

Correlation & Regression

The same point of time, if a variable is changed due to change in another variable then both are Co-related.

* TYPES :-

i) Positive Correlation :-

If both variables move in same direction or have +ve relation.

eg: Height and weight / Day temperature and Sale of cold drink.

ii) Negative Correlation :-

If both variable move in opposite direction or have -ve relation.

eg: Price & demand / No of Claims or Profit of Insurance company.

iii) Zero correlation :-

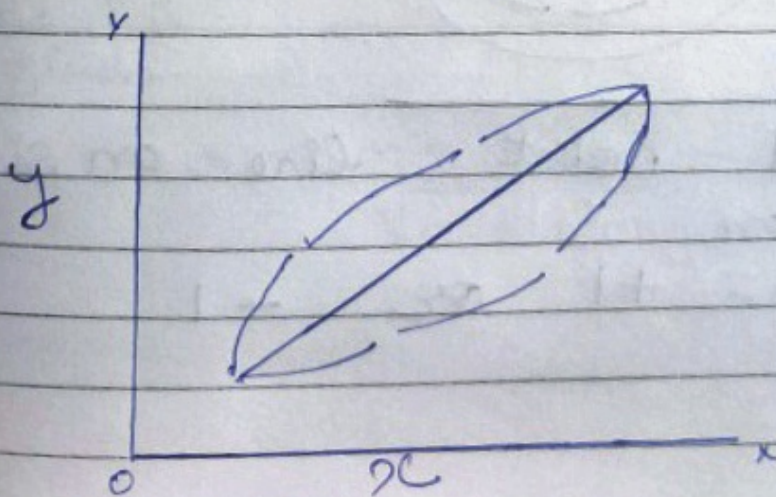
Variable does ^{not} move or produce

any difference in the number.
eg: Height & IQ / Rain and
space etc.

* Measures of Correlation

i) Scatter Diagram: This is the simple diagrammatic method that represent correlation coefficient between 2 variables. This is used when both variable are in linear or curvilinear relationship. It measures only nature of correlation but fails to find the value of correlation.

- In +ve correlation, plotted points lie from lower left corner to upper right corner.



- In -ve correlation, plotted points lie from upper left corner to lower right corner.



- In zero correlation, plotted points are equal and evenly distributed.



If plotted points lie on a single line then,
 $r = +1$ or -1

② Karl Pearson Method / Product Moment Method

This is used when variable are in linear relationship so this is shown by

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

$$\sigma_x = \text{SD of } x$$

$$\sigma_y = \text{SD of } y$$

$$\text{Cor}(x, y) = \text{Covariance}$$

$$(ii) r = \frac{\sum dx dy}{n \sigma_x \sigma_y}$$

$\sum dx dy \rightarrow$ Sum of product of

$$(iii) r = \frac{N \sum xy - \sum x \sum y}{\sqrt{N \sum x^2 - (\sum x)^2} \times \sqrt{N \sum y^2 - (\sum y)^2}}$$

When table of x & y is given

||
Huge mo.

$$(iv) r = \frac{m \sum dx dy - \sum dx \sum dy}{\sqrt{m \sum d^2 m - (\sum dx)^2} \times \sqrt{m \sum dy - (\sum dy)^2}}$$

③ Spearman's Rank Correlation Coefficient

Is used when qualitative data (attributes) are given. This method is also used to find the level of agreement between two judges. eg (Singing, dancing, beauty etc)

$$r_s = 1 - \frac{6 \sum D^2}{n(n^2-1)}$$

$\sum D^2$ = Sum of sq of Difference in Rank
 n = total no. of observation

least imp

$$r_s = 1 - \frac{6 \left[\sum d^2 + \frac{\sum t^3 - t}{12} \right]}{n(n^2-1)}$$

↓
In case of tie

t = length of tie

* Properties:

• $-1 \leq r_s \leq 1$

• $\sum D = 0$, sum of diff in Rank is 0

- IF ranks are given by two judges are in reverse order
 R_x 1 2 3 4 5
 R_y 5 4 3 2 1 $r_R = -1$

④ Coefficient of Concurrent Deviation

It is used when we are not concerned about the magnitude of variables

$$r_c = \pm \sqrt{\frac{2c - m}{m}}$$

c = no. of concurrent deviation

m = total no of observation

* Properties:

- If ~~$2c - m$~~ $(2c - m) > 0$ then we take the sign inside as well as outside
- If $(2c - m) < 0$ then we take -ve sign

$$-1 \leq r_c \leq 1$$

This is the quickest method

② Properties of Correlation

- r is unit free measurement
- r always lies b/w -1 to $+1$
(Both)
Inclusive

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

- r is not affected due to change in origin & scale but affected by sign.

$$r_{xy} = 0.5 \quad r_{uv} = ?$$

But like sign can be used

$$++ + + = +$$

$$+- ++ = -$$

$$-++- = +$$

$$+-+- = +$$

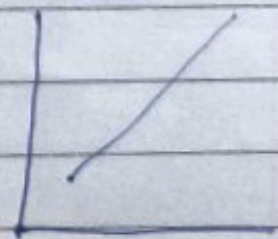
- Coefficient of determination is equal to r^2
- Coefficient of non-determination is equal to $1 - r^2$
- Percentage of unaccounted variation $(1 - r^2) \times 100$
- $b_x \times b_y \geq \text{Cor}(x, y)$

* Regression :-

In this concept we estimate or predict the value of one variable from the given value of another variable. It give exact value/Relation b/w x & y .

It tells what is the linear relation b/w the variable

eg $2x + 3y = 8$
 $y = 7 + 2x$



- There are two regression line:

↓
line of
y on x

$$y = a + bx \quad \left[\begin{array}{l} \text{Hindi} \\ \text{for eq.} \end{array} \right]$$

a = Intercept

bx = slope or
Regg. Coeff.

$$y - \bar{y} = b_{yx}(x - \bar{x})$$

for ques.
solving

b_{yx} = coeff of
regg. of y on
x

$$= r \cdot \frac{\sigma_y}{\sigma_x}$$

↓
line of
x on y

$$x = a + by$$

$$x - \bar{x} = b_{xy}(y - \bar{y})$$

b_{xy} = coeff of regg.
of x on y

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

* Properties of Regression Line

- Intersection point of both regression line is called mean (\bar{x}, \bar{y})

- Least square method is the best method to find regression line

$$\Sigma Y = Na + b \Sigma x$$

$$\Sigma xy = a \Sigma x + b \Sigma x^2$$

- If $r = (+1)$ or (-1) so both regression line are identical or co-incident

- If $r = 0$ then both line are perpendicular or bisect Right angle

- r is the geometric mean of both regression coefficient

$$r = \sqrt{b_{yx} \times b_{xy}}$$

- r , b_{yx} and b_{xy} always have same sign

- Product of b_{yx} & b_{xy} is always less than or equal to unity (1).
- Regression Coefficient is unaffected due to change in origin but affected by change in scale.

If $U = a_1 + b_1x$ & $V = a_2 + b_2y$

then $b_{UV} = \frac{p}{q} \times b_{xy}$ (Regg Coeff of UV)
or

$$\frac{\text{Regg Coeff of } UV}{\text{of } xy} = \frac{p}{q} \times b_{yx}$$

$$p = \frac{\text{coeff of } U}{\text{coeff of } x}$$

$$q = \frac{\text{coeff of } V}{\text{coeff of } y}$$