

Correlation & Regression

correlation

→ Degree of association between two variable

→ correlation can be negative positive or zero

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

→ for Direction Relation b/w $x \neq y$

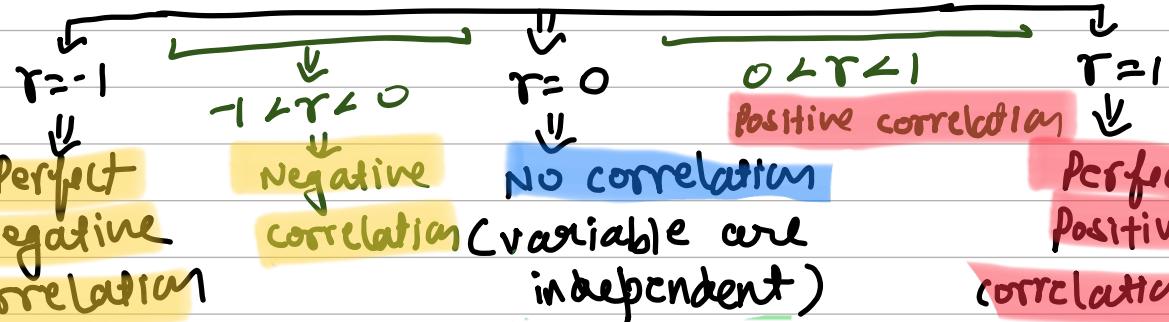
$$x \uparrow y \uparrow \text{ or } x \downarrow y \downarrow$$

r is positive

→ for inverse Relation b/w $x \neq y$

$$x \downarrow y \uparrow \text{ or } x \uparrow y \downarrow$$

r is negative



methods of calculating correlation

Graphical method

Scattered Diagram method

non graphical method

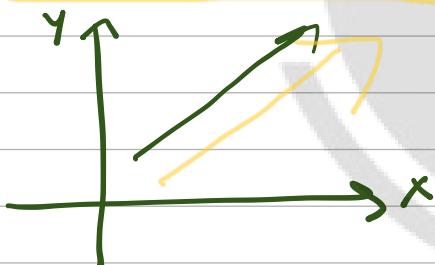
Karl Pearson method

Spearman method

Concurrent Deviations

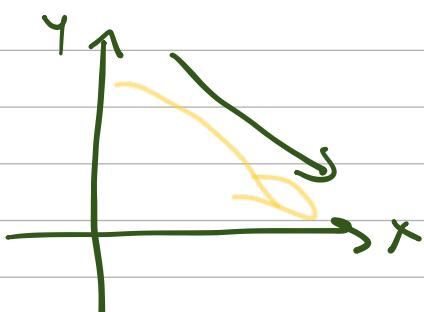
Scattered Diagram method

Plot the values of two variable on graph



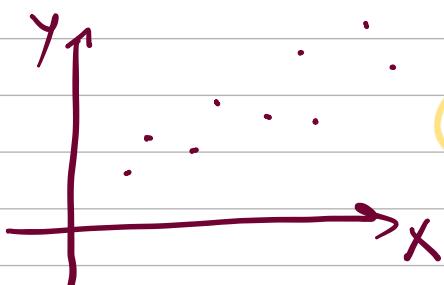
for straight line
(lower left to upper right)

$$r=1$$



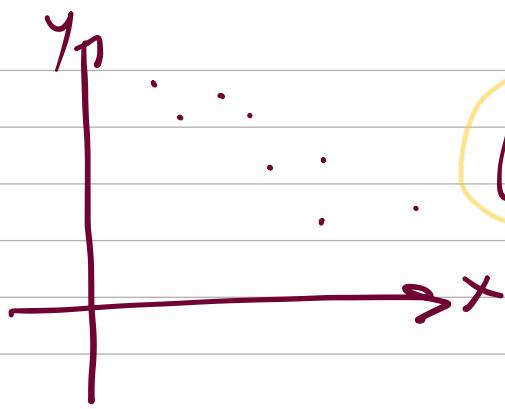
for straight line
(upper left to lower right)

$$r=-1$$



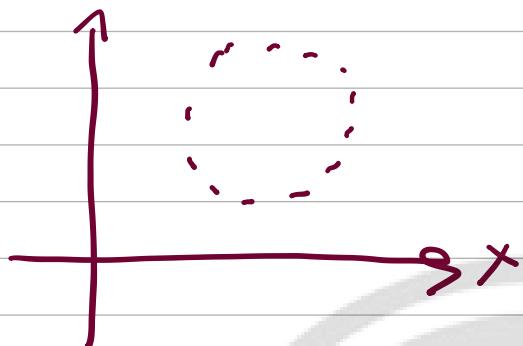
$$0 < r < 1$$

(no straight line)

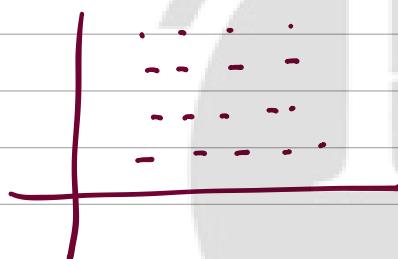


$$-1 < \gamma < 0$$

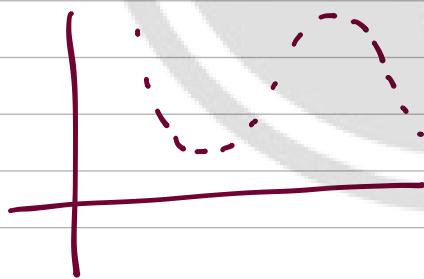
No straight Line



$$\gamma = 0$$



$$\overline{\gamma} = 0$$



$$\overline{\gamma} = 0$$

→ Exact magnitude can not be calculated in scattered diagram method

Karl Pearson's coefficient correlation

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

non-optional

→ covariance
 $= \text{cov}(x, y)$

optional

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{N s_x s_y}$$

→ $\text{cov}(x, y)$ can be
 only real number
 (negative, positive)
 zero

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \times \sqrt{\sum (y_i - \bar{y})^2}}$$

→ $\text{cov}(x, y)$ does not
 change with origin
 but it changes with
 change of scale

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

→ Correlation does not change with change
of origin

i.e. $U_i = x_i - A \quad \& \quad r_i = y_i - B$

then

$$r = \frac{\sum u_i v_i}{\sqrt{\sum u_i^2} \sqrt{\sum v_i^2}}$$
$$\sqrt{\sum u_i^2} = \sqrt{\frac{(x_i - A)^2}{N}} \quad \sqrt{\sum v_i^2} = \sqrt{\frac{(y_i - B)^2}{N}}$$

→ Correlation does not change with change
of scale (provided scale is positive)

→ When scale is negative, the sign of
correlation may change but magnitude
does not change

g $r(x,y) = 0.5$

if $v_i = 2x_i \quad \& \quad r_i = 4y_i$

then $r(v_i, r_i) = (+)(+) 0.5 = 0.5$

g $r(x,y) = 0.6$

if $v_i = -2x_i \quad \& \quad r_i = 3y_i$

then $r(v_i, r_i) = (-)(+) 0.6 = -0.6$

g $r(x,y) = 0.7$

if $v_i = -2x_i \quad \& \quad r_i = -4y_i$

then $r(v_i, r_i) = (-)(-) 0.7 = (+) 0.7$

Spearman's Rank Correlation

→ This method is used for qualitative characters & level of agreements & disagreements b/w opinions of judges

→ When no numbers repeat

$$\rho = 1 - \frac{6 \sum D^2}{N(N^2 - N)}$$

when some numbers repeat

$$\rho = 1 - \frac{6 \left[\sum D^2 + \frac{1}{12} (m_1^3 - m_1) + \frac{1}{12} (m_2^3 - m_2) \right]}{N^3 - N}$$

→ $\sum D = 0$ → $\sum d^2$ is always zero.

X	Y	R_1	R_2	$D = R_1 - R_2$	D^2
0	1	1	2	-1	1

#

Concurrent Deviation method

concurrent deviation



when n & y both increase
or both decrease

$$\gamma = \pm \sqrt{\pm \left(\frac{2c-m}{m} \right)}$$

where c = Total no of concurrent deviations

$$\underline{m = n - 1}$$

X	Y	concurrent deviations
+ 10	12	Yes (+)
+ 12	15	No (-)
+ 15	14	No (-)
- 14	15	Yes (+)
+ 16	18	Yes (+)
- 12	20	No (-)

$$c = 2$$

$$m = 6 - 1 = 5$$

$$\gamma = \sqrt{\frac{2(2) - 5}{5}}$$

→ when $\frac{2c-m}{m}$ is negative

$$\gamma = \sqrt{-\frac{1}{5}}$$

then (-) sign will be taken
out from $\sqrt{ }$ sign

$$\gamma = -\sqrt{\frac{1}{5}}$$

Bivariate Frequency Distribution Table

maths \ stats	0-5	5-10	10-15	15-20	Total
0-10	2	5	1	0	8
10-20	3	2	3	4	12
20-30	4	1	2	6	13
Total	9	8	6	10	33

Frequency distribution of marks of maths

marks	No of Students
0-10	8
10-20	12
20-30	13
	33

This is called marginal distribution

Frequency distribution of marks in maths when score in stats is '5-10'

marks	No of Students
0-10	5
10-20	2
20-30	1
	8

This is called conditional distribution

for Bivariate distribution table of "m×n"

Total cells = $m \times n$

Total marginal distribution = 2

Total conditional Distribution = $m+n$

coefficient of Determination

$$= \gamma^2 = \frac{\text{Explained variance}}{\text{Total variance}}$$

coefficient of Non Determination

$$= 1 - \gamma^2$$

Regressions

- Establishing mathematical relation b/w two variables (Independent & Dependent)
- Prediction of dependent variable
- for Linear Regression Least Square method is used

There two Linear Regression Lines

When y depends on x

when x depends on y

Standard form

$$y = a + bx$$

$$x = a + by$$

How to find = ?

use formula

$$y - \bar{y} = bxn (x - \bar{x})$$

use formula

$$x - \bar{x} = bny (y - \bar{y})$$

slope of line

$$\text{Slope} = bxn$$

$$\text{Slope} = \frac{1}{bny}$$

bxn & bny are known as regression coefficients

Calculation of Regression coefficients

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

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

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

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

$$b_{yx} = \frac{\sum xy - \frac{\sum x \sum y}{N}}{\sum x^2 - \frac{(\sum x)^2}{N}}$$

$$b_{xy} = \frac{\sum xy - \frac{\sum x \sum y}{N}}{\sum y^2 - \frac{(\sum y)^2}{N}}$$

$$b_{yx} = r \frac{\sigma_y}{\sigma_x}$$

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

#

$$\gamma = \pm \sqrt{b_{yx} \times b_{xy}}$$

$\rightarrow \gamma$ will be positive if both b_{yx} & b_{xy} are positive

$\rightarrow \gamma$ will be negative if both b_{yx} & b_{xy} are negative.

\rightarrow

$$(b_{yx} \times b_{xy}) \leq 1$$

#

$$\gamma \leq \frac{byx + bxy}{2}$$

If Regression Line y on x is: $ax + by + c = 0$

then

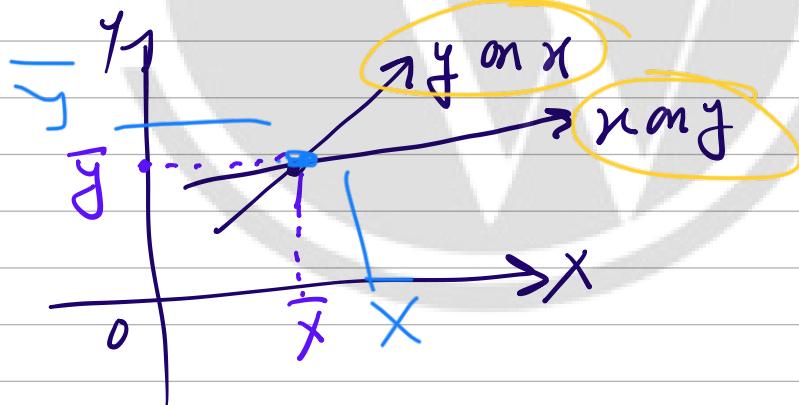
$$byx = -\frac{a}{b}$$

If Regression Line x on y is $Ax + By + C = 0$

$$bxy = -\frac{B}{A}$$

#

Two Regression Lines intersect each other at (\bar{x}, \bar{y})



#

Regression coefficients does not change with change of origin

$$u_i = x_i - A$$

$$v_i = y_i - B$$

$$\text{then } byx = bvu \quad \text{and} \quad bxy = buv$$

$$byx = bvu \quad \text{and} \quad bxy = buv$$

#

Regression coefficients change with
change of scale

$$U_i = a x_i + b \quad f$$

$$v_i = c y_i + d$$

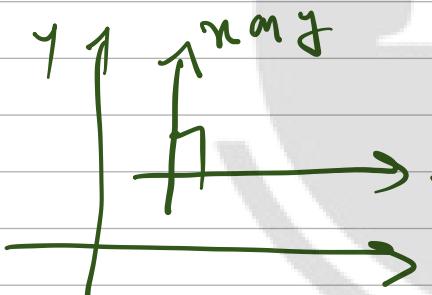
So

$$\delta v_u = \frac{c}{a} \times \delta y_n$$

or

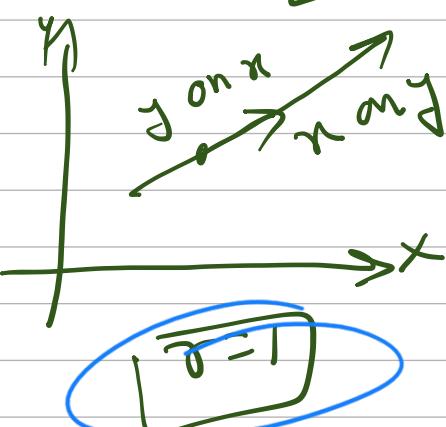
$$\delta v_u = \frac{\text{scale of } y}{\text{scale of } n} \times \delta y_n$$

when two Lines are perpendicular



$$\gamma = 0$$

when two lines are coincident



$$\gamma = 1$$

