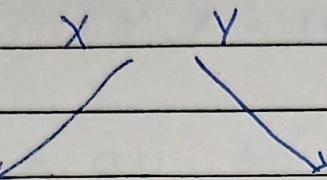


Unit - 6 : Correlation and Regression

#

Correlation -

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



Positive / Direct

Negative / Inverse

$$X \uparrow Y \uparrow$$

$$X \downarrow Y \downarrow$$

$$X \uparrow Y \downarrow$$

$$X \downarrow Y \uparrow$$

If change in one variable effects change in another variable than two variables are said to be correlation.

- $r = -1 \rightarrow$ perfect positive negative
- $r = 1 \rightarrow$ perfect positive
- $r = 0 \rightarrow$ no correlation

#

Results on correlation coefficient -

- correlation coefficient is independent of change of origin and change of scale.

$$r(x, y)$$

$$U = X - 10$$



origin

$$V = Y + 2$$

$$r(U, V) = r(X, Y)$$

$$U = 2X$$



scale

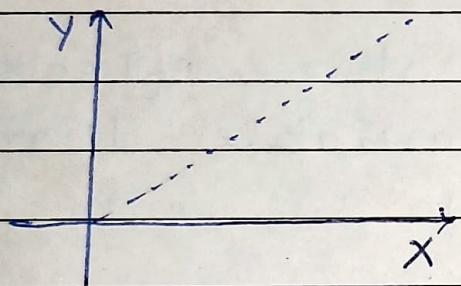
$$Y = \frac{V}{10}$$

$$u(U, V) = u(x, y)$$

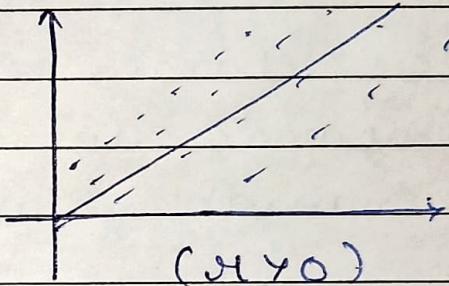
Independent

- Two variables are un-correlated.

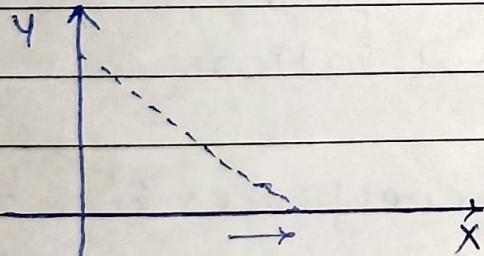
Scatter plot / diagrams -



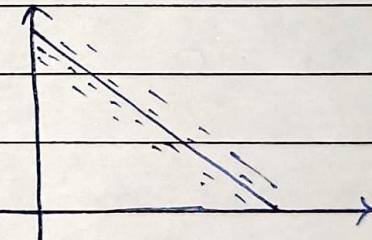
($\mu = 1$)



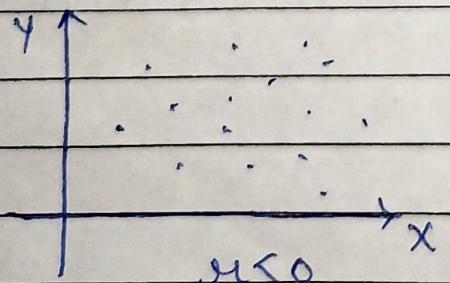
$(\mu > 0)$



$(\mu = -1)$



$(\mu < 0)$



#

Karl Pearson's correlation coefficient -

$$\rho(x, y) = \rho_{xy} = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{\sigma_{xy}}{\sigma_x \sigma_y}$$

where, $\sigma_{xy} = \frac{1}{n} \sum xy_i - \bar{x}\bar{y}$

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

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

- Q. Calculate the correlation coefficient for x and y.

UV	X	Y	$U = X - 68$	$V = Y - 69$	U^2	V^2
6	65	67	-3	-2	9	4
2	66	68	-2	-1	4	1
4	67	65	-1	-4	1	16
1	67	68	-1	-1	1	1
0	68	72	0	3	0	9
3	69	72	1	3	1	9
0	70	69	2	0	4	0
8	72	71	<u>4</u>	<u>2</u>	<u>16</u>	<u>4</u>
<u>24</u>			<u>0</u>	<u>0</u>	<u>36</u>	<u>44</u>

$$\bar{U} = 0 \quad \sigma_U^2 = \frac{1}{8} (36 - 0) = 4.5$$

$$\bar{V} = 0 \quad \sigma_V^2 = \frac{1}{8} (44 - 0) = 5.5$$

$$\Gamma_{UV} = \frac{1}{8}(24) - 0 = 3$$

$$\rho(U, V) = \frac{3}{\sqrt{4.5 \times 5.5}} = 0.604$$

Since correlation coefficient is independent of change of origin therefore -

$$\rho(X, Y) = \rho(U, V) = 0.604$$

#

correlation coefficient for Bi-variate data
and probability distribution -

Q.

The following table gives according to age,
the frequency of marks obtained by
100 students -

Age (x) →	18	19	20	21	Total	
Marks (y) ↓ u	-1	8	0	1	2	$h(v)$

-2	10 - 20	- 15	4	8	2	6	2	9	-	8
-1	20 - 30	- 25	5	5	4	6	6	6	4	19
0	30 - 40	- 35	6	0	8	0	10	6	11	0
1	40 - 50	- 45	4	4	4	0	6	6	8	16
2	50 - 60	- 55	-		2	0	4	8	4	16
3	60 - 70	- 65	-		2	0	3	9	1	6
	EUVf(u,v)		9		0		13		30	52
	Total G(u)		19		22		31		28	100
	ug(u)		-19		0		31		56	68
	$u^2g(u)$		19		0		31		56	162

$$u = x - 19$$

$$v = \frac{y - 35}{10}$$

$$\rho(u, v) = \frac{\Gamma_{uv}}{\Gamma_u \Gamma_v} = \frac{0.35}{\sqrt{1.1576 \times 1.6075}} = 0.26 \quad \text{Ans.}$$

$$\Gamma_{uv} = \frac{1}{N} \sum_u \sum_v uv f(u, v) - \bar{u}\bar{v}$$

$$= \frac{52}{100} - (0.68)(0.25) = 0.35$$

$$\sigma_u^2 = \frac{1}{N} \sum_u u^2 g(u) - \bar{u}^2 = \frac{162}{100} - (0.68)^2 = 1.157$$

$$\sigma_v^2 = \frac{1}{N} \sum_v v^2 g(v) - \bar{v}^2 = \frac{167}{100} - (0.25)^2 = 1.6075$$

$$\bar{u} = \frac{1}{N} \sum_u u g(u) = \frac{1}{100} (68) = 0.68$$

$$\bar{v} = \frac{1}{N} \sum_v v h(v) = \frac{25}{100} = 0.25$$

$h(v)$	$v h(v)$	$v^2 h(v)$	$\sum_{u,v} f(u,v)$
8	-16	32	4
19	-19	19	-9
35	0	0	0
22	22	22	18
10	20	40	24
<u>6</u>	<u>18</u>	<u>54</u>	<u>15</u>
<u>100</u>	<u>25</u>	<u>167</u>	<u>52</u>

#

Rank correlation -

① Spearman's Rank correlation coefficient -

$$f = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} \quad -1 \leq f \leq 1$$

$$d_i = x_i - y_i$$

$$\sum d_i = 0$$

Q. The ranks of same 16 students in maths, and physics are given

x	y	di	di^2
1	1	0	0
2	10	-8	64
3	3	0	0
4	4	0	0
5	5	0	0
6	7	-1	1
7	2	5	25
8	6	2	4
9	8	1	1
10	11	-1	1
11	15	-4	16
12	9	3	9
13	14	-1	1
14	12	2	4
15	6	-1	1
16	13	3	<u>9</u> <u>136</u>

Calculate the rank correlation coefficient for expertise of this group in mathematics and physics.

$$Sol. \rho = 1 - \frac{6 \sum di^2}{16(16^2 - 1)}$$

$$= 1 - \frac{6(136)}{16(16^2 - 1)}$$

$$\rho = 0.8 \quad \text{Ans.}$$

Q. 10 competitors in a musical test were ranked by three judges

Ranked by A(x) Rank by B (y) Rank by C (z)

1	3	6
6	5	5
5	8	9
10	4	8
3	7	1
2	10	2
4	2	3
9	1	10
7	6	5
8	9	7

Discuss which pair of judges has the nearest approach to common liking in music.

Solution of previous question -

$$d_1 = x - y$$

$$d_2 = y - z$$

$$d_3 = x - z$$

d_1	d_2	d_3	d_1^2	d_2^2	d_3^2
-2	-3	-5	4	9	
1	0	1	1	0	
-3	-1	-4	9	1	
6	-4	2	36	16	
-4	6	2	16	36	
-8	8	0	64	64	
2	-1	1	4	1	
8	-9	-1	64	81	
1	1	2	1	1	
-1	2	<u>1</u>	<u>1</u>	<u>4</u>	
			200	213	

$$\sum d_1^2 = 200$$

$$\sum d_2^2 = 214$$

$$\sum d_3^2 =$$

$$f(x,y) = 1 - \frac{6 \sum d_1^2}{n(n^2-1)} = 1 - \frac{6(200)}{10(99)} = -0.21$$

$$f(y,z) = 1 - \frac{6 \sum d_2^2}{n(n^2-1)} = 1 - \frac{6(213)}{990} = -0.29$$

$$f(x,z) = 1 - \frac{6 \sum d_3^2}{n(n^2-1)} = 1 - \frac{6(57)}{990} = 0.65$$

Since rank $g(x, z)$ is maximum, the judges A and C has the nearest approach to common liking in music.

Tied Ranks / Repeated Ranks -

Q. Obtained the rank correlation coefficient for the following data -

x	y	x	y	$d_i = x_i - y_i$	d_i^2
68	62	4	5	-1	1
64	58	6	7	-1	1
75	68	2.5	3.5	-1	1
50	45	9	10	-1	1
64	81	6	1	5	25
80	60	1	6	-5	25
75	68	2.5	3.5	-1	1
40	48	10	9	1	1
55	50	8	8	0	0
64	70	6	2	4	<u>16</u>
					<u>72</u>

$$\sum d_i^2 = 72$$

→ correction factor

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

$$n(n^2 - 1)$$

$$64 \rightarrow 3$$

$$68 \rightarrow 2$$

$$75 \rightarrow 2$$

$$f = 1 - \frac{6}{12} \left(72 + 3(3^2-1) + 2(2^2-1) + 2(2^2-1) \right) \\ 10(99)$$

$$f = 0.55$$

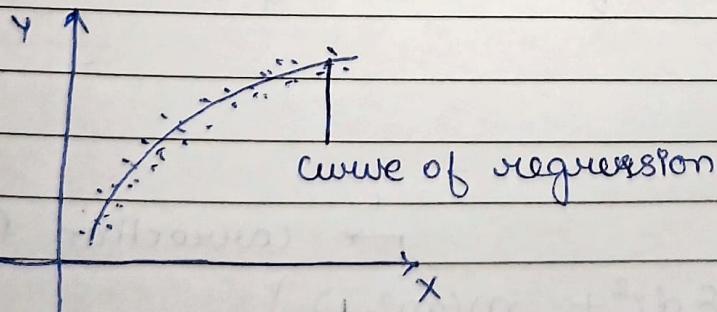
#

Regression

Regression is the study of nature of relationship between variables so that one may be able to predict the unknown value of one variable for a known value of another variable.

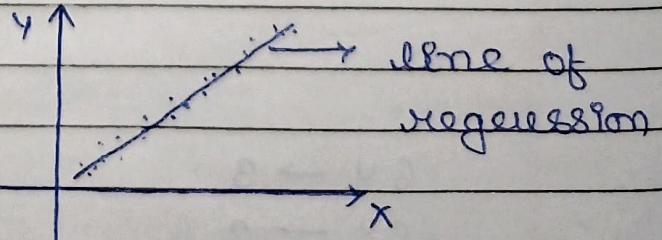
One variable is considered as independent and another is taken as dependent variable.

We predict the value of dependent variable using the value of independent variable.



#

Linear regression



Line of Regression of Y on X

$$Y = a + bX$$

The line of Y on X is given by -

$$Y - \bar{Y} = b_{yx}(X - \bar{x}), \text{ where } b_{yx} = \text{ or } \frac{\sigma_y}{\sigma_x} \text{ is regression coefficient of Y on X.}$$

Line of Regression of X on Y

$$X = a + bY$$

The line of X on Y is -

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

$$b_{xy} = \text{ or } \frac{\sigma_x}{\sigma_y} \text{ is regression coefficient of}$$

X on Y.

Q. Obtain the equations of two lines of regression for the following data

X	Y
65	67
66	68
67	65
67	68
68	72
69	72
70	69
72	71

Also obtain the estimate of x when $y = 70$.

Sol. From previous solved question in correlation.

$$\mu = 0.604$$

$$\sigma_x^2 = \sigma_u^2 = 4.5$$

$$\sigma_y^2 = \sigma_v^2 = 5.5$$

$$\bar{x} = 68 \quad \bar{y} = 69$$

$$\sigma_{x+b}^2 = \sigma_x^2$$

$$b_{yx} = \mu \frac{\sigma_y}{\sigma_x} = 0.604 \frac{\sqrt{5.5}}{\sqrt{4.5}} = 0.67$$

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

$$y - 69 = 0.67(x - 68)$$

$$y = 69 + 0.67x - 0.67 \times 68$$

$$y = 0.67x + 23.66$$

x on y

$$x = 0.55y + 30.05$$



when $y = 70$

$$x = 0.55 \times 70 + 30.05$$

$$x = 68.55$$

Properties of regression coefficient -

1) The correlation coefficient is the geometric mean between the regression coefficients denoted by γ

$$b_{xy} = \gamma \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

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

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

$$\gamma^2 = b_{xy} + b_{yx}$$

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

γ, b_{xy}, b_{yx}

2) γ, b_{xy}, b_{yx} all have the same signs.
(signs of regression coefficient should be same)

3) If one regression coefficient is greater than 1 than another must be less than 1.

4) Regression coefficients are independent of change of origin but not of scale.

$$U = \frac{x-a}{h}, \quad V = \frac{y-b}{k}$$

$$b_{xy} = n \frac{b_{uv}}{k}$$

$$b_{yx} = k \frac{b_{vu}}{h}$$

5) The point of intersection of the regression lines is ordered pair (\bar{x}, \bar{y}) .

6) The modulus of the arithmetic mean between two regression coefficients is greater than the modulus of the correlation coefficients.

$$\left| \frac{b_{xy} + b_{yx}}{2} \right| > |\gamma|$$

Angle between two regression line

$$\theta = \tan^{-1} \left[\frac{1-\gamma^2}{|\gamma|} \left(\frac{\rho_x \rho_y}{\rho_x^2 + \rho_y^2} \right) \right]$$

$$\text{If } \gamma=0, \quad \theta = \frac{\pi}{2}$$

If two variables are uncorrelated, the lines of regression are perpendicular to each other.

$$\text{If } \gamma = \pm 1, \quad \theta = 0 \text{ or } \pi$$

In the case of perfect correlation, the or -ve the two lines of regression coincide.

Q In a partially destroyed laboratory only the following result are readable. The variance of x is 9 and the regression lines are $8x - 10y + 66 = 0$
 $40x - 18y = 214$

- (i) what are the means of x and y - \bar{x} and \bar{y} ?
- (ii) What is the correlation coefficient between y and x
- (iii) What are standard deviation of y ($\sigma_y = ?$)?

Sol.

$$8\bar{x} - 10\bar{y} + 66 = 0$$

$$40\bar{x} - 18\bar{y} = 214$$

$$\bar{y} = 17 \quad \bar{x} = 13$$

x on y

$$x = \frac{10}{8}y - \frac{66}{8}$$

$$b_{xy} = \frac{10}{8}$$

y on x

$$y = \frac{40}{18}x - \frac{214}{18}$$

$$b_{yx} = \frac{40}{18}$$

$$\frac{10}{8} \times \frac{40}{18} = \frac{50}{18}$$

$$y = \frac{8x}{10} + \frac{66}{10}$$

$$b_{yx} = \frac{8}{10}$$

$$b_{xy} = \frac{18}{40}$$

$$\gamma^2 = \frac{8}{10} \times \frac{18}{40} = \frac{9}{25}$$

$$\boxed{\gamma = \pm \frac{3}{5}}$$

- $\gamma = \frac{3}{5} \rightarrow \text{positive}$

(iii)

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

$$\frac{8}{10} = \frac{3}{5} \times \frac{\sigma_y}{9}$$

$$\boxed{\sigma_y = 4}$$

Q.

Find the most likely price in mumbai corresponding to the price of rupees 70 at kolkata from the following

kolkata (X) mumbai (Y)

Average Price 65

67

Standard deviation 2.5

3.5

The correlation coefficient between the prices in two cities is : 0.8.

So,

$$\rho = 0.8$$

$$\bar{x} = 65, \bar{y} = 67$$

$$\sigma_x = 2.5, \sigma_y = 3.5$$

Y on X

$$Y - 67 = 0.8 \left(\frac{3.5}{2.5} \right) (X - 65)$$

$$Y = 1.12X - 5.8$$

$$Y = 1.12(70) - 5.8$$

$$\hat{Y} = 72.6$$

Fitting of a curve -

The def

Fitting of a straight line :

$$Y = a + bx$$

Normal Equation :

$$\sum Y = \sum a + \sum bx$$

$$\sum Y = na + b \sum X - \text{multiple by } \pm$$

$$\left[\begin{array}{l} \sum XY = \sum aX + \sum bx^2 \end{array} \right]$$

- Normal Equations

$$\left[\begin{array}{l} \sum XY = a \sum X + b \sum X^2 - \text{multiple by } X \end{array} \right]$$

Fitting a parabola -

$$Y = a + b_1 x + b_2 x^2$$

$$\begin{cases} EY = a + b_1 EX + b_2 EX^2 & \text{--- (i)} \\ EXY = a EX + b_1 EX^2 + b_2 EX^3 & \text{--- (ii)} \\ EX^2Y = a EX^2 + b_1 EX^3 + b_2 EX^4 & \text{--- (iii)} \\ \vdash \text{Normal Equation} \end{cases}$$

Q. For 10 randomly selected observations the following data were recorded were values of x and y -

x	y
1	2
1	7
2	7
2	10
3	8
3	12
4	10
5	14
6	11
7	14

Find the regression coefficient and the regression equation using the non linear form $y = a + b_1x + b_2x^2$

Sol.

$$\mathbb{E}Y = na + b_1\mathbb{E}x + b_2\mathbb{E}x^2$$

$$\mathbb{E}XY = a\mathbb{E}x + b_1\mathbb{E}x^2 + b_2\mathbb{E}x^3$$

$$\mathbb{E}X^2Y = a\mathbb{E}x^2 + b_1\mathbb{E}x^3 + b_2\mathbb{E}x^4$$

$$\mathbb{E}x = 34$$

$$\mathbb{E}XY = 377$$

$$\mathbb{E}x^2 = 154$$

$$\mathbb{E}X^2Y = 1849$$

$$\mathbb{E}Y = 95$$

$$\mathbb{E}x^3 = 820$$

$$\mathbb{E}x^4 = 4774$$

$$10a + 34b_1 + 154b_2 = 95$$

$$34a + 154b_1 + 820b_2 = 377$$

$$154a + 820b_1 + 4774b_2 = 1849$$

$$a = 1.8$$

$$b_1 = 3.49$$

$$b_2 = -0.27$$

$$Y = 1.8 + 3.49x - 0.27x^2$$

Q. Fit an exponential curve $y = ab^x$ to the following data

x	y	x^2	$v = \log_{10} y$
1	1		0
2	1.2		0.0792
3	1.8		0.2553
4	2.5		0.3979
5	3.6		0.5563
6	4.7		0.6721
7	6.6		0.8195
8	9.1		0.9590

Sol. $\log_{10} Y = \log_{10} a + x \log_{10} b$

$$v = A + BX$$

$$v = \log_{10} Y \quad B = \log_{10} b$$

$$A = \log_{10} a$$

$$a = 10^A \quad b = 10^B$$

$$U = A + BX$$

$$\Sigma U = nA + B \Sigma X$$

$$\Sigma UX = A \Sigma X + B \Sigma X^2$$

$$\Sigma X = 36$$

$$\Sigma X^2 = 204$$

$$U = \log_{10} Y$$

$$8A + 36B = 3.7393$$

$$36A + 204B = 22.7385$$

$$A = -0.16632 \Rightarrow a = 10^{-0.1662} = 0.682$$

$$B = 0.1408 \Rightarrow b = 10^{0.1408} = 1.383$$

$$Y = (0.682)(1.383)^X$$