

$$\therefore r = \pm \sqrt{b_{xy} \cdot b_{yx}}$$

If the regression coefficients are positive correlation coefficient is positive.
If the regression coefficients are negative

$|r|$ is negative.

2. If one of the regression co-efficient is greater than unity the other must be less than unity

(ie) If $b_{yx} > 1$ then $b_{xy} < 1$

Let $b_{yx} > 1$

$$\frac{1}{b_{yx}} < 1$$

$$r^2 \leq 1 \Rightarrow b_{yx} \times b_{xy} \leq 1$$

Dividing by (b_{yx})

$$b_{xy} \leq \frac{1}{b_{yx}} < 1$$

(ie) $b_{yx} > 1 \Rightarrow b_{xy} < 1$

3. Arithmetic mean of the regression coefficient is greater than the correlation coefficient we have to prove that

$$\frac{1}{2} (b_{yx} + b_{xy}) \geq r$$

$$\Rightarrow b_{yx} + b_{xy} \geq 2r$$

$$\Rightarrow b_{yx} + b_{xy} \geq 2 (\pm \sqrt{b_{xy} + b_{yx}})$$

$$b_{yx} + b_{xy} \pm 2\sqrt{b_{xy} + b_{yx}} \geq 0$$

$[\sqrt{b_{yx}} + \sqrt{b_{xy}}]^2 \geq 0$ which is always true since the square of a real quantity is always non negative

4. Regression coefficients are independent of the origin but not scale.

$$\text{Let } u = \frac{x-a}{h} \text{ and } v = \frac{y-b}{k}$$

$$x = a + uh \text{ and } y = b + vk$$

a, b, h, k are constant

$$E(x) = a + hE(u) \text{ and } E(y) = b + kE(v)$$

$$x - E(x) = h\{u - E(u)\} \text{ and } y - E(y) = k\{v - E(v)\}$$

$$\begin{aligned} \sigma_x^2 &= E[(x - E(x))]^2 \text{ and } \sigma_y^2 = E[(y - E(y))]^2 \\ &= E[h\{u - E(u)\}]^2 = E[k\{v - E(v)\}]^2 \\ &= h^2 E[u - E(u)]^2 = k^2 E[v - E(v)]^2 \\ &= h^2 \cdot \sigma_v^2 \end{aligned}$$

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

$$\begin{aligned} b_{yx} &= \frac{\text{cov}(x, y)}{\sigma_y^2} = \frac{hk \cdot \text{cov}(u, v)}{k^2 \sigma_v^2} \\ &= \frac{h}{k} \frac{\text{cov}(u, v)}{\sigma_v^2} = \frac{h}{k} b_{uv} \end{aligned}$$

Angles between two lines of regression

Regression line of y on x is

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

$$\text{line is } r \cdot \frac{\sigma_y}{\sigma_x}$$

Regression line of x on y is

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

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

$$\text{Slope is } \frac{\sigma_y}{r \cdot \sigma_x}$$

If θ is the acute angle between the two regression lines, then

$$\tan \theta = \frac{r \cdot \frac{\sigma_y}{\sigma_x} - \frac{\sigma_y}{r \cdot \sigma_x}}{1 + r \frac{\sigma_y}{\sigma_x} \cdot \frac{\sigma_y}{r \cdot \sigma_x}}$$

$$= \frac{r^2 \sigma_y \sigma_x - \sigma_x \sigma_y}{r \sigma_x^2}$$

$$1 + \frac{\sigma_y^2}{\sigma_x^2}$$

$$= \frac{\sigma_x \sigma_y (r^2 - 1)}{r \sigma_x^2}$$

$$\frac{\sigma_x^2 + \sigma_y^2}{\sigma_x^2}$$

$$= \frac{\sigma_x \sigma_y (1 - r^2)}{r}$$

$$\sigma_x^2 + \sigma_y^2$$

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

$$\therefore \theta = \tan^{-1} \left[\frac{1-r^2}{r} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2} \right]$$

If $r=0$, $\tan \theta = \infty \Rightarrow \theta = \pi/2$
 (i.e.) If the two variables are uncorrelated, the regression lines are perpendicular to each other.

If $r = \pm 1$, $\tan \theta = 0 \Rightarrow \theta = 0$ (or) $\theta = \pi$.

(i.e.) The two regression lines 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.

Define Correlation:-

(The term correlation refers to the relationship between the variables.)

$$r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

Types of correlation:-

- * positive and negative
- * simple (or) partial (or) multiple
- * linear (or) non-linear (or) no correlation

Positive correlation:-

When the value of the two variables changes in the same direction there is a positive correlation between

the two variables.

Negative correlation:-

When the value of the two variables changes in the opposite direction, there is negative correlation between the two variables.

Simple correlation:-

When only two variables are considered as under positive or negative correlation the correlation between them is called simple correlation.

Partial correlation:-

When more than two variables are considered, the correlation between two of them when all other variables are held constant (i.e.) when the linear effects of all other variables on them are removed is called partial correlation.

Multiple correlation:-

When more than two variables are considered, the correlation between one of them and its estimate based on the group consisting of the other variables is called multiple correlation.

Linear correlation:-

When all the points lie on a line or scattered around a line, there is a linear correlation between the two variables.

Non-linear correlation:-

When all the points lie exactly on a curve or scattered around the curve, there is a non-linear correlation between the two variables.

NO correlation:-

When the points are scattered neither around a line nor around a curve, there is no correlation between the two variables.

