

Unit - 2 [Predictive Analytics]

$\text{COPA} \Rightarrow \text{Collection, Organization, Presentation, Analysis, Interpretation.}$

Predictive Analytics: - Predictive modeling & Analysis. -
Regression analysis, correlation analysis, Rank correlation coefficient.

BIVARIATE DATA

The above example x_i may represent height and y_i weight of a group of persons such as data $(x_{ij}, y_j) \quad i = 1, 2, \dots, n$ is called BD.

Application of predictive modeling

- Predictive modeling is the process of using known results to create process, & validate a model that can be used to make future prediction.
- Two of the most widely used predictive modeling technique are regression & neural network.
- Companies can use predictive technique modeling to forecast event, customer behaviour as well as financial income economic & market risk.
- Other predictive modeling techniques used by

Correlation

Consider a set of bivariate of data
 $(x_{ij}, y_j), \quad i = 1, 2, \dots, n.$

The variables are said to be correlated, if there exist a change in one variable corresponding, to a change in the other variable.

Types of Correlation

i) Positive & negative Correlation

Correlation is positive [Direct] if the variables are direct / vary in the same direction . If they increase or decrease together correlation is negative [Inverse] if the variables are deviate / vary in the opposite direction . If one variable is increasing , the other variable is decreasing or vice versa.

ii) Simple, Partial & Multiple Correlation

Simple correlation are concerned with two variable only , while partial & multiple correlations are concerned with three or more related variables .

Eg.

Rice production (x), rain fall (y) & fertility (z) are partially connected .

iii) Linear & Nonlinear [curvilinear] correlation

If the ~~most~~ amount of change in one variable tends to near a constant ratio to the amount of change in the other variable , then the correlation is said to be linear otherwise nonlinear .

Eg.

If we double the amount of rainfall , the production of rice or wheat etc. would not necessarily be doubled .

Note : If the variables do not show related variation , they are said to non-correlated .

Eg.

The value of a grape (x) & atmospheric temperature (y) are non-correlated .

Measure of Correlation

There are various method of ascertaining whether two variables are correlated or not of these methods, the following three are

- Scatter Diagram [Dot diagram]
- Karl Pearson's [Product moments] coefficient of correlation
- Spearman's coefficient of rank correlation.

Scatter Diagram :-

It is a graphic representation of bivariate data. Here bivariate data with n pairs of values is represented by n points or dots on the xy -plane. The two variables are taken along the two axes and every pair of values in the data is represented by a point of graph.

Degrees of Correlation.

None	Low	High	Perfect

Types of correlation



Karl Pearson's coefficient of correlation

This is a measure of linear correlation

ship b/w the two variables. It indicates the degree of correlations b/w the two variables.
It is denoted by ' r '

$$r = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2] * [n \sum y^2 - (\sum y)^2]}}$$

Let $u = x - a$ here 'a' is assume value from 'x'
Let $v = y - b$ here 'b' is assume value from 'y'

$$r = \frac{n \sum uv - \sum u \sum v}{\sqrt{[n \sum u^2 - (\sum u)^2] * [n \sum v^2 - (\sum v)^2]}}$$

16/07/25
Thursday

Properties of coefficient of correlation

- ① The coefficient of correlation is independent of the units of measurement of the variables.
- ② The coefficient of correlation is independent of the origin and scale of measurement of the variable.
- ③ The coefficient of correlation is a value b/w -1 & +1.

Problem ① Calculate the product moments coefficient of correlation b/w the following marks (out of 10) in Computer science, mathematics of 5 marks.

Student	1	2	3	4	5
Computer	4	7	8	3	4
Mathematics	5	8	6	3	5

Sol:- Let x & y be the marks of computer science and mathematics
here $n = 5$

x	y	xy	x^2	y^2	$\sum x = 26$
4	5	20	16	25	$\sum y = 27$
7	8	56	49	64	$\sum xy = 153$
8	6	48	64	36	$\sum x^2 = 154$
3	3	9	9	9	$\sum y^2 = 159$
4	5	20	16	25	
$\sum x = 26$		$\sum y = 27$	$\sum xy = 153$	$\sum x^2 = 154$	$\sum y^2 = 159$

$$r^2 = \frac{n \sum xy - \sum x \sum y}{\sqrt{[n \sum x^2 - (\sum x)^2] * [n \sum y^2 - (\sum y)^2]}}$$

$$= \frac{5 \times 153 - 26 \times 27}{\sqrt{5 \times 154 - 26^2 * 5 \times 159 - 27^2}}$$

$$= \frac{765 - 702}{\sqrt{770 - 676 * 295 - 729}}$$

$$= \frac{63}{\sqrt{94 * 66}}$$

$$= \frac{63}{\sqrt{6204}}$$

$$= \frac{63}{78.765} = 0.7998$$

=

(2)

Husband	164	176	178	184	175	167	173	190
Wife	168	174	178	181	173	166	173	179

S.F.

$$\text{Let } x \& y \rightarrow u \& v$$

$$\text{Let } u = x - a =$$

$$= x - 170 \quad (164 - 184)$$

$$\text{Let } v = y - b \quad (161 - 171)$$

$$= y - 170$$

x	y	u	v	uv	u ²	v ²
164	168	-6	-2	12	36	4
176	174	6	4	24	36	16
178	175	8	5	40	64	25
184	181	14	11	154	196	121
175	173	5	3	15	25	9
167	166	-3	-4	12	9	16
173	173	3	3	9	9	9
180	179	10	9	90	100	81
		$\Sigma u = 37$	$\Sigma v = 29$	$\Sigma uv = 356$	$\Sigma u^2 = 475$	$\Sigma v^2 = 281$

$$r^2 = \frac{n \sum uv - \sum u \sum v}{\sum u^2 - (\sum u)^2} \times \frac{n \sum v^2 - (\sum v)^2}{\sum v^2 - (\sum v)^2}$$

$$\sum u = 37$$

$$\sum v = 29$$

$$\sum uv = 356$$

$$\sum u^2 = 475$$

$$= \frac{8 \times 356 - 37 \times 29}{8 \times 475 - 37^2}$$

$$\sum v^2 = 281$$

$$\sqrt{\frac{8 \times 356 - 37 \times 29}{8 \times 475 - 37^2} \times \frac{8 \times 475 - 37^2}{8 \times 281}}$$

$$n^2 = 8$$

$$= \sqrt{2,848 - 473 \times 1073}$$

$$= \sqrt{2,848 - 1369 \times 2,248 - 841}$$

$$= \sqrt{2,848 - 3,800}$$

$$= \sqrt{3,800 - 2,248}$$

= ~~1786~~ 1775

$$\sqrt{2431 * 1407} = 2431 * 1407$$

~~2431 * 1407~~
~~859 * 2959~~

PAGE NO.:

DATE:

$$= \frac{1786}{\sqrt{3420417}} = \frac{1786}{\sqrt{18494369}} = \frac{1786}{13597}$$

~~3420417~~
~~18494369~~

(3)

$$\begin{array}{ccccccc} x & 218 & 220 & 236 & 225 & 220 & 227 & 228 \\ y & 12.3 & 12.7 & 12.0 & 12.2 & 12.7 & 12.1 & 12.0 \end{array}$$

Sd Let $u = x - a$

$$\Rightarrow x = u + a$$

$$\bar{y} = y - b$$

$$= y - 12$$

10.7622

x	y	u	v	uv	u ²	v ²
218	12.3	-2	0.3	-0.6	4	0.09
220	12.7	0	0.7	0	0	0.49
236	12.0	16	0	0	256	0
225	12.2	5	0.2	1	25	0.04
220	12.7	0	0.7	0	0	0.49
227	12.1	7	0.1	0.7	49	0.01
228	12.0	8	0	0	64	0
$\Sigma u = 34$		$\Sigma v = 2$		$\Sigma uv = 1.1$	$\Sigma u^2 = 398$	$\Sigma v^2 = 1.12$

$$r = \frac{n \sum uv - \sum u \sum v}{\sqrt{[n \sum u^2 - (\sum u)^2] * [n \sum v^2 - (\sum v)^2]}}$$

$$7x_{1.1} - 34x_2$$

$$\sqrt{7x_{398} - 34^2 * 7x_{1.12} - 2^2}$$

$$7.7 - 68$$

$$\sqrt{2786 - 1156 * 7.84 - 4}$$

$$-60.3$$

$$\sqrt{1630 * 3.84} = 6.259.2$$

Spearman's coefficient of rank of correlation

Suppose that a group of n individual are arranged in the order of merit or efficiency with respect to some characteristics.

In a bivariate data, if the values of the variables are ranked in the decreasing (or increasing) order, the correlation b/w these ranks is rank correlation. The coefficient of correlation computed for these ranks is Spearman's coefficient of rank correlation. Denoted by ρ . (Read as Rho).

If R_1 and R_2 are the ranks in the two characteristics and $d = R_1 - R_2$ is the diff b/w the ranks, coefficient of correlation is

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

Problem 0

Marks in Statistics	25	43	27	35	54	61	37	45
Marks in Maths	35	47	20	37	63	54	28	40

x	y	R_1	R_2	$d = R_1 - R_2$	d^2
25	35	8	6	$8 - 6 = 2$	4
43	47	4	3	$4 - 3 = 1$	1
27	20	7	8	$7 - 8 = -1$	1
35	37	6	5	$6 - 5 = 1$	1
54	63	2	1	$2 - 1 = 1$	1
61	54	1	2	$1 - 2 = -1$	1
37	28	5	7	$5 - 7 = -2$	4
45	40	3	4	$3 - 4 = -1$	1
				$\sum d^2 = 14$	

$$P = 1 - \frac{6 \cdot \sum d^2}{n^3 - n}$$

$$= 1 - 6(14) = 1 - 84 = 1 - 0.1666$$

$$8^3 - 8 = 512 - 8$$

$$= 0.8333$$

~~H/W~~ Karl Pearson's problem

$$\textcircled{1} \quad \begin{array}{ccccccccc} x & 65 & 66 & 67 & 67 & 68 & 69 & 70 & 73 \\ y & 67 & 68 & 65 & 68 & 72 & 72 & 69 & 71 \end{array}$$

$$\textcircled{2} \quad \begin{array}{l} \text{Let } u = x - a \\ = x - 68 \end{array}$$

$$\begin{array}{l} v = y - b \\ = y - 69 \end{array}$$

x	y	u	v	uv	u ²	v ²
65	67	-3	-2	6	9	4
66	68	-2	-1	2	4	1
67	65	-1	-4	4	1	16
67	68	-1	-1	1	1	1
68	72	0	3	0	0	9
69	72	1	3	3	1	9
70	69	2	0	0	4	0
73	71	4	2	8	16	4
		$\Sigma u = 0$	$\Sigma v = 0$	$\Sigma uv = 24$	$\Sigma u^2 = 36$	$\Sigma v^2 = 44$

$$r^2 = \frac{n \sum uv - \sum u \sum v}{\sqrt{[n \sum u^2 - (\sum u)^2] * [n \sum v^2 - (\sum v)^2]}}$$

$$= \frac{8 \times 24 - 0 \times 0}{\sqrt{8 \times 36 - 0^2 * 8 \times 44 - 0^2}} = 0.6030$$

$$\begin{aligned} & \frac{192 - 0}{\sqrt{384 * 352}} = \frac{192}{\sqrt{138,336}} = \frac{192}{318,3959} \\ & \boxed{0.6030} \end{aligned}$$

PAGE NO.:

DATE:

0.8333

50/1

17/01/25
Friday

Spearman's coefficient of rank correlation

Problem ①

mathematics	38	50	42	61	43	55	67	46	72
statistics	41	64	70	75	44	58	62	56	60

$$\text{Sol.} - n = 9$$

different data

$$P_r = \frac{6 \cdot \sum d^2}{n^3 - n}$$

x	y	R ₁	R ₂	d ₂ R ₁ - R ₂	d ²
38	41	9	9	0	0
50	64	5	3	2	4
42	70	8	2	6	36
61	75	3	1	2	4
43	44	7	8	-1	1
55	56	4	6	-3	9
67	62	2	4	-2	4
46	58	6	6	0	0
72	60	1	5	-4	16
$\sum d^2 = 74$					

$$P_r = \frac{6(74)}{9^3 - 9}$$

$$P_r = \frac{444}{729 - 9}$$

$$P_r = \frac{444}{720}$$

$$= 0.3833$$

$$\frac{3}{2} \times 6$$

216

96

3

PAGE NO.:

DATE:

(2)

Height	165	167	166	170	169	172
Weight	61	60	63.5	63	61.5	64

85:- n = 6

x	y	R ₁	R ₂	d = R ₁ - R ₂	d ²
165	61	6	5	1	1
167	60	64	6	-2	4
166	63.5	45	2	3	9
170	63	2	3	-1	1
169	61.5	3	4	-1	1
172	64	1	1	0	0

$\sum d^2 = 16$

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

$$= 1 - \frac{16}{216 - 6} = 1 - \frac{16}{210} = 0.5428$$

Repeated Rank [Data with Ties]

If one rank repeat m₁ times, another rank repeat m₂ times, and third rank repeat m₃ times and so that the correlation factor [C.R] is

$$C.R = \frac{m_1^3 - m_1}{12} + \frac{m_2^3 - m_2}{12} + \frac{m_3^3 - m_3}{12} + \dots$$

The coefficient of rank correlation

$$P = 1 - \frac{6 \left[\sum d^2 + C.R \right]}{n^3 - n}$$

Problem ① The following table using repeated ranks

X	43	96	74	38	35	43	22	56	35	80
y	30	94	84	13	30	18	30	41	48	95

Sol:-

$$n = 10$$

There are similar data ranks so we use

$$P = 1 - \frac{6}{n^3 - n} [\sum d^2 + C.P]$$

$$\text{Correlation Factor} = C.P = m_1^3 - m_1 + m_2^3 - m_2 + m_3^3 - m_3 + \dots$$

$$\frac{5+6+5+5}{4} = 5.5 \quad \frac{8+9+2+7+8.5}{5} = 8.5 \quad 12 \quad 12 \quad 12$$

X	Y	R ₁	R ₂	d = R ₁ - R ₂	d ²
43	30	5.5	7	-1.5	2.25
96	94	1	2	-1	1
74	84	3	3	0	0
35	13	7	10	-3	9
35	30	8.5	7	1.5	2.25
43	18	5.5	9	-3.5	12.25
22	30	10	7	3	9
56	41	4	5	-1	1
35	48	8.5	4	4.5	20.25
80	95	2	1	1	1
					$\sum d^2 = 58$

$$\frac{6+7+8}{3} = \frac{21}{3} = 7$$

43 is repeated 2 times so
 $m_1 = 2$

35 is repeated 2 times so

$$m_2 = 2$$

30 is repeated 3 times so

$$m_3 = 3$$

$$= \frac{8-2}{12} + \frac{8-2}{12} + \frac{27-3}{12}$$

$$= \frac{6+6+24}{12} = \frac{36}{12} = 3$$

$$P = 1 - \frac{6(58) + 3}{10^3 - 10}$$

$$= 1 - \frac{366}{990} \rightarrow 0.6454 \quad 0.6303$$

~~00/01/2015~~
Monday

Introduction Regression:-

The statistical tool with the help of which it is possible to estimate [or predict] the unknown values of one variable from the known values of other variable [when the variables are correlated] is called Regression.

The Regression theory was developed by Sir Francis Galton.

$$y \text{ on } x\text{-axis}$$

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

$$\text{where } \bar{x} = a + \frac{\sum u}{n}, \quad \bar{y} = b + \frac{\sum v}{n}$$

b_{yx} is coefficient of y on x -axis

$$b_{yx} = \frac{n \sum uv - \sum x \sum v}{n \sum u^2 - (\sum u)^2}$$

x on y -axis

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

$$\text{where } \bar{x} = a + \frac{\sum u}{n}, \quad \bar{y} = b + \frac{\sum v}{n}$$

b_{xy} is coefficient of x on y -axis

$$b_{xy} = \frac{n \sum uv - \sum u \sum v}{n \sum v^2 - (\sum v)^2}$$

PAGE NO.:

DATE:

$$= \frac{\sum (x-\bar{x})(y-\bar{y})}{\sum (y-\bar{y})^2} = r, \frac{ax}{ay}$$

Properties of Regression Line

- The regression line intersect at (\bar{x}, \bar{y}) .
- If there is perfect correlation, the regression line coincide [there will be only one regression line]
- The geometric mean of the regression coefficient is equal to coefficient of correlation [numerically]. That is

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

- The regression coefficients cannot be of opposite sign.
- * If r is positive, both the regression coefficient will be positive.
- * If r is negative, both the regression coefficient will be negative.
- * If r is zero, both the regression coefficient will be zero.

Problem ① Obtain two regression equations [height - x] [weight - y]

x	153	157	162	160	170	163
y	48	50	50	49	54	59

i) Height of a person whose weight is 60 kg

ii) Weight of a person whose height is 165 kg

Sol. y on x -axis

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

where

$$\bar{x} = a + \frac{\sum u}{n} \quad \text{&} \quad \bar{y} = b + \frac{\sum v}{n}$$

$$b_{yx} = \frac{n \sum uv - \sum u \sum v}{n \sum v^2 - (\sum v)^2}$$

$$x \text{ or } y \text{ axis}$$

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

PAGE NO.:

DATE:

where $\bar{x} = a + \frac{\sum u}{n}$ & $\bar{y} = b + \frac{\sum v}{n}$

$$bxy = \frac{n \sum uv - \sum u \sum v}{n \sum v^2 - (\sum v)^2}$$

21/10/25
Monday

Let $u = x - a$ Let $v = y - b$
 $x = 160$ $y = 50$ $a = 160$ $b = 50$

x	y	u	v	u^2	v^2	uv
158	48	-7	-2	49	4	-14
157	50	-3	0	9	0	0
168	50	8	0	64	0	0
160	49	0	-1	0	1	0
170	54	10	4	100	16	40
163	53	3	3	9	9	9
				$\sum u^2 = 259$	$\sum v^2 = 90$	$\sum uv = 6463$
$\bar{x} = 160$	$\bar{y} = 50$	$\sum u = 1811$	$\sum v = 4$	$\sum u^2 = 259$	$\sum v^2 = 90$	$\sum uv = 6463$

$$\bar{x} = a + \frac{\sum u}{n} \quad \bar{y} = b + \frac{\sum v}{n}$$

$$= 160 + \frac{11}{6} \quad = 50 + \frac{4}{6}$$

$$= 161.83 \quad = 50.66$$

$$byx = \frac{n \sum uv - \sum u \sum v}{n \sum u^2 - (\sum u)^2} \quad bxy = \frac{n \sum uv - \sum u \sum v}{n \sum v^2 - (\sum v)^2}$$

$$= \frac{6(63) - 11 \times 4}{6(83) - (11)^2} = \frac{384}{6(30) - 4^2}$$

$$= \frac{378 - 44}{1886 - 121} = \frac{334}{1865}$$

$$= \frac{334}{1265} = 0.264$$

=

$$= 0.03$$

y on x -axis

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

PAGE NO.:

DATE:

$$50.66 = 0.264 (\cancel{x} - 161.83)$$

$$\begin{aligned} - 50.66 &\rightarrow 0.264 (\cancel{x} - 161.83) \\ - 50.66 &\rightarrow -0.264 \cancel{x} + 0.264 \cdot 161.83 \end{aligned}$$

$$(y - 50.66) = 0.264 (x - 161.83)$$

$$\cancel{y - 50.66} = y - 0.264x - 42.72 + 50.66$$

$$y = 0.264x + 7.94 \quad \text{--- (1)}$$

x on y -axis

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

$$(x - 161.83) = 2.03 (y - 50.66)$$

$$x = 2.03y - 102.83 + 161.83$$

$$x = 2.03y + 59 \quad \text{--- (2)}$$

i] To find weight of a person whose height is 165

$$y = ? , x = 165$$

$$y = 0.264(165) + 7.94$$

$$= 43.56 + 7.94$$

$$= 51.5$$

ii] To find height of a person whose weight is 60

$$x = ?, y = 60$$

$$x = 2.03(60) + 59$$

$$= 121.8 + 959$$

$$= 180.8$$

22/01/25
Wednesday

② Find the co-efficient of correlation & equation of the lines of regression

$$a = 3$$

$$b = 5$$

x 5 2 1 3 4

Coefficient of correlation = $\sqrt{b_{xy} * b_{yx}}$

PAGE NO.:

DATE:

Qd.

x	y	u	v	u^2	v^2	uv
5	4	2	2	4	4	4
2	5	-1	0	1	0	0
1	2	-2	-3	4	9	6
3	3	0	-2	0	4	0
4	8	1	3	1	9	3
		$\Sigma u=0$	$\Sigma v=0$	$\Sigma u^2=10$	$\Sigma v^2=26$	$\Sigma uv=13$

$$\bar{x} = a + \frac{\sum u}{n}$$

$$\bar{y} = b + \frac{\sum v}{n}$$

$$= 3 + \frac{0}{5}$$

$$= 5 + \frac{0}{5}$$

$$= 3$$

$$= 5$$

$$b_{yx} = \frac{n \sum uv - \sum u \sum v}{n \sum u^2 - (\sum u)^2}$$

$$b_{xy} = \frac{n \sum uv - \sum u \sum v}{n \sum v^2 - (\sum v)^2}$$

$$= \frac{5(13) - 0 \times 0}{5(10) - (0)^2}$$

$$= \frac{65}{5(26) - 0^2}$$

$$= \frac{65 - 0}{60 - 0}$$

$$= \frac{65}{130 - 0}$$

$$= \frac{65}{60}$$

$$= \frac{65}{130}$$

$$= 0.8$$

$$= 0.6$$

y on x-axis

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

$$y - 3 = 1.3(x - 3)$$

$$y = 1.3x - 3.9 + 3$$

$$y = 1.3x + 1.1 - 0$$

x on y-axis

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

$$x - 3 = 0.6(y - 5)$$

$$x = 0.6y - 3 + 3$$

$$x = 0.6y + 0.6 - 0$$

Coefficient of correlation

PAGE NO.:

DATE:

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

$$= \sqrt{0.5 \times 1.3}$$

$$= \sqrt{0.65}$$

$$= 0.8$$

∴ The coefficient of correlation
is highly +ve.

(3)

Calculate the coefficient of correlation & line of regression. Also estimate y corresponding to x = 6.2

x	1	2	3	4	5	6	7	8
y	9	8	10	12	11	13	14	16

Sol:

Find y on x-axis

$$\frac{4+5}{2} = \frac{9}{2} \rightarrow a$$

$$y - a = \frac{u}{2} \quad u = y - a$$

x	y	$x-4.5$	$y-11$	u^2	v^2	uv
1	9	-3.5	-2	12.25	4	7
2	8	-2.5	-3	6.25	9	7.5
3	10	-1.5	-1	2.25	1	1.5
4	12	-0.5	1	0.25	1	-0.5
5	11	0.5	0	0.25	0	0
6	13	1.5	2	2.25	4	3
7	14	2.5	3	6.25	9	7.5
8	16	3.5	5	12.25	25	17.5
		$\sum u = 0$	$\sum v = 5$	$\sum u^2 = 42$	$\sum v^2 = 53$	$\sum uv = 43.5$

$$\bar{x} = a + \frac{\sum u}{n}$$

$$\bar{y} = b + \frac{\sum v}{n}$$

$$= 4.5 + \frac{0}{8}$$

$$= 11 + \frac{5}{8}$$

$$= 4.5$$

$$= 11.625$$

$$b_{yx} = \frac{n \sum uv - \sum u \sum v}{n \sum u^2 - (\sum u)^2}$$

$$= \frac{8(43.5) - 0 \times 5}{8(42) - 0^2}$$

$$= \frac{348}{336}$$

$$= 1.0357$$

$$b_{xy} = \frac{n \sum uv - \sum u \sum v}{n \sum v^2 - (\sum v)^2}$$

$$= \frac{348}{53 - 25}$$

$$= 0.8721$$

$$= \frac{348}{424 - 25} = 0.837$$

$$= \frac{348}{399} = 0.87$$

y on x -axis

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

$$y - 11.625 = 1.0357 (x - 4.5)$$

$$y = 1.0357x - 4.6606 + 11.625 \quad \text{④}$$

$$y = 1.0357x + 6.9644 \quad \text{— ①}$$

Coefficient of correlation

$$r = \frac{b_{xy}}{\sqrt{b_{xx} b_{yy}}}$$

$$= \frac{\sqrt{0.8721 \times 1.0357}}{0.9032}$$

$$= \sqrt{0.9032}$$

$$= 0.9503$$

=?

∴ The coefficient of correlation is highly +ve.

$$\therefore x = 6.2$$

$$\begin{aligned}y &= 1.0357x + 6.9644 \\&= 1.0357(6.2) + 6.9644 \\&= 6.4213 + 6.9644 \\&= 13.3857\end{aligned}$$

∴

⇒ In a bivariate data, the regression lines are

$$4x - 5y + 33 = 0 \quad \text{①}$$

$$20x - 9y = 107 \quad \text{②}$$

find the mean of x & y coefficient of correlation [or find \bar{x}, \bar{y} & r]

Qd:

$$\text{Given } 4x - 5y = -33 \quad \text{— ①}$$

$$20x - 9y = 107 \quad \text{— ②}$$

from ① & ② (in a coefficient compared [20 & 4])

∴ 20x is greater so

x on y -axis

$$20x - 9y = 107$$

$$20x = 9y + 107$$

$$x = \frac{9}{20}y + \frac{107}{20}$$

y on axis

$$4x - 5y = -33$$

$$-5y = -33 - 4x$$

$$y = \frac{4}{5}x + \frac{33}{5}$$

Revenue

₹

make +ve

∴ The coefficient of x on y -axis

$$b_{xy} = \frac{9}{20}$$

∴ The coefficient of y on x -axis

$$b_{yx} = \frac{4}{5} \frac{4}{5}$$

The coefficient of correlation

$$r = \frac{b_{xy} \times b_{yx}}{\sqrt{\frac{4}{5}} \times \sqrt{\frac{9}{25}}}$$

$$= \sqrt{\frac{9}{25}}$$

$$= \frac{9}{5}$$

$$= 0.6$$

∴ The coefficient of correlation is highly +ve.

To find \bar{x} & \bar{y} , WKT the two regression lines are intersect at (\bar{x}, \bar{y})

$$\text{Now } 4x - 5y = -83 \quad \textcircled{1}$$

$$20x - 9y = 107 \quad \textcircled{2}$$

$$\textcircled{1} \text{ mult by 20} \Rightarrow 80x - 100y = -1660$$

$$\textcircled{2} \text{ mult by 4} \Rightarrow 80x - 36y = -428$$

$$\begin{array}{ccc} (-) & (+) & (+) \end{array}$$

$$-64y = -1088$$

$$y = \frac{-1088}{-64} = 17$$

$$y = 17$$

Substitute y with equation $\textcircled{1}$

$$4x - 5y = -83$$

$$4x - 5(17) = -83$$

$$4x - 85 = -83$$

$$4x = -83 + 85$$

$$x = \underline{5}$$

$$4$$

$$\boxed{x = 13}$$

The mean of x & y is

$$\bar{x} = x = 13$$

$$\bar{y} = y = 17$$

~~Q~~ Find the regression lines are

$$3x + 2y = 26 \quad \text{Eq}$$

$$6x + y = 31$$

Find \bar{x}, \bar{y} & r

Ques. Given $3x + 2y = 26$ - ①
 $6x + 2y = 31$ - ②

6 is greater than equation ① [3x] so

x on y-axis y on x-axis
 $6x + y = 31$
 $6x \downarrow$

23/01/25
Thursday

Introduction - Curve fit, Least square, Goodness of fit.

One of the fields of elementary mathematics / statistics where the digital computer has a good amount of contribution is that of the curve.

The basic problem in curve-fitting can be described as follows:

Let $(x_1, v_1), (x_2, v_2), \dots, (x_n, v_n)$ be a set of n observations which can be represented on xy -plane. The problem of finding a functional relation of the form $y = f(x)$ which is satisfied by the given set of points is called curve fitting & the curve $y = f(x)$ is called the curve of the best fit.

In fitting a curve to given data points, there are two possible approaches.

i] To have a graph of the approximating functions that passes exactly through the given data points.

ii] To have an approximating function that has a smooth curve, and serves as the best approximation to the actual curve. The second curve is generally employed for curve fitting and is referred as the method of least squares and this is discussed in this chapter.

i) A Straight Line $y = a + bx$
whose normal equations are

$$\Sigma \rightarrow \Sigma y = na + b \Sigma x$$

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

from this two equation 'a' & 'b' are determined.
 Substituting these values of 'a' & 'b' in $y = a + bx$
 we get 1st equation of the (Straight) line of best fit
 for the given data.

Fitting a straight line problem

- ① Fit in a straight line for the following data

x	6	7	7	8	8	8	9	9	10
y	5	5	4	5	4	3	4	3	3

Sol:- $y = a + bx$

Let $X = x - \text{center value}$

∴ Straight line becomes $y = a + bX$

Normal Equation

$$\sum \rightarrow \sum y = na + b \sum X$$

$$\sum X \rightarrow \sum Xy = a \sum X + b \sum X^2$$

∴ Let $X = x - 8$

∴ $y = a + bX$

26/01/25
Friday

X	y	$X = x - 8$	$\sum y$	$\sum X^2$
6	5	-2	36 - 16	4
7	5	-1	-5	1
7	4	-1	-4	1
8	5	0	0	0
8	4	0	0	0
9	3	1	0	0
9	4	1	4	1
9	3	1	3	1
10	3	2	6	4

$$\sum y = 36 \quad \sum X = 0 \quad \sum Xy = -6 \quad \sum X^2 = 12$$

$$\sum y = na + b \sum x$$

$$36 = 9a + b(0)$$

$$9a = 36$$

$$a = \frac{36}{9} = 4$$

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

$$12 = a(0) + b(12)$$

~~$$-6 = 12b$$~~

~~$$12b = -6$$~~

~~$$b = \frac{-6}{12} = -0.5$$~~

$$y = a + bx$$

$$= 4 + (-0.5)(x - 8)$$

$$= 4 - 0.5x + 4$$

$$y = 0.5x + 8$$

=,

② Find a straight line \rightarrow mid-term

x 1941 1921 1931 1941 1951

y 15 23 28 32 39

Sol:- Let $X = x - 1931$

$$\therefore y = a + bX$$

X	y	$X = x - 1931$	XY	X^2
1941	15	-10	-300	100
1921	23	-10	-230	100
1931	28	0	0	0
1941	32	10	320	100
1951	39	20	780	400

$$\sum y = 137 \quad \sum X = 0 \quad \sum XY = 870 \quad \sum X^2 = 1000$$

$$\Sigma y = na + b \sum X$$

$$137 = 5a + b(0)$$

$$137 = 5a$$

PAGE NO.:

DATE:

$$a = \frac{137}{5} = 27.4$$

$$\Sigma XY = a \sum X + b \sum X^2$$

$$570 = 0a(0) + b(1000)$$

$$570 = 1000b$$

$$b = \frac{570}{1000} = 0.57$$

$$y = a + bX$$

$$= 27.4 + 0.57(x - 1931)$$

$$= 27.4 + 0.57x - 1,100.67$$

$$y = 0.57x - 1,073.27$$

=,

③ Find a straight line \rightarrow mid term

x	1	2	3	4	5	6	7	8	9
y	3	4	6	5	10	9	10	12	11

S.f.

$$\text{Let } x = x - 5$$

$$\therefore y = a + bX$$

x	y	$x-a-5$	Xy	x^2
1	3	-4	-12	16
2	4	-3	-12	9
3	6	-2	-12	4
4	5	-1	-5	1
5	10	0	0	0
6	9	1	9	1
7	10	2	20	4
8	12	3	36	9
9	11	4	44	16

$$\Sigma y = 70 \quad \Sigma x = 0 \quad \Sigma xy = 68 \quad \Sigma x^2 = 60$$

$$\begin{aligned}\Sigma y &= na + b \sum x \\ 70 &= 9a + b(0) \\ 9a &= 70\end{aligned}$$

PAGE NO.:

DATE:

$$a = \frac{70}{9} = 7.77$$

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

$$68 = a(0) + b(60)$$

$$68 = 60b$$

$$b = \frac{68}{60} = 1.13$$

$$y = a + bx$$

$$= 7.77 + 1.13(x - 5)$$

$$= 7.77 + 1.13x - 5.65$$

$$= 1.13x + 2.12$$

Ans.

- ④ Find the straight line with reference to the following data.

x 20 5 10 8 15 10 6 12 10

y 24 15 17 1.22 12 6 8 10 8