

PDFZilla – Unregistered

PDFZilla - Unregistered

PDFZilla - Unregistered



Module-4

Part - I

Curve Fitting

It is the method of finding equation of a curve that approximates a given set of n data points and this equation is called best fitting equation.

Curve Fitting By The Method of Least squares

1. Fitting of a straight line
2. Fitting of a Second degree curve
(Parabola)
3. Fitting of Exponential curve

Fitting of a straight line $y = ax + b$

The normal equations of fitting $y = ax + b$ are

$$\Sigma y = a \Sigma x + nb \quad \text{and}$$

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

The normal equations for fitting a straight line $y = a + bx$ are

$$\Sigma y = na + b \Sigma x \quad \text{and}$$

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

Working procedure to fit a straight line $y = ax + b$ or $y = a + bx$

- Write the normal equations of the given curve.
- Prepare the relevant table and find the value of summation present in the normal equation and substitute.
- We find the parameters ‘ a ’ and ‘ b ’ by solving and then substitute in the given equation.

WORKED EXAMPLES

Fit a straight line $y = ax + b$ for the data

x	5	10	15	20	25
y	16	19	23	26	30

The normal equations of fitting $y = ax + b$ are

$$\Sigma y = a\Sigma x + nb \quad \text{and}$$

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

We prepare a relevant table as follows

x	y	xy	x^2
5	16	80	25
10	19	190	100
15	23	345	225
20	26	520	400
25	30	750	625
Σ	75	114	1375
		1885	

Here, $n = 5$, $\Sigma x = 75$, $\Sigma y = 114$, $\Sigma xy = 1885$ and $\Sigma x^2 = 1375$

Substituting these values in the above normal equations, we get

$$114 = 75a + 5b$$

$$98 = 1375a + 75b$$

$$114 = 75a + 5b$$

$$98 = 1375a + 75b$$

Solving these equations, we get

$$a = 0.7, \quad b = 12.3$$

∴ The equation of the best fitting straight line is $y = 0.7x + 12.3$

Fit a straight line $y = ax + b$ for the data

x	0	1	2	3	4	5	6
y	-4	-2	0	2	4	6	8

The normal equations of fitting $y = ax + b$ are

$$\Sigma y = a\Sigma x + nb \quad \text{and}$$

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

We prepare a relevant table as follows

	x	y	xy	x^2
	0	-4	0	0
	1	-2	-2	1
	2	0	0	4
	3	2	6	9
	4	4	16	16
	5	6	30	25
	6	8	48	36
Σ	21	14	98	91

Here, $n = 7$, $\Sigma x = 21$, $\Sigma y = 14$, $\Sigma xy = 98$ and $\Sigma x^2 = 91$

Substituting these values in the above normal equations, we get

$$14 = 21a + 7b$$

$$98 = 91a + 21b$$

Solving these equations, we get

$$a = 2, \quad b = -4$$

\therefore The equation of the best fitting straight line is $y = 2x - 4$

Fit a straight line $y = a + bx$ for the data

x	0	1	3	6	8
y	1	3	2	5	4

The normal equations for fitting a straight line $y = a + bx$ are

$$\Sigma y = na + b\Sigma x \quad \text{and}$$

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

We prepare a relevant table as follows

x	y	xy	x^2
0	1	0	0
1	3	3	1
3	2	6	9
6	5	30	36
8	4	32	64
Σ	18	15	110

Here, $n = 5$, $\Sigma x = 18$, $\Sigma y = 15$, $\Sigma xy = 71$ and $\Sigma x^2 = 110$

Substituting these values in the above normal equations, we get

$$15 = 5a + 18b$$

$$71 = 18a + 110b$$

$$15 = 5a + 18b$$

$$71 = 18a + 110b$$

Solving these equations, we get

$$a = 1.6, \quad b = 0.38$$

∴ The equation of the best fitting straight line is $y = 1.6 + 0.38x$

If P is the pull required to lift a load W by means of a pulley block, find a linear law of the form $P = mW + c$ connecting P and W , using the following data:

P	12	15	21	25
W	50	70	100	120

Also find P when $W = 150$ kgs

The normal equations of fitting $P = mW + c$ are

$$\Sigma P = m \Sigma W + nc \quad \text{and}$$

$$\Sigma WP = m \Sigma W^2 + c \Sigma x$$

We prepare a relevant table as follows

	W	P	WP	W^2
	50	12	600	2500
	70	15	1050	4900
	100	21	2100	10000
	120	25	3000	14400
Σ	340	73	6750	31800

Here, $n = 4$, $\Sigma W = 340$, $\Sigma P = 73$, $\Sigma WP = 6750$ and $\Sigma W^2 = 31800$

Substituting these values in the above normal equations, we get

$$73 = 340m + 4c$$

$$6750 = 31800m + 340c$$

Solving these equations, we get

$$m = 0.1879, \quad c = 2.2759$$

\therefore The best fit of a line is $P = 0.1879W + 2.2759$

When $W = 150$ kg

$$\Rightarrow P = 0.1879(150) + 2.2759 = 30.4635 \text{ kg}$$

EXERCISE

1. Fit a straight line to the following data

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

2. Fit a straight line to the following data

x	1	3	4	6	8	9	11	14
y	1	2	4	4	5	7	8	9

(VTU 2011)

3. Fit a straight line to the following data

Year x	1961	1971	1981	1991	2001
Production (in tons) y	8	10	12	10	16

and find the expected production in 2006.

4. A simply supported beam carries a concentrated load P at its mid-point. Corresponding to various values of P , the maximum deflection Y is measured. The data are given below:

P	100	120	140	160	180	200
y	0.45	0.55	0.60	0.70	0.80	0.85

Find a law of the form $Y = a + bP$

5. The results of measurement of electric resistance R of a copper bar at various temperatures $t^{\circ}\text{C}$ are listed below:

t	19	25	30	36	40	45	50
R	76	77	79	80	82	83	85

6. Fit a straight line to the following data

x	0	0.2	0.5	0.7	0.9	1.1	1.3
y	1.5	1.08	0.45	0.03	-0.39	-0.81	-1.23



Module-4

Part - 2

Fitting of a Second degree curve (Parabola) $y = ax^2 + bx + c$

The normal equations of fitting $y = ax^2 + bx + c$ are

$$\Sigma y = a\Sigma x^2 + b\Sigma x + nc$$

$$\Sigma xy = a\Sigma x^3 + b\Sigma x^2 + c\Sigma x \quad \text{and}$$

$$\Sigma x^2 y = a\Sigma x^4 + b\Sigma x^3 + c\Sigma x^2$$

The normal equations for fitting a straight line $y = a + bx + cx^2$ are

$$\Sigma y = na + b\Sigma x + c\Sigma x^2$$

$$\Sigma xy = a\Sigma x + b\Sigma x^2 + c\Sigma x^3 \quad \text{and}$$

$$\Sigma x^2 y = a\Sigma x^2 + b\Sigma x^3 + c\Sigma x^4$$

Working procedure to fit a second degree curve $y = ax^2 + bx + c$ or

$$y = a + bx + cx^2$$

- **Write the normal equations of the given curve.**
- **Prepare the relevant table and find the value of summation present in the normal equation and substitute.**
- **We find the parameters ‘ a ’, ‘ b ’ and ‘ c ’ by solving and then substitute in the given equation.**

WORKED EXAMPLES

Fit a second degree curve $y = ax^2 + bx + c$ for the data

x	1.0	1.5	2.0	2.5	3.0	3.5	4.0
y	1.1	1.3	1.6	2.0	2.7	3.4	4.1

The normal equations of fitting $y = ax^2 + bx + c$ are

$$\Sigma y = a \Sigma x^2 + b \Sigma x + nc$$

$$\Sigma xy = a \Sigma x^3 + b \Sigma x^2 + c \Sigma x \quad \text{and}$$

$$\Sigma x^2 y = a \Sigma x^4 + b \Sigma x^3 + c \Sigma x^2$$

We prepare a relevant table as follows

x	y	xy	x^2	x^3	x^4	x^2y
1.0	1.1	1.1	1	1	1	1.1
1.5	1.3	1.95	2.25	3.375	5.0625	2.925
2.0	1.6	3.2	4.0	8.0	16.0	6.4
2.5	2.0	5.0	6.25	15.625	39.0625	12.5
3.0	2.7	8.1	9.0	27.0	81.0	24.3
3.5	3.4	11.9	12.25	42.875	150.0625	41.65
4.0	4.1	16.4	16	64	256	65.6
Σ	17.5	16.2	47.65	50.75	161.875	154.475

Here, $n = 7$, $\Sigma x = 17.5$, $\Sigma y = 16.2$, $\Sigma xy = 47.65$, $\Sigma x^2 = 50.75$, $\Sigma x^3 = 161.875$, $\Sigma x^4 = 548.1875$ and $\Sigma x^2y = 154.475$

Substituting these values in the above normal equations, we get

$$16.2 = 50.75a + 17.5b + 7c$$

$$47.65 = 161.875a + 50.75b + 17.5c$$

$$154.475 = 548.1875a + 161.875b + 50.75c$$

$$16.2 = 50.75a + 17.5b + 7c$$

$$47.65 = 161.875a + 50.75b + 17.5c$$

$$154.475 = 548.1875a + 161.875b + 50.75c$$

Solving these equations, we get

$$a = 0.243, \ b = -0.193 \text{ and } c = 1.04$$

∴ The equation of the best fitting curve is $y = 0.243x^2 - 0.193x + 1.04$

Fit a second degree curve $y = a + bx + cx^2$ for the data

x	0	1	2	3	4
y	1	1.8	1.3	2.5	6.3

The normal equations for fitting a straight line $y = a + bx + cx^2$ are

$$\Sigma y = na + b\Sigma x + c\Sigma x^2$$

$$\Sigma xy = a\Sigma x + b\Sigma x^2 + c\Sigma x^3 \quad \text{and}$$

$$\Sigma x^2 y = a\Sigma x^2 + b\Sigma x^3 + c\Sigma x^4$$

We prepare a relevant table as follows

x	y	xy	x^2	x^3	x^4	x^2y
0	1	0	0	0	0	0
1	1.8	1.8	1	1	1	1.8
2	1.3	2.6	4	8	16	5.2
3	2.5	7.5	9	27	81	22.5
4	6.3	25.2	16	64	256	100.8
Σ	10	12.9	37.1	30	100	354
						130.3

Here, $n = 5$, $\Sigma x = 10$, $\Sigma y = 12.9$, $\Sigma xy = 37.1$, $\Sigma x^2 = 30$, $\Sigma x^3 = 100$, $\Sigma x^4 = 354$ and $\Sigma x^2 y = 130.3$

Substituting these values in the above normal equations, we get

$$12.9 = 4a + 10b + 30c$$

$$37.1 = 10a + 30b + 100c$$

$$130.3 = 30a + 100b + 354c$$

Solving these equations, we get

$$a = 1.42, \ b = -1.07 \text{ and } c = 0.55$$

∴ The equation of the best fitting curve is $y = 1.42 - 1.07x + 0.55 x^2$

Fit a second degree curve $y = a + bx + cx^2$ for the data

x	-3	-2	-1	0	1	2	3
y	4.63	2.11	0.67	0.09	0.63	2.15	4.58

The normal equations for fitting a straight line $y = a + bx + cx^2$ are

$$\Sigma y = na + b\Sigma x + c\Sigma x^2$$

$$\Sigma xy = a\Sigma x + b\Sigma x^2 + c\Sigma x^3 \quad \text{and}$$

$$\Sigma x^2 y = a\Sigma x^2 + b\Sigma x^3 + c\Sigma x^4$$

We prepare a relevant table as follows

x	y	xy	x^2	x^3	x^4	x^2y
-3	4.63	-13.89	9	-27	81	41.67
-2	2.11	-4.22	4	-8	16	8.44
-1	0.67	-0.67	1	-1	1	0.67
0	0.09	0	0	0	0	0
1	0.63	0.63	1	1	1	0.63
2	2.15	4.3	4	8	16	8.6
3	4.58	13.74	9	27	81	41.22
Σ	0	14.86	-0.11	28	0	196
						101.23

Here, $n = 7$, $\Sigma x = 0$, $\Sigma y = 14.86$, $\Sigma xy = -0.11$, $\Sigma x^2 = 28$, $\Sigma x^3 = 0$, $\Sigma x^4 = 196$ and $\Sigma x^2y = 101.23$

Substituting these values in the above normal equations, we get

$$14.86 = 7a + 0b + 28c$$

$$-0.11 = 0a + 28b + 0c$$

$$101.23 = 28a + 0b + 196c$$

$$\begin{aligned}14.86 &= 7a + 0b + 28c \\-0.11 &= 0a + 28b + 0c \\101.23 &= 28a + 0b + 196c\end{aligned}$$

Solving these equations, we get

$$a = 0.1329, \quad b = -0.00393 \text{ and } c = 0.4975$$

\therefore The equation of the best fitting curve is

$$y = 0.1329 - 0.00393x + 0.4975x^2$$

Example

Find the best values of a , b and c if the equation

$y = a + bx + cx^2$ is to fit most closely to the following observations:

x	-2	-1	0	1	2
y	-3.150	-1.390	0.620	2.880	5.378

The normal equations for fitting a straight line $y = a + bx + cx^2$ are

$$\Sigma y = na + b\Sigma x + c\Sigma x^2$$

$$\Sigma xy = a\Sigma x + b\Sigma x^2 + c\Sigma x^3 \quad \text{and}$$

$$\Sigma x^2 y = a\Sigma x^2 + b\Sigma x^3 + c\Sigma x^4$$

We prepare a relevant table as follows

x	y	xy	x^2	x^3	x^4	x^2y
-2	-3.150	6.3	4	-8	16	-12.6
-1	-1.390	1.39	1	-1	1	-1.39
0	0.620	0	0	0	0	0
1	2.880	2.880	1	1	1	2.880
2	5.378	10.756	4	8	16	21.512
Σ	0	4.328	21.326	10	0	34
						10.402

Here, $n = 5$, $\Sigma x = 0$, $\Sigma y = 4.328$, $\Sigma xy = 21.326$, $\Sigma x^2 = 10$, $\Sigma x^3 = 0$, $\Sigma x^4 = 34$ and $\Sigma x^2y = 10.402$

Substituting these values in the above normal equations, we get

$$4.338 = 5a + 0b + 10c$$

$$21.326 = 0a + 10b + 0c$$

$$10.402 = 10a + 0b + 34c$$

$$4.338 = 5a + 0b + 10c$$

$$21.326 = 0a + 10b + 0c$$

$$10.402 = 10a + 0b + 34c$$

Solving these equations, we get

$$a = 0.621, \ b = 2.1326 \text{ and } c = 0.1233$$

\therefore The best values of a , b and c for fitting of $y = a + bx + cx^2$ are

$$a = 0.621, \ b = 2.1326 \text{ and } c = 0.1233$$

Fit a second degree curve $y = ax^2 + bx + c$ for the data

x	1	2	3	4	5
y	10	12	13	16	19

The normal equations for fitting a straight line $y = a + bx + cx^2$ are

$$\Sigma y = na + b\Sigma x + c\Sigma x^2$$

$$\Sigma xy = a\Sigma x + b\Sigma x^2 + c\Sigma x^3 \quad \text{and}$$

$$\Sigma x^2 y = a\Sigma x^2 + b\Sigma x^3 + c\Sigma x^4$$

We prepare a relevant table as follows

x	y	xy	x^2	x^3	x^4	x^2y
1	10	10	1	1	1	10
2	12	24	4	8	16	48
3	13	39	9	27	81	117
4	16	64	16	64	256	256
5	19	95	25	125	625	475
Σ	15	61	232	55	225	906

Here, $n = 5$, $\Sigma x = 15$, $\Sigma y = 61$, $\Sigma xy = 232$, $\Sigma x^2 = 55$, $\Sigma x^3 = 225$, $\Sigma x^4 = 979$ and $\Sigma x^2y = 906$

Substituting these values in the above normal equations, we get

$$61 = 55a + 15b + 5c$$

$$232 = 225a + 55b + 15c$$

$$906 = 979a + 225b + 55c$$

$$61 = 55a + 15b + 5c$$

$$232 = 225a + 55b + 15c$$

$$906 = 979a + 225b + 55c$$

Solving these equations, we get

$$a = 0.29, \ b = 0.46 \text{ and } c = 9.43$$

∴ The equation of the best fitting curve is $y = 0.29x^2 + 0.46x + 9.43$

Fit a second degree Parabola for the following data

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

The normal equations of fitting $y = ax^2 + bx + c$ are

$$\Sigma y = a\Sigma x^2 + b\Sigma x + nc$$

$$\Sigma xy = a\Sigma x^3 + b\Sigma x^2 + c\Sigma x \quad \text{and}$$

$$\Sigma x^2 y = a\Sigma x^4 + b\Sigma x^3 + c\Sigma x^2$$

We prepare a relevant table as follows

x	y	xy	x^2	x^3	x^4	x^2y
1	2	2	1	1	1	2
2	6	12	4	8	16	24
3	7	21	9	27	81	63
4	8	32	16	64	256	128
5	10	50	25	125	625	250
6	11	66	36	216	1296	396
7	11	77	49	343	2401	539
8	10	80	64	512	4096	640
9	9	81	81	729	6561	729
Σ	45	74	421	285	2025	15333
						2771

Here, $n = 9$, $\Sigma x = 45$, $\Sigma y = 74$, $\Sigma xy = 421$, $\Sigma x^2 = 285$, $\Sigma x^3 = 2025$, $\Sigma x^4 = 15333$ and $\Sigma x^2y = 2771$

Substituting these values in the above normal equations, we get

$$74 = 285a + 45b + 9c$$

$$421 = 2025a + 285b + 45c$$

$$2771 = 15333a + 2025b + 285c$$

Solving these equations, we get

$$a = -0.2673, \quad b = 3.523 \text{ and } c = -0.9282$$

∴ The equation of the best fitting curve is

$$y = -0.2673x^2 + 3.523x - 0.9282$$

EXERCISE

1. Fit a second degree Parabola to the following data:

x	1	2	3	4	5	6	7	8	9	10
y	124	129	140	159	228	289	315	302	263	210

2. Fit a second degree Parabola to the following data:

x	2	4	6	8	10
y	3.07	12.85	31.47	57.38	91.29

(VTU 2011)

3. Find the best values of a, b, c if the equation $y = a + bx + cx^2$ is to fit most closely to the following observation.

x	1	2	3	4	5
y	10	12	13	16	19

(VTU 2013)

4. The velocity V of a liquid is known to vary with temperature according to a quadratic law $V = a + bT + cT^2$. Find the best values of a , b and c for the data:

T	1	2	3	4	5	6	7
V	2.31	2.01	3.80	1.66	1.55	1.47	1.41

5. The following table gives the results of the measurements of train resistances, V is the velocity in miles per hour, R is the resistance in pounds per ton:

V	20	40	60	80	100	120
R	5.5	9.1	14.9	22.8	33.3	46.0

If R is related to V by the relation $R = a + bV + cV^2$, find a , b and c .

6. Fit a second degree Parabola $y = ax^2 + bx + c$ to the following data:

x	10	20	30	40	50	60
y	157	179	210	252	302	361



Module-4

Part - 3

Fitting of a curve of the form $y = ax^b$

$$y = ax^b \quad \dots\dots (1)$$

Taking \log on both sides

$$\log y = \log(ax^b)$$

$$\log y = \log a + b \log x$$

$$Y = A + BX \quad \dots\dots (2)$$

where $Y = \log y$, $A = \log a$, $B = b$ and $X = \log x$

The normal equations of (2) are

$$\Sigma Y = nA + B\Sigma X \quad \text{and}$$

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

Solving these equations we obtain A and B from which $a = e^A$ and $b = B$ can be found.

Substituting these values of a and b in (1), we obtain the equation of the best fitting curve.

WORKED EXAMPLES

Fit a curve of the form $y = ax^b$ for the data

x	1	2	3	4	5	6
y	2.98	4.26	5.21	6.1	6.8	7.5

$$y = ax^b \quad \text{----- (1)}$$

Taking \log on both sides

$$\log y = \log(ax^b)$$

$$\log y = \log a + b \log x$$

$$Y = A + BX \quad \text{----- (2)}$$

where $Y = \log y$, $A = \log a$, $B = b$ and $X = \log x$

The normal equations of (2) are

$$\Sigma Y = nA + B\Sigma X \quad \text{and}$$

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

where $Y = \log y$, $A = \log a$, $B = b$ and $X = \log x$

We prepare a relevant table as follows

x	y	$X = \log x$	$Y = \log y$	XY	X^2
1	2.98	0	1.0919	0	0
2	4.26	0.6931	1.4492	1.0044	0.4804
3	5.21	1.0986	1.6506	1.8133	1.2069
4	6.1	1.3863	1.8083	2.5068	1.9218
5	6.8	1.6094	1.9169	3.0851	2.5909
6	7.5	1.7918	2.0149	3.6103	3.2105
Σ	21	6.5792	9.9318	12.0199	9.4098

Here, $n = 6$, $\Sigma X = 6.5792$, $\Sigma Y = 9.9318$, $\Sigma XY = 12.0199$ and $\Sigma X^2 = 9.4098$

Substituting these values in the above normal equations, we get

$$9.9318 = 6A + 6.5792B$$

$$12.0199 = 6.5792A + 9.4098B$$

Solving these equations, we get

$$A = 1.0912, \quad B = 0.5144$$

$$\Rightarrow a = e^A = e^{1.0912} = 2.9778 \quad \text{and } b = B = 0.5144$$

\therefore The equation of the best fitting curve is $y = (2.9778)x^{0.5144}$

Fit a curve of the form $y = ax^b$ for the data

x	1	2	3	4	5
y	0.5	2	4.5	8	12.5

$$y = ax^b \quad \text{---- (1)}$$

Taking \log on both sides

$$\log y = \log(ax^b)$$

$$\log y = \log a + b \log x$$

$$Y = A + BX \quad \text{---- (2)}$$

where $Y = \log y$, $A = \log a$, $B = b$ and $X = \log x$

The normal equations of (2) are

$$\Sigma Y = nA + B\Sigma X \quad \text{and}$$

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

where $Y = \log y$, $A = \log a$, $B = b$ and $X = \log x$

We prepare a relevant table as follows

x	y	$X = \log x$	$Y = \log y$	XY	X^2
1	0.5	0	-0.6931	0	0
2	2	0.6931	0.6931	0.4804	0.4804
3	4.5	1.0986	1.5041	1.6524	1.2069
4	8	1.3863	2.0794	2.8827	1.9218
5	12.5	1.6094	2.5257	4.0649	2.5903
Σ	15	4.7874	6.1092	9.0804	6.1993

Here, $n = 5$, $\Sigma X = 4.7874$, $\Sigma Y = 6.1092$, $\Sigma XY = 9.0804$ and $\Sigma X^2 = 6.1993$

Substituting these values in the above normal equations, we get

$$6.1092 = 5A + 4.7874B$$

$$9.0804 = 4.7874A + 6.1993B$$

Solving these equations, we get

$$A = -0.6929, \quad B = 1.9998$$

$$\Rightarrow a = e^A = e^{-0.6929} = 0.5 \quad \text{and} \quad b = B = 1.9998$$

\therefore The equation of the best fitting curve is $y = (0.5)x^{1.9998}$

An experiment gave the following values:

v (ft/min)	350	400	500	600
t (min)	61	26	7	2.6

It is known that v and t are connected by the relation $v = at^b$. Find the best possible values of a and b .

$$v = at^b \quad \text{---- (1)}$$

Taking \log on both sides

$$\log v = \log(at^b)$$

$$\log v = \log a + b \log t$$

$$Y = A + BX \quad \text{---- (2)}$$

where $Y = \log v$, $A = \log a$, $B = b$ and $X = \log t$

The normal equations of (2) are

$$\Sigma Y = nA + B\Sigma X \quad \text{and}$$

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

where $Y = \log v$, $A = \log a$, $B = b$ and $X = \log t$

We prepare a relevant table as follows

v	t	$X = \log t$	$Y = \log v$	XY	X^2
350	61	4.1109	5.8579	24.0812	16.8995
400	26	3.2581	5.9915	19.5209	10.6152
500	7	1.9459	6.2146	12.0930	3.7865
600	2.6	0.9555	6.3969	6.1122	0.9130
Σ	1850	10.2704	24.4609	61.8073	29.2142

Here, $n = 4$, $\Sigma X = 10.2704$, $\Sigma Y = 24.4609$, $\Sigma XY = 61.8073$ and $\Sigma X^2 = 29.2142$

Substituting these values in the above normal equations, we get

$$24.4609 = 4A + 10.2704B$$

$$61.8073 = 10.2704A + 29.2142B$$

Solving these equations, we get

$$A = 7.0167, \quad B = -0.3511$$

$$\Rightarrow a = e^A = e^{7.0167} = 1115.10 \quad \text{and} \quad b = B = -0.3511$$



EXERCISE

1. Fit a curve of the form $y = ax^b$ for the data

x	20	16	10	11	14
y	22	41	120	89	56

(VTU 2015)

2. Fit a curve of the form $y = ax^b$ for the data

x	1	2	4	6
y	6	4	2	2

3. Predict y at $x = 3.75$, by fitting a curve $y = ax^b$ for the data

x	1	2	3	4	5	6
y	2.98	4.26	5.21	6.10	6.80	7.50

4. Fit a curve of the form $y = ax^b$ for the data

x	1	2	3	4	5
y	0.5	2	4.5	8	12.5



Module-4

Part – 4

STATISTICAL METHODS

Mean (Arithmetic mean): If $x_1, x_2, x_3, \dots, x_n$ are set of n values of a variate, then the mean is given by

$$\bar{x} = \frac{x_1 + x_2 + x_3 + \dots + x_n}{n} = \frac{\sum x_i}{n}.$$

In a frequency distribution, if $x_1, x_2, x_3, \dots, x_n$ be the mid-values of the class-intervals having frequencies $f_1, f_2, f_3, \dots, f_n$ respectively, we have

$$\bar{x} = \frac{x_1 f_1 + x_2 f_2 + x_3 f_3 + \dots + x_n f_n}{f_1 + f_2 + f_3 + \dots + f_n} = \frac{\sum x_i f_i}{\sum f_i}$$

Variance: If a variate x take the values $x_1, x_2, x_3, \dots, x_n$ then the variance V is defined as follows:

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

Also for a grouped data,

$$V = \frac{\sum f(x - \bar{x})^2}{\sum f}$$

Standard deviation (S.D):

$$\sigma = \sqrt{V} \text{ or } \sigma^2 = V$$

\Rightarrow

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} = \sqrt{\frac{\sum X^2}{n}} \text{ where } X = x - \bar{x}$$

or

$$\sigma^2 = \frac{\sum (x - \bar{x})^2}{n} = \frac{\sum X^2}{n}$$

Alternative formula for σ^2

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

Correlation

- The two variables x and y are related in such a way that an increase in one is accompanied by an increase or decrease in the other is called co-variation.
- Co-variation of two independent magnitudes is known as correlation.

Correlation

- **Correlation is Positive when the values increase together and Correlation is Negative when one value decreases as the other increases.**

Coefficient of correlation

The numerical measure of correlation between two variables x and y is known as **Karl Pearson's coefficient of correlation** or **simply coefficient of correlation** and is denoted by r and is defined by

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

$$r = \frac{\Sigma XY}{n\sigma_x\sigma_y}$$

Where X = deviation from the mean = $x - \bar{x}$, Y = deviation from the mean = $y - \bar{y}$,
 n = number of values of the two variables, σ_x = S.D. of x -series,
 σ_y = S.D. of y -series.

Substituting the value of σ_x and σ_y in the above formula, we get

$$r = \frac{\Sigma XY}{\sqrt{\Sigma X^2 \Sigma Y^2}}$$

Another form of the above formula is

$$r = \frac{n\Sigma xy - \Sigma x \Sigma y}{\sqrt{\left\{n\Sigma x^2 - (\Sigma x)^2\right\} \times \left\{n\Sigma y^2 - (\Sigma y)^2\right\}}}$$

Property of coefficient of correlation

The coefficient of correlation numerically does not exceed '1',

i.e., $-1 \leq r \leq 1$.

Formula for correlation coefficient

The formula to compute coefficient of correlation is

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

Here,

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

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

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

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

$$\sigma_{x-y}^2 = \frac{\sum (x-y)^2}{n} - \left(\bar{(x-y)}\right)^2$$

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

Regression

It is an estimation of one independent variable in terms of the other.

Example: The best fitting straight line of the form $y = a + bx$ is called the regression of y on x and $x = a + by$ is called the regression of x on y .

Equation of the Regression lines

The regression line of y on x is

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

The regression line of x on y is

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

Here $r \frac{\sigma_y}{\sigma_x}$ and $r \frac{\sigma_x}{\sigma_y}$ are the regression coefficients. Their product is r^2

Here $r \frac{\sigma_y}{\sigma_x}$ and $r \frac{\sigma_x}{\sigma_y}$ are the regression coefficients. Their product is r^2

Thus 'r' is the geometric mean of regression coefficients,

$$\text{i.e., } r = \pm \sqrt{(\text{coefficient of } x) \times (\text{coefficient of } y)}$$

Note: If both the coefficients are positive then we have to consider *r value as* positive and if both coefficients are negative then we have to consider *r value as* negative.

Angle between the lines of regression

The angle between the lines of regression is

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

Proof:

If m_1 and m_2 are the slopes of two lines then the angle between the lines is given by

$$\tan \theta = \left| \frac{m_1 - m_2}{1 + m_1 m_2} \right| \quad \text{---- (1)}$$

The lines of regression are

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

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

The slopes of regression lines are $m_1 = \frac{r\sigma_y}{\sigma_x}$ and $m_2 = \frac{\sigma_y}{r\sigma_x}$

Using these in (1), we get

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

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

\therefore The angle between the lines of regression is

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

Its significance when $r = 0$ and $r = \pm 1$

Case – 1: If $r = 0$ then $\tan \theta = \infty$

\Rightarrow

$$\theta = \frac{\pi}{2}$$

\therefore The lines are perpendicular

Case – 2: If $r = \pm 1$ then $\tan \theta = 0$

\Rightarrow

$$\theta = 0$$

\therefore The lines are parallel

Example 1: Find the coefficient of correlation and the lines of regression for the data

x	1	2	3	4	5
y	2	5	3	8	7

We have,

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

We prepare a relevant table as follows

x	y	$x - y$	x^2	y^2	$(x - y)^2$
1	2	-1	1	4	1
2	5	-3	4	25	9
3	3	0	9	9	0
4	8	-4	16	64	16
5	7	-2	25	49	4
Σ	15	-10	55	151	30

Here, $n = 5$, $\Sigma x = 15$, $\Sigma y = 25$, $\Sigma(x-y) = -10$, $\Sigma x^2 = 55$, $\Sigma y^2 = 151$ and $\Sigma(x-y)^2 = 30$.

$$\bar{x} = \frac{\Sigma x}{n} = \frac{15}{5} = 3 \quad \bar{y} = \frac{\Sigma y}{n} = \frac{25}{5} = 5$$

$$\overline{(x-y)} = \frac{\Sigma(x-y)}{n} = -\frac{10}{5} = -2$$

$$\sigma_x^2 = \frac{\Sigma x^2}{n} - (\bar{x})^2 = \frac{55}{5} - (3)^2 = 2$$

$$\sigma_y^2 = \frac{\Sigma y^2}{n} - (\bar{y})^2 = \frac{151}{5} - (5)^2 = 5.2$$

$$\sigma_{x-y}^2 = \frac{\Sigma(x-y)^2}{n} - (\overline{(x-y)})^2 = \frac{30}{5} - (-2)^2 = 2$$

Substituting these values in r , we get

$$r = \frac{2 + 5.2 - 2}{2(\sqrt{2})(\sqrt{5.2})} = 0.81$$

Also, we have the equations of the regression lines are

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

$$y - 5 = (0.81) \frac{\sqrt{5.2}}{\sqrt{2}} (x - 3) \text{ and } x - 3 = (0.81) \frac{\sqrt{2}}{\sqrt{5.2}} (y - 5)$$

$$y - 5 = 1.306(x - 3) \text{ and } x - 3 = 0.502(y - 5)$$

$$y = 1.306x + 1.082 \text{ and } x = 0.502y + 0.49$$

These are the lines of regression.

Example 2: Find the coefficient of correlation and the lines of regression for the data

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

We have,

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

We prepare a relevant table as follows

x	y	$x - y$	x^2	y^2	$(x - y)^2$
1	9	-8	1	81	64
2	8	-6	4	64	36
3	10	-7	9	100	49
4	12	-8	16	144	64
5	11	-6	25	121	36
6	13	-7	36	169	49
7	14	-7	49	196	49
Σ	28	-49	140	875	347

Here, $n = 7$, $\Sigma x = 28$, $\Sigma y = 77$, $\Sigma(x-y) = -49$, $\Sigma x^2 = 140$, $\Sigma y^2 = 875$ and $\Sigma(x-y)^2 = 347$

$$\bar{x} = \frac{\Sigma x}{n} = \frac{28}{7} = 4 \quad \bar{y} = \frac{\Sigma y}{n} = \frac{77}{7} = 11$$

$$\overline{(x-y)} = \frac{\Sigma(x-y)}{n} = -\frac{49}{7} = -7$$

$$\sigma_x^2 = \frac{\Sigma x^2}{n} - (\bar{x})^2 = \frac{140}{7} - (4)^2 = 4$$

$$\sigma_y^2 = \frac{\Sigma y^2}{n} - (\bar{y})^2 = \frac{875}{7} - (11)^2 = 4$$

$$\sigma_{x-y}^2 = \frac{\Sigma(x-y)^2}{n} - (\overline{(x-y)})^2 = \frac{347}{7} - (-7)^2 = 0.57$$

Substituting these values in r , we get

$$r = \frac{4 + 4 - 0.57}{2(2)(2)} = 0.93$$

Also, we have the equations of the regression lines are

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

$$y - 11 = (0.93) \frac{2}{2} (x - 4) \text{ and } x - 4 = (0.93) \frac{2}{2} (y - 11)$$

$$y - 11 = 0.93(x - 4) \text{ and } x - 4 = 0.93(y - 11)$$

$$y = 0.93x + 7.28 \text{ and } x = 0.93y - 6.23$$

These are the lines of regression.

Example 3: Find the coefficient of correlation and the lines of regression for the data

x	1	2	3	4	5	6	7	8	9	10
y	10	12	16	28	25	36	41	49	40	50

We have,

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

We prepare a relevant table as follows

x	y	x - y	x²	y²	(x - y)²
1	9	-9	1	100	81
2	12	-10	4	144	100
3	16	-13	9	256	169
4	28	-24	16	784	576
5	25	-20	25	625	400
6	36	-30	36	1296	900
7	41	-34	49	1681	1156
8	49	-41	64	2401	1681
9	40	-31	81	1600	961
10	50	-40	100	2500	1600
Σ	55	-252	385	11387	7624

Here, $n = 10$, $\Sigma x = 55$, $\Sigma y = 307$, $\Sigma(x-y) = -252$, $\Sigma x^2 = 385$, $\Sigma y^2 = 11387$ and $\Sigma(x-y)^2 = 7624$

$$\bar{x} = \frac{\Sigma x}{n} = \frac{55}{10} = 5.5 \quad \bar{y} = \frac{\Sigma y}{n} = \frac{307}{10} = 30.7$$

$$\overline{(x-y)} = \frac{\Sigma(x-y)}{n} = -\frac{252}{10} = -25.2$$

$$\sigma_x^2 = \frac{\Sigma x^2}{n} - (\bar{x})^2 = \frac{385}{10} - (5.5)^2 = 8.25$$

$$\sigma_y^2 = \frac{\Sigma y^2}{n} - (\bar{y})^2 = \frac{11387}{10} - (30.7)^2 = 196.21$$

$$\sigma_{x-y}^2 = \frac{\Sigma(x-y)^2}{n} - (\overline{(x-y)})^2 = \frac{7624}{10} - (-25.2)^2 = 127.36$$

Substituting these values in r , we get

$$r = \frac{8.25 + 196.21 - 127.36}{2(\sqrt{8.25})(\sqrt{196.21})} = 0.96$$

Also, we have the equations of the regression lines are

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

$$y - 30.7 = (0.96) \frac{\sqrt{196.21}}{\sqrt{8.25}} (x - 5.5) \text{ and } x - 5.5 = (0.96) \frac{\sqrt{8.25}}{\sqrt{196.21}} (y - 30.7)$$

$$y = 4.686x + 4.927 \text{ and } x = 0.197y - 0.548$$

These are the lines of regression.

Example**Find the lines of regression for the following data**

Ages of cars (in years)	2	4	6	7	8	10	12
Annual Maintenance cost (in hundreds)	16	15	18	19	17	21	20

Hence estimate the maintenance cost if the age of a car is 9 years and the age of a car if the maintenance cost is Rs.1550/-

We prepare a relevant table as follows

x	y	$x - y$	x^2	y^2	$(x - y)^2$
2	16	-14	4	256	196
4	15	-11	16	225	121
6	18	-12	36	324	144
7	19	-12	49	361	144
8	17	-9	64	289	81
10	21	-11	100	441	121
12	20	-8	144	400	64
Σ	49	126	413	2296	871

Here, $n = 7$, $\Sigma x = 49$, $\Sigma y = 126$, $\Sigma(x - y) = -77$, $\Sigma x^2 = 413$, $\Sigma y^2 = 2296$ and $\Sigma(x - y)^2 = 871$

$$\bar{x} = \frac{\Sigma x}{n} = \frac{49}{7} = 7$$

$$\bar{y} = \frac{\Sigma y}{n} = \frac{126}{7} = 18$$

$$\overline{(x - y)} = \frac{\Sigma(x - y)}{n} = -\frac{77}{7} = -11$$

$$\sigma_x^2 = \frac{\Sigma x^2}{n} - (\bar{x})^2 = \frac{413}{7} - (7)^2 = 10$$

$$\sigma_y^2 = \frac{\Sigma y^2}{n} - (\bar{y})^2 = \frac{2296}{7} - (18)^2 = 4$$

$$\sigma_{x-y}^2 = \frac{\Sigma(x - y)^2}{n} - (\overline{(x - y)})^2 = \frac{871}{7} - (-11)^2 = 3.429$$

Substituting these values in r , we get

$$r = \frac{10 + 4 - 3.429}{2(\sqrt{10})(\sqrt{4})} = 0.8357$$

Also, we have the equations of the regression lines are

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

$$y - 18 = (0.8357) \frac{2}{\sqrt{10}} (x - 7) \text{ and } x - 7 = (0.8357) \frac{\sqrt{10}}{2} (y - 18)$$

$$y = 0.52857x + 14.3 \text{ and } x = 1.321429y - 16.786$$

These are the lines of regression.

Also,

- (i) If the age of the car $x = 9$ then the maintenance cost y is given by the line of regression of y on x is

$$y = 0.52857 (9) + 14.3 = 19.06 \text{ hundreds.}$$

∴ The maintenance cost = Rs.1906/-

- (ii) If the maintenance cost y is Rs.1550/- i.e., 15.5 hundreds then the age of the car x is given by the line of regression of x on y is

$$x = 1.321429 (15.5) - 16.786 = 3.7 \text{ years}$$

∴ Age of car = 3.7 years.

Example : In the following table are recorded data showing the test scores made by salesmen on an intelligence test and their weekly sales:

Salesmen	1	2	3	4	5	6	7	8	9	10
Test scores	40	70	50	60	80	50	90	40	60	60
Sales	2.5	6.0	4.5	5.0	4.5	2.0	5.5	3.0	4.5	3.0

Calculate the regression line of sales on test scores and estimate the most probable weekly sales volume if a salesman makes a score of 70.

Example In a partially destroyed laboratory record of correlation data, the following results only are available:

Variance of x is 9 and the regression equations are $4x - 5y + 33 = 0$ and $20x - 9y = 107$. Calculate (i) the mean values of x and y , (ii) standard deviation of y and (iii) the coefficient of correlation between x and y .

We know that the regression line passes through (\bar{x}, \bar{y})

$$\therefore 4\bar{x} - 5\bar{y} = -33$$

$$20\bar{x} - 9\bar{y} = 107$$

Solving these equations, we get

$$\bar{x} = 13 \text{ and } \bar{y} = 17$$

We rewrite the given regression equations as

$$y = \frac{4}{5}x + \frac{33}{5} \quad \text{and} \quad x = \frac{9}{20}y + \frac{107}{20} \quad \text{---- (1)}$$

We have, the coefficient of correlation between x and y is

$$r = \pm \sqrt{(\text{coefficient of } x)(\text{coefficient of } y)}$$

$$r = \pm \sqrt{\left(\frac{4}{5}\right)\left(\frac{9}{20}\right)} = \pm \frac{3}{5}$$

$$r = \frac{3}{5} \quad (\because \text{both the coefficients are positive})$$

Given that, $\sigma_x^2 = 9 \Rightarrow \sigma_x = 3$

We have,

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

Comparing (2) with (1),

$$\text{we have } r \frac{\sigma_y}{\sigma_x} = \frac{4}{5}$$

we have $r \frac{\sigma_y}{\sigma_x} = \frac{4}{5}$

$$\sigma_y = \frac{4}{5} \times \frac{\sigma_x}{r}$$

$$= \frac{4}{5} \times \frac{3}{(3/5)} = 4$$

Example

Compute \bar{x} , \bar{y} and r from the equations of the regression lines $2x + 3y + 1 = 0$ and $x + 6y = 4$.

We know that the regression line passes through (\bar{x}, \bar{y})

$$\therefore 2\bar{x} + 3\bar{y} = -1$$

$$\bar{x} + 6\bar{y} = 4$$

Solving these equations, we get

$$\bar{x} = -2 \quad \text{and} \quad \bar{y} = 1$$

We rewrite the given regression equations as

$$y = -\frac{2}{3}x - \frac{1}{3} \quad \text{and} \quad x = -6y + 4$$

We have, the coefficient of correlation between x and y is

$$r = \pm \sqrt{(\text{coefficient of } x)(\text{coefficient of } y)}$$

$$r = \pm \sqrt{\left(-\frac{2}{3}\right)(-6)} = \pm \sqrt{4} = \pm 2$$

which is not possible because $-1 \leq r \leq 1$.

Hence, we rewrite the given regression equations in other form as

$$x = -\frac{3}{2}y - \frac{1}{2} \quad \text{and} \quad y = -\frac{1}{6}x + \frac{2}{3}$$

We have, the coefficient of correlation between x and y is

$$r = \pm \sqrt{(\text{coefficient of } x)(\text{coefficient of } y)}$$

$$r = \pm \sqrt{\left(-\frac{3}{2}\right)\left(-\frac{1}{6}\right)} = \pm \frac{1}{2}$$

$$r = -\frac{1}{2} \quad (\because \text{both the coefficients are negative})$$

Example If $2x - 3y = 0$ and $3x - 2y = 5$ are the lines of regression of the variables x and y . Find the following:

- (i) Mean of x and y
- (ii) Coefficient of correlation between x and y .
- (iii) Angle between the lines of regression.
- (iv) Standard deviation of x when variance of y is 2.

We know that the regression line passes through (\bar{x}, \bar{y})

$$\therefore 2\bar{x} - 3\bar{y} = 0$$

$$3\bar{x} - 2\bar{y} = 5$$

Solving these equations, we get

$$\bar{x} = 3 \quad \text{and} \quad \bar{y} = 2$$

We rewrite the given regression equations as

$$y = \frac{2}{3}x \quad \text{and} \quad x = \frac{2}{3}y + 5 \quad \text{-----(1)}$$

We have, the coefficient of correlation between x and y is

$$r = \pm \sqrt{(\text{coefficient of } x)(\text{coefficient of } y)}$$

$$r = \pm \sqrt{\left(\frac{2}{3}\right)\left(\frac{2}{3}\right)} = \pm \frac{2}{3}$$

$$r = \frac{2}{3} \quad (\because \text{both the coefficients are positive})$$

The angle between the lines of regression is $\tan \theta = \left| \frac{m_1 - m_2}{1 + m_1 m_2} \right|$

Now, $m_1 = -\frac{\text{coefficient of } x}{\text{coefficient of } y} = -\frac{2}{-3} = \frac{2}{3}$ and

$$m_2 = -\frac{\text{coefficient of } x}{\text{coefficient of } y} = -\frac{3}{-2} = \frac{3}{2}$$

$$\therefore \tan \theta = \left| \frac{(2/3) - (3/2)}{1 + (2/3)(3/2)} \right| = \left| \frac{-(5/6)}{2} \right| = \frac{5}{12}$$

$$\Rightarrow \theta = \tan^{-1} \left(\frac{5}{12} \right)$$

Given that, $\sigma_y^2 = 2 \Rightarrow \sigma_y = \sqrt{2}$

We have, $y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$ ----- (2)

Comparing (2) with (1),

we have $r \frac{\sigma_y}{\sigma_x} = \frac{2}{3}$

we have $r \frac{\sigma_y}{\sigma_x} = \frac{2}{3}$

$$\left(\frac{2}{3}\right) \frac{\sqrt{2}}{\sigma_x} = \frac{2}{3}$$

$$\sigma_x = \sqrt{2}$$

Example While calculating correlation coefficient between two variables x and y from 25 pairs of observations, the following results were obtained: $n = 25$, $\Sigma x = 125$, $\Sigma y = 100$, $\Sigma x^2 = 650$, $\Sigma y^2 = 460$, $\Sigma xy = 508$. Later it was

$x \quad y$

discovered at the time of checking that the pairs of values 8 12 were

6 8

$x \quad y$

copied down as 6 14 Obtain the correct value of correlation coefficient.

8 6

To get the correct results, we subtract the incorrect values and add the corresponding correct values.

∴ The correct results are

$$n = 25,$$

$$\Sigma x = 125 - 6 - 8 + 8 + 6 = 125,$$

$$\Sigma y = 100 - 14 - 6 + 12 + 8 = 100,$$

$$\Sigma x^2 = 650 - 6^2 - 8^2 + 8^2 + 6^2 = 650,$$

$$\Sigma y^2 = 460 - 14^2 - 6^2 + 12^2 + 8^2 = 436,$$

$$\Sigma xy = 508 - 6 \times 14 - 8 \times 6 + 8 \times 12 + 6 \times 8 = 520$$

We have,

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

$$\Rightarrow r = \frac{25 \times 520 - 125 \times 100}{\sqrt{\{25 \times 650 - (125)^2\} \times \{25 \times 436 - (100)^2\}}} = \frac{20}{\sqrt{25 \times 36}} = \frac{2}{3}$$

EXERCISE

- Find two lines of regression and coefficient of correlation for the data
 $n=18, \Sigma x=12, \Sigma y=18, \Sigma x^2=6, \Sigma y^2=96, \Sigma xy=48.$
- Find the coefficient of correlation and the lines of regression for the data

x	2	4	6	8	10
y	5	7	9	8	11

- Find the coefficient of correlation from the following data

x	78	89	97	69	59	79	68	57
y	125	137	156	112	107	138	123	108

- Find the coefficient of correlation for the data

x	21	23	30	54	57	58	72	78	87	90
y	60	71	72	83	110	84	100	92	113	135

- If the coefficient of correlation between two variables x and y is 0.5 and the acute angle between their lines of regression is $\tan^{-1}\left(\frac{3}{8}\right)$, show that

$$\sigma_x = \frac{1}{2} \sigma_y.$$

(VTU 2004)

6. Two random variables have the regression lines with equations $3x + 2y = 26$ and $6x + y = 31$. Find the mean values and the correlation coefficient between x and y .
7. The regression equations of two variables x and y are $x = 0.7y + 5.2$ and $y = 0.3x + 2.8$. Find the means of the variables and the coefficient of correlation between them.
8. In a partially destroyed laboratory data, only the equations giving the two lines of regression of y on x and x on y are available and are respectively, $7x - 16y + 9 = 0$ and $5y - 4x - 3 = 0$. Calculate the coefficient of correlation and mean values of x and y .
9. The two regression equations of the variables x and y are $x = 19.13 - 0.87y$ and $y = 11.64 - 0.50x$. Find the mean values of x and y and the correlation coefficient between x and y .

(VTU 2004)

10. Psychological tests of intelligence and of engineering ability were applied to 10 students. Here is a record of ungrouped data showing intelligence ratio (I.R) and engineering ratio (E.R). Calculate the coefficient of correlation.

Student	A	B	C	D	E	F	G	H	I	J
I.R	105	104	102	101	100	99	98	96	93	92
E.R	101	103	100	98	95	96	104	92	97	94



Module-4

Part – 5

Rank Correlation

Rank Correlation

- A group of individuals may be arranged in order to merit with respect to some characteristic. The same group would give different orders for different characteristics.
- Considering the orders corresponding to two characteristics A and B, the correlation between these n pairs of ranks is called the rank correlation in the characteristics A and B for that group of individuals.

Let x_i and y_i be the ranks of the i^{th} individuals in A and B respectively. Assuming that no two individuals are bracketed equal in either case, each of the variables taking the values $1, 2, \dots, n$,

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

If X and Y be the deviations of x and y from their means, then

$$\begin{aligned}\sum X_i^2 &= \sum (x_i - \bar{x})^2 = \sum x_i^2 + n(\bar{x})^2 - 2\bar{x}\sum x_i \\ &= \sum n^2 + \frac{n(n+1)^2}{4} - 2\left(\frac{n+1}{2}\right)\sum n \\ &= \frac{n(n+1)(2n+1)}{6} + \frac{n(n+1)^2}{4} - \frac{n(n+1)^2}{2} \\ &= \frac{1}{12}(n^3 - n)\end{aligned}$$

$$\text{Similarly, } \Sigma Y_i^2 = \frac{1}{12} (n^3 - n)$$

Now, let $d_i = x_i - y_i$ so that $d_i = (x_i - \bar{x}) - (y_i - \bar{y}) = X_i - Y_i$

$$\therefore \Sigma d_i^2 = \Sigma X_i^2 + \Sigma Y_i^2 - 2 \Sigma X_i Y_i$$

$$\Rightarrow \Sigma X_i Y_i = \frac{1}{2} [\Sigma X_i^2 + \Sigma Y_i^2 - \Sigma d_i^2]$$

$$\Rightarrow \Sigma X_i Y_i = \frac{1}{12} (n^3 - n) - \frac{1}{2} \Sigma d_i^2$$

Hence the correlation coefficient between these variables is

$$r = \frac{\sum X_i Y_i}{\sqrt{\sum X_i^2 \sum Y_i^2}} = \frac{\frac{1}{12}(n^3 - n) - \frac{1}{2} \sum d_i^2}{\frac{1}{12}(n^3 - n)}$$

\Rightarrow

$$r = 1 - \frac{6 \sum d_i^2}{(n^3 - n)}$$

This formula is called the rank correlation coefficient or Spearman's Rank Correlation Coefficient and is denoted by ρ .

$$\text{i.e., } \rho = 1 - \frac{6 \sum d_i^2}{(n^3 - n)}$$

- Spearman's rank is probably one of the most useful statistical tests that we can do in Geography to prove a relationship between two different sets of data.
- The Spearman rank correlation coefficient, *is the non-parametric* version of the Pearson correlation coefficient. Data must be ordinal, interval or ratio. Spearman's returns a value from –1 to 1, where:

- +1 = a perfect positive correlation between ranks
- 1 = a perfect negative correlation between ranks
- 0 = no correlation between ranks.

Note: To assign the rank for the given set of values, order the scores from greatest to smallest; assign the rank 1 to the highest score, 2 to the next highest and so on. If ranks are tied, i.e., Tied ranks are where two items in a column have the same rank then each tied data point assigned a mean rank.

Example 4.4.1: Calculate the rank correlation coefficient and comment on the following data on sunflower:

Height of a sunflower (in cm)	183	134	234	256	190	89	112
Width of the stem (in mm)	21	14	24	32	29	18	20

We have, the rank correlation coefficient is

$$\rho = 1 - \frac{6\sum d_i^2}{(n^3 - n)}$$

We prepare a relevant table as follows

Height of a sunflower (in cm)	Rank x_i	Width of the stem (in mm)	Rank y_i	$d_i = x_i - y_i$	d_i^2
183	4	21	4	0	0
134	5	14	7	-2	4
234	2	24	3	-1	1
256	1	32	1	0	0
190	3	29	2	1	1
89	7	18	6	1	1
112	6	20	5	1	1

$\Sigma \quad 8$

Here, $n = 7$, $\Sigma d_i^2 = 8$

$$\rho = 1 - \frac{6 \Sigma d_i^2}{(n^3 - n)} = 1 - \frac{6(8)}{(7^3 - 7)} = 1 - 0.143 = 0.857$$

As ρ is close to one, we can conclude that the wider the stem to higher the sunflower grows.

Example 4.4.2: The scores for 9 students in Physics and Maths are as follows:

Physics:	35	23	47	17	10	43	9	6	28
Mathematics:	30	33	45	23	8	49	12	4	31

Compute the ranks of students in the two subjects and compute the Spearman's rank correlation.

We have, the rank correlation coefficient is

$$\rho = 1 - \frac{6\sum d_i^2}{(n^3 - n)}$$

We prepare a relevant table as follows

Physics	Rank x_i	Mathematics	Rank y_i	$d_i = x_i - y_i$	d_i^2
35	3	30	5	-2	4
23	5	33	3	2	4
47	1	45	2	-1	1
17	6	23	6	0	0
10	7	8	8	-1	1
43	2	49	1	1	1
9	8	12	7	1	1
6	9	4	9	0	0
28	4	31	4	0	0
					Σ 12

Here, $n = 9$, $\Sigma d_i^2 = 12$

$$\rho = 1 - \frac{6 \Sigma d_i^2}{(n^3 - n)} = 1 - \frac{6(12)}{(9^3 - 9)} = 1 - 0.1 = 0.9$$

Example 4.4.3: Ten participants in a contest are ranked by two judges as follows:

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

Calculate the rank correlation coefficient.

We have, the rank correlation coefficient is

$$\rho = 1 - \frac{6\sum d_i^2}{(n^3 - n)}$$

We prepare a relevant table as follows

Rank x_i	Rank y_i	$d_i = x_i - y_i$	d_i^2
1	6	-5	25
6	4	2	4
5	9	-4	16
10	8	2	4
3	1	2	4
2	2	0	0
4	3	1	1
9	10	-1	1
7	5	2	4
8	7	1	1
Σ			60

Here, $n = 10$, $\Sigma d_i^2 = 60$

$$\rho = 1 - \frac{6 \Sigma d_i^2}{(n^3 - n)} = 1 - \frac{6(60)}{(10^3 - 10)} = 1 - 0.36 = 0.64$$

Example 4.4.4: Three judges A, B and C give the following ranks. Find which pair of judges has common approach

A	1	6	5	10	3	2	4	9	7	8
B	3	5	8	4	7	10	2	1	6	9
C	6	4	9	8	1	2	3	10	5	7

We have, the rank correlation coefficient is

$$\rho = 1 - \frac{6\sum d_i^2}{(n^3 - n)}$$

We prepare a relevant table as follows

Ranks by A x_i	Ranks by B y_i	Ranks by C z_i	d_1 $x_i - y_i$	d_2 $y_i - z_i$	d_3 $z_i - x_i$	d_1^2	d_2^2	d_3^2
1	3	6	-2	-3	5	4	9	25
6	5	4	1	1	-2	1	1	4
5	8	9	-3	-1	4	9	1	16
10	4	8	6	-4	-2	36	16	4
3	7	1	-4	6	-2	16	36	4
2	10	2	-8	8	0	64	64	0
4	2	3	2	-1	-1	4	1	1
9	1	10	8	-9	1	64	81	1
7	6	5	1	1	-2	1	1	4
8	9	7	-1	2	1	1	4	1
Σ			200	214	60			

Here, $n = 10$

$$\rho(x, y) = 1 - \frac{6 \Sigma d_1^2}{\left(n^3 - n\right)} = 1 - \frac{6(200)}{\left(10^3 - 10\right)} = -0.2$$

$$\rho(y, z) = 1 - \frac{6 \sum d_2^2}{(n^3 - n)} = 1 - \frac{6(214)}{(10^3 - 10)} = -0.3$$

$$\rho(z, x) = 1 - \frac{6 \sum d_3^2}{(n^3 - n)} = 1 - \frac{6(60)}{(10^3 - 10)} = 0.6$$

Since $\rho(z, x)$ is maximum, the pair of judges A and C have the nearest common approach.

EXERCISE

1. Find the rank correlation for the following data:

x	56	42	72	36	63	47	55	49	38	42	68	60
y	147	125	160	118	149	128	150	145	115	140	152	155

2. Calculate the rank correlation coefficient from the following data showing ranks of 10 students in two subjects:

Maths:	3	8	9	2	7	10	4	6	1	5
Physics:	5	9	10	1	8	7	3	4	2	6

3. Find the rank correlation coefficient for the following data

x	68	64	75	50	64	80	75	40	55	64
y	62	58	68	45	81	60	68	48	50	70

(VTU 2018)

Mat 41

Module - 4

Statistical Methods

Mean (Arithmetic mean) :-

If x_1, x_2, \dots, x_n be a set of n values of a variate x_i , the mean denoted by \bar{x} is defined as follows.

$$\bar{x} = \frac{\sum x}{n} \quad \text{or} \quad \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

for a grouped data in the form of a frequency distribution,

$$\bar{x} = \frac{\sum f_i x_i}{\sum f_i} \quad \text{or} \quad \bar{x} = \frac{\sum f_i x_i}{\sum f_i}$$

where f_i 's are the frequency of the classes having corresponding midpoint x_i .

Variance (V) and Standard deviation (SD) :-

If a variate x_i take values x_1, x_2, \dots, x_n the variance (V) is defined as follows.

$$V = \frac{\sum (x - \bar{x})^2}{n} \quad \text{or} \quad V = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$$

Also for a grouped data

$$V = \frac{\sum f_i (x - \bar{x})^2}{\sum f_i} \quad \text{or} \quad V = \frac{\sum_{i=1}^n f_i (x_i - \bar{x})^2}{\sum_{i=1}^n f_i}$$

Standard deviation (SD), $\sigma = \sqrt{V}$ or $\sigma^2 = V$

Alternative expression for σ^2 :-

$$\text{consider, } \sigma^2 = \frac{1}{n} \sum (x - \bar{x})^2$$

$$\begin{aligned} \sigma^2 &= \frac{1}{n} \sum [x^2 + 2x\bar{x} + (\bar{x})^2] \\ &= \frac{\sum x^2}{n} - 2(\bar{x})\bar{x} + \frac{n(\bar{x})^2}{n} \end{aligned}$$

Here, $\frac{\sum x}{n} = \bar{x}$ & $(\bar{x})^2$ being a constant added
n times gives $n(\bar{x})^2$.

$$\text{i.e. } \sigma^2 = \frac{\sum x^2}{n} - 2(\bar{x})^2 + (\bar{x})^2$$

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

for a grouped data the expression will be of
the form

$$\sigma^2 = \frac{\sum f x^2}{\sum f} - (\bar{x})^2$$

Example :-

① we shall find the mean and standard deviation of a set of observations 6, 8, 7, 5, 4, 9, 3

$$\text{Soln:- } \bar{x} = \frac{\sum x}{n}$$

$$\bar{x} = \frac{6+8+7+5+4+9+3}{7} = \frac{42}{7} = 6$$

$$\text{Thus mean } (\bar{x}) = 6$$

$$\checkmark \text{ or } \sigma^2 = \frac{\sum (x-\bar{x})^2}{n} = \frac{1}{7} \left\{ (6-6)^2 + (8-6)^2 + (7-6)^2 + (5-6)^2 + (4-6)^2 + (9-6)^2 + (3-6)^2 \right\}$$

$$\checkmark = \frac{28}{7} = 4$$

$$\therefore \sigma = \sqrt{4} = 2$$

$$\begin{aligned} \text{Alternate: } \sigma^2 &= \frac{\sum x^2}{n} - (\bar{x})^2 \\ &= \frac{6^2 + 8^2 + 7^2 + 5^2 + 4^2 + 9^2 + 3^2}{7} - (6)^2 \\ &= \frac{280}{7} - 36 \\ &= 4 \end{aligned}$$

$$\text{Thus } SD = \sigma = 2$$

② Let us find the mean & SD for the following grouped data

class	1-10	11-20	21-30	31-40	41-50	51-60
Frequency	3	16	26	31	16	8

class	f	x	fx	$(x-\bar{x})$	$(x-\bar{x})^2$	$\sum f(x-\bar{x})^2$
1-10	3	5.5	16.5		702.25	2106.75
11-20	16	15.5	248.0		272.25	4356.00
21-30	26	25.5	663.0		42.25	1098.50
31-40	31	35.5	1100.5		12.25	379.75
41-50	16	45.5	728.0		182.25	2916.00
51-60	8	55.5	444.0		552.25	4418.00
Totals	100		3200			15275

$$\bar{x} = \frac{\sum f x}{\sum f} = \frac{3200}{100} = 32$$

$$S^2 = \frac{\sum f (x - \bar{x})^2}{\sum f} = \frac{15275}{100} = 152.75$$

$$\therefore S = \sqrt{152.75} = 12.36$$

Curve fitting :-

* Fitting of a straight line : $y = ax + b$
 consider a set of n given values (x, y) for fitting the straight line $y = ax + b$ where a & b are parameters to be determined. The residual $R = y - (ax + b)$ is the difference between the observed and estimated values of y . By the method of least squares we find parameters a & b such that the sum of squares of the residuals is minimum.

$$a \sum x + n b = \sum y$$

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

problems:-

- ① fit a straight line $y = ax + b$ for the following data.

x	1	3	4	6	8	9	11	14
y	1	2	4	4	5	7	8	9

The normal eqn for fitting the straight line.

$$y = ax + b \text{ are } \sum y = a \sum x + nb \quad (n=8)$$
$$\sum xy = a \sum x^2 + b \sum x$$

x	y	xy	x^2
1	1	1	1
3	2	6	9
4	4	16	16
6	4	24	36
8	5	40	64
9	7	63	81
11	8	88	121
14	9	126	196

$\sum x = 56 \quad \sum y = 40 \quad \sum xy = 364 \quad \sum x^2 = 524$

The normal eqns become

$$56a + 8b = 40.$$

$$524a + 56b = 364$$

$$\therefore a = 0.63 \approx 0.64, b = 0.54 \approx 0.55$$

∴ Thus by substituting these values in
 $y = ax + b$ we obtain the eqn.

$$y = 0.64x + 0.55$$

- ② find the eqn of the best fitting straight line for the following data & hence estimate the value of the dependent variable corresponding to the value 30 of the independent variable.

x	5	10	15	20	25
y	16	19	23	26	30

$$\therefore y = ax + b$$

$$\sum y = a \sum x + nb$$

$$\sum xy = a \sum x^2 + b \sum x \quad (n=5)$$

(3)

x	y	xy	x^2
5	16	80	25
10	19	190	100
15	23	345	225
20	26	520	400
25	30	750	625

$$\sum x = 75 \quad \sum y = 114 \quad \sum xy = 1885 \quad \sum x^2 = 1375$$

$$\therefore 75a + 5b = 114$$

$$1375a + 75b = 1885$$

$$\therefore a = 0.7, \quad b = 12.3$$

$$y = ax + b$$

$$y = 0.7x + 12.3$$

when $x = 30$, we obtain $y = 0.7(30) + 12.3 = 33.3$

(3) A simply supported beam carries a concentrated load P at its mid point.

corresponding to various values of P the maximum deflection y is measured & is given in the following table

	100	120	140	160	180	200
P	0.45	0.55	0.60	0.70	0.80	0.85
y						

find a law of the form $y = a + bp$ & hence estimate y when P is 150.

Soln: The normal equations associated with $y = a + bp$ are as follows.

$$\sum y = na + b \sum P \quad (n=6)$$

$$\sum Py = a \sum P + b \sum P^2$$

P	y	Py	P^2
100	0.45	45	10000
120	0.55	66	14400
140	0.60	84	19600
160	0.70	112	25600
180	0.80	144	32400
200	0.85	170	40000

$$\sum P = 900 \quad \sum y = 3.95 \quad \sum Py = 621 \quad \sum P^2 = 142000$$

$\therefore \text{eqns} \Rightarrow$

$$6a + 900b = 3.95$$

$$900a + 142000b = 621$$

$$\therefore a = 0.0476, \quad b = 0.0041$$

Thus the required Law is $y = 0.0476 + 0.0041P$

Also when $P = 150$, $y = 0.6626 \approx 0.66$

④ Fit a straight line to the following data.

year	1961	1971	1981	1991	2001
production (in tons)	8	10	12	10	16

Also find the expected production in the year 2006.

Soln: let $x = x - 1981$. & the line of fit with be $y = a + bx$

$$\sum y = na + b \sum x \quad (n=5)$$

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

(4)

x	y	xy	x^2
-20	8	-160	400
-10	10	-100	100
0	12	0	0
10	10	100	100
20	16	320	400

$$\sum x = 0 \quad \sum y = 56 \quad \sum xy = 160 \quad \sum x^2 = 1000$$

The normal equations become,

$$5a = 56 \quad \text{&} \quad 1000b = 160$$

$$a = 11.2 \quad \text{&} \quad b = 0.16$$

Hence $y = a + bx$, with $x = t - 1981$ becomes

$$y = 11.2 + 0.16(t - 1981)$$

Thus $\underline{y = -305.76 + 0.16x}$ is the required line of fit.

Also when $x = 2006$, $y = -305.76 + 0.16(2006)$

$$y = 15.2$$

Expected production in the year 2006 is 15.2 ton.

(5) Find the eqn of the best fitting straight line for the following data.

i)	x	1	2	3	4	5
	y	14	13	9	5	2

ii)	x	0	1	2	3	4	5
	y	9	8	24	28	26	20

iii)	x	62	64	65	69	70	71	72
	y	65.7	66.8	67.2	69.3	69.8	70.5	70.9

iv)	x	1	2	3	4	5	6	7
	y	80	90	92	83	94	99	92

⑥	Year (x)	1911	1921	1931	1941	1951
	productivity (in thousand tons)	8	10	12	10	6

Soln:- Let $x = x - 1931$ & the line of fit will be $y = a + bx$

The normal eqns associated with $y = a + bx$ are as follows.

$$\sum y = na + b \sum x \quad (n=5)$$

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

x	y	xy	x^2
-20	8	-160	400
-10	10	-100	100
0	12	0	0
10	10	100	100
20	6	120	400

$$\sum x = 0 \quad \sum y = 46 \quad \sum xy = -40 \quad \sum x^2 = 1000$$

The normal eqns become,

$$46 = 5a \rightarrow ①$$

$$-40 = 1000b \rightarrow ②$$

$$\text{Eqn } ① \quad a = 9.2$$

$$\text{Eqn } ② \quad b = -0.04$$

Hence $y = a + bx$, with $x = x - 1931$

$$y = 9.2 + (-0.04)(x - 1931)$$

$$= 9.2 - 0.04x + 77.24$$

Thus, $y = 86.44 - 0.04x$ is the required line of fit.

* fitting of a second degree parabola

$$y = ax^2 + bx + c$$

consider a set of n given values (x, y) for fitting the curve $y = ax^2 + bx + c$. The residual $R = y - (ax^2 + bx + c)$ is the difference b/w the observed & estimated values of y . we have to find parameters a, b, c such that the sum of the squares of the residuals is the least.

$$\sum y = a \sum x^2 + b \sum x + nc$$

$$\sum xy = a \sum x^3 + b \sum x^2 + c \sum x$$

$$\sum x^2 y = a \sum x^4 + b \sum x^3 + c \sum x^2$$

- * ① fit a best fitting parabola $y = ax^2 + bx + c$ for the following data:

$x \quad 1 \quad 2 \quad 3 \quad 4 \quad 5$ and hence estimate
 $y \quad 10 \quad 12 \quad 13 \quad 16 \quad 19$ at $x=6$.

Soln:- The normal eqn associated with $y = ax^2 + bx + c \rightarrow *$ are as follows.

$$y = ax^2 + bx + c \rightarrow ①$$

$$\sum y = a \sum x^2 + b \sum x + nc \rightarrow ②$$

$$\sum xy = a \sum x^3 + b \sum x^2 + c \sum x \rightarrow ③$$

$$\sum x^2 y = a \sum x^4 + b \sum x^3 + c \sum x^2 \rightarrow ④$$

x	y	xy	x^2	$x^2 y$	x^3	x^4
1	10	10	1	10	1	
2	12	24	4	48	8	16
3	13	39	9			
4	16	64	16	256	64	256
5	19	95	25	475	125	625

$$\sum x = 15 \quad \sum y = 70 \quad \sum xy = 232 \quad \sum x^2 = 55 \quad \sum x^2 y = 906 \quad \sum x^3 = 225 \quad \sum x^4 = 979$$

$$\text{Eqn } ① \Rightarrow 70 = 55a + 15b + 5c$$

$$\text{Eqn } ② \Rightarrow 232 = 225a + 55b + 15c$$

$$\text{Eqn } ③ \Rightarrow 906 = 979a + 225b + 55c$$

$$\therefore a = 0.2857 \approx 0.29, b = 0.4857 \approx 0.49, c = 9.4$$

Thus the required second degree parabola is

$$y = 0.29x^2 + 0.49x + 9.4 \text{ also at } x=6$$

$$y = 22.78$$

- (2) * fit a parabola $y = ax^2 + bx + c$ for the data.
- | x | 0 | 1 | 2 | 3 | 4 |
|---|---|-----|-----|-----|-----|
| y | 1 | 1.8 | 1.3 | 2.5 | 2.3 |

Soln:- The normal eqn associated $y = ax^2 + bx + c$

$$\text{are } \sum y = a\sum x^0 + b\sum x^1 + c\sum x^2 \rightarrow ①$$

$$\sum xy = a\sum x^1 + b\sum x^2 + c\sum x^3 \rightarrow ②$$

$$\sum x^2 y = a\sum x^2 + b\sum x^3 + c\sum x^4 \rightarrow ③$$

x	y	xy	x^2	$x^2 y$	x^3	x^4
0	1	0	0	0	0	0
1	1.8	1.8	1	1.8	1	1
2	1.3	2.6	4	5.2	8	16
3	2.5	7.5	9	22.5	27	81
4	2.3	9.2	16	36.8	64	256

$$\sum x = 10 \quad \sum y = 8.9 \quad \sum xy = 21.1 \quad \sum x^2 = 30 \quad \sum x^2 y = 66.3 \quad \sum x^3 = 100 \quad \sum x^4 = 354$$

$$\text{Eqn } ① \Rightarrow 8.9 = 5a + 10b + 30c$$

$$21.1 = 10a + 30b + 100c$$

$$\text{Eqn } ② \Rightarrow 66.3 = 30a + 100b + 354c$$

$$\text{Eqn } ③ \Rightarrow a = 1.0771, \quad b = 0.4157, \quad c = -0.0214$$

Thus the parabola of fit is

$$y = 1.0771 + 0.4157x - 0.0214x^2$$

- (3) * fit a second degree parabola to the following data,

x	1.0	1.5	2.0	2.5	3.0	3.5	4.0
y	1.1	1.3	1.6	2.0	2.7	3.4	4.1

- (4) * fit a second degree parabola to the following data:

x	0	1	2	3	4
y	1	1.8	1.3	2.5	6.3

* fitting of a curve of the form $y = ae^{bx}$.

consider, $y = ae^{bx}$. Taking logarithm (to the base e) on both sides we get.

$$\begin{aligned}\log_e y &= \log_e (ae^{bx}) \\ &= \log_e a + \log_e e^{bx} \\ &= \log_e a + bx \log_e e\end{aligned}$$

$$\left| \begin{array}{l} \log mn = \log m + \log n \\ \log m^n = n \log m \\ \log_e e = 1 \end{array} \right.$$

$$\log_e y = \log_e a + bx \quad (1)$$

$$y = A + BX \rightarrow (1)$$

where $y = \log_e y$, $A = \log_e a$, $B = b$, $X = x$

It is evident that eqn (1) is the eqn of a straight line & the associated normal eqns are as follows.

$$\sum y = nA + B \sum X \rightarrow (2)$$

$$\sum xy = A \sum X + B \sum X^2 \rightarrow (3)$$

Solving eqn (2) & (3) we obtain 'A' & 'B'. But

$\log_e a = A \Rightarrow a = e^A$. Also $b = B$. Substituting these values in $y = ae^{bx}$ we get the curve of best fit, in the required form.

* fit a curve of the form $y = ae^{bx}$ to the following data.

x	7.7	100	185	239	285
y	2.4	3.4	7.0	11.1	19.6

Soln: consider, $y = ae^{bx} \rightarrow *$

$$y = A + bx$$

The normal eqns are as follows.

$$\sum y = nA + b \sum X \rightarrow (1)$$

$$\sum xy = A \sum X + B \sum X^2 \rightarrow (2)$$

where $y = \log_e y$, $A = \log_e a$

x	y	$y = \log_e y$	xy	x^2
77	2.4	0.8754	67.4058	5929
100	3.4	1.2237	122.37	10000
185	7.0	1.9459	359.9915	34225
239	11.1	2.4069	575.2491	57121
285	19.6	2.9755	848.0175	81225
$\sum x = 886$		$\sum y = 43.5$	$\sum y = 9.4274$	$\sum xy = 1973.0339$
				$\sum x^2 = 188500$

$$\text{Eqn ①} \Rightarrow 9.4274 = 5A + 886b$$

$$\text{Eqn ②} \Rightarrow 1973.0339 = 886A + 188500b$$

$$A = \log_e a = 0.1838 \quad b = 9.6028 \times 10^{-3}$$

$$a = e^A = e^{0.1838}$$

$$a = e^{0.1838} = 1.2017 \quad b = 9.6028 \times 10^{-3}$$

The curve of fit $y = a e^{bx}$ is the curve of fit.

$$\text{Thus } y = (1.2017)e^{9.6028 \times 10^{-3}x}$$

- ② Fit an exponential curve of the form $y = a e^{bx}$ by the method of least squares for the following data.

No of petals	5	6	7	8	9	10
No of flowers	133	55	23	7	2	2

(7)

x	y	$y = \log_e y$	xy	x^2
5	133	4.8903	24.4515	25
6	55	4.0073	24.0438	36
7	23	3.1355	21.9485	49
8	7	1.9459	15.5672	64
9	2	0.6931	6.2379	81
10	2	0.6931	6.9310	100

$$\sum x = 45 \quad \sum y = 15.3652 \quad \sum xy = 99.1799 \quad \sum x^2 = 355$$

The normal eqn becomes

$$6A + 45b = 15.3652$$

$$45A + 355b = 99.1799$$

$$A = 9.4433 \quad b = -0.9177$$

Thus the required curve of fit is,

$$y = (126.23.3) e^{-0.9177x}$$

- ③ Fit a curve of the form $y = ax^b$ for the data
- | x | 1 | 2 | 3 | 4 | 5 | 6 |
|---|------|------|------|-----|-----|-----|
| y | 2.98 | 4.26 | 5.21 | 6.1 | 6.8 | 7.5 |

- ④ Find the eqn of the best fitting curve in the form $y = ae^{bx}$ for the data

x	0	2	4
y	5.02	10	31.62

correlation and correlation co-efficient:-

co-variation of two independent magnitudes is known as correlation. If two variables x & y are related in such a way that $\uparrow \text{ or } \downarrow$ in one of them corresponds to $\uparrow \text{ or } \downarrow$ in the other, we say that the variables are truly correlated. Also if increase or decrease in one of them corresponds to decrease or \uparrow in the other, the variables are said to be very correlated.

The numerical measure of correlation b/w two variables x & y is known as Pearson's coefficient of correlation usually denoted by r is defined as follows.

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{n \sqrt{x} \sqrt{y}} \rightarrow ①$$

This can be put in an alternative form as follows. If $x = x - \bar{x}$, $y = y - \bar{y}$ we can write.

$$\sqrt{x}^2 = \frac{\sum (x - \bar{x})^2}{n} = \frac{\sum x^2}{n},$$

$$\sqrt{y}^2 = \frac{\sum (y - \bar{y})^2}{n} = \frac{\sum y^2}{n}$$

$$\therefore \sqrt{x} \sqrt{y} = \frac{\sqrt{\sum x^2} \sqrt{\sum y^2}}{\sqrt{n^2}} = \frac{\sqrt{\sum x^2} \sqrt{\sum y^2}}{n}$$

$$\therefore n \sqrt{x} \cdot \sqrt{y} = \sqrt{\sum x^2} \cdot \sqrt{\sum y^2}$$

Thus eqn ① becomes

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

Note :- * The coefficient of correlation numerically does not exceed unity. i.e $-1 \leq r \leq +1$

* If $\gamma = \pm 1$ we say that x & y are perfectly correlated & if $\gamma = 0$ we say that x & y are non correlated. (8)

Alternative formula for the correlation co-efficient γ :

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

Proof: Let $z = x-y$

$$\therefore \frac{\sum z}{n} = \frac{\sum x}{n} - \frac{\sum y}{n} \quad (6) \quad \bar{z} = \bar{x} - \bar{y}$$

$$\text{Hence, } (z - \bar{z}) = (x - \bar{x}) - (\bar{x} - \bar{y})$$

$$\text{i.e. } (z - \bar{z}) = (x - \bar{x}) - (y - \bar{y})$$

Squaring both sides, taking summation &
dividing by n we have,

$$\begin{aligned} \frac{\sum (z - \bar{z})^2}{n} &= \frac{\sum [(x - \bar{x}) - (y - \bar{y})]^2}{n} \\ &= \frac{\sum (x - \bar{x})^2}{n} + \frac{\sum (y - \bar{y})^2}{n} - 2 \frac{\sum (x - \bar{x})(y - \bar{y})}{n} \end{aligned}$$

$$\text{i.e. } \sigma_z^2 = \sigma_x^2 + \sigma_y^2 - 2\gamma \sigma_x \sigma_y$$

$$\text{i.e. } \sigma_{x-y}^2 = \sigma_x^2 + \sigma_y^2 - 2\gamma \sigma_x \sigma_y$$

$$\text{Thus } \gamma = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x \sigma_y}$$

Note: In general if $z = ax+by$ we can obtain as before

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

$$\text{i.e. } \sigma_{ax+by}^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab \sigma_x \sigma_y$$

Regression :-

Regression is an estimation of one independent variable in terms of the other. If x & y are correlated, the best fitting straight line in the least square sense give reasonably a good relation b/w x & y .

The best fitting straight line of the form $y = ax + b$ (x being the independent variable) is called the regression line of y on x & $x = ay + b$ (y being the independent variable) is called the regression line of x on y . $y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$, $x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$

Note :- 1) The lines of regression (3) a.

$$y = \frac{\sum xy}{\sum x^2} (x) \text{ and } x = \frac{\sum xy}{\sum y^2} (y)$$

where $x = x - \bar{x}$ & $y = y - \bar{y}$.

This form will be useful to find out the coefficient of correlation by first obtaining the lines of regression as we have deduced that

$$r = \pm \sqrt{(\text{coeff. of } x)(\text{coeff. of } y)}$$

2) To compute the coefficient of correlation we prefer to use the formula.

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

where SDs can be found by applying the formula.

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

If \bar{x} & \bar{y} are integers computation of r by the formula.

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} \text{ is convenient where } x = x - \bar{x}, y = y - \bar{y}.$$

problems:-

- ① compute the coefficient of correlation and the equation of the lines of regression for the data.

$x \quad 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \quad 7$

$y \quad 9 \quad 8 \quad 10 \quad 12 \quad 11 \quad 13 \quad 14$

Soln:- we have $\gamma = \frac{5x^2 + 5y^2 - 6\bar{x}\bar{y}}{2\sqrt{5x} \sqrt{5y}} \rightarrow ①$

x	y	$z = x-y$	x^2	y^2	z^2
1	9	-8	1	81	64
2	8	-6	4	64	36
3	10	-7	9	100	49
4	12	-8	16	144	64
5	11	-6	25	121	36
6	13	-7	36	169	49
7	14	-7	49	196	49
$\sum x = 28$		$\sum y = 77$	$\sum z = -49$	$\sum x^2 = 140$	$\sum y^2 = 875$
					$\sum z^2 = 347$

$$\bar{x} = \frac{\sum x}{n}, \quad \bar{y} = \frac{\sum y}{n}, \quad \bar{z} = \frac{\sum z}{n}$$

$$\bar{x} = \frac{28}{7} = 4, \quad \bar{y} = \frac{77}{7} = 11, \quad \bar{z} = \frac{-49}{7} = -7$$

$$5x^2 = \frac{\sum x^2}{n} - (\bar{x})^2, \quad 5y^2 = \frac{\sum y^2}{n} - (\bar{y})^2, \quad 5z^2 = \frac{\sum z^2}{n} - (\bar{z})^2$$

$$= \frac{140}{7} - (4)^2 = 4 \quad = \frac{875}{7} - (11)^2 = 4 \quad = \frac{347}{7} - (-7)^2 = 0.5:$$

we have, $5x^2 = 4, 5y^2 = 4, 5z^2 = 5x^2 - 5y^2 = 0.57$

$$\text{Eqn } ① \Rightarrow \gamma = \frac{4+4-0.57}{2\sqrt{4}\sqrt{4}} = 0.92875 \approx 0.93,$$

Thus $\gamma = 0.93$

The lines of regression are given by

$$y - \bar{y} = \gamma \frac{5y}{5x} (x - \bar{x}) \quad \& \quad x - \bar{x} = \gamma \frac{5x}{5y} (y - \bar{y})$$

$$y - 11 = \frac{(0.93)12}{2} (x - 4), \quad x - 4 = \frac{(0.93)2}{2} (y - 11)$$

$$y - 11 = 0.93(x - 4), \quad x - 4 = 0.93(y - 11)$$

Thus lines of regression, $y = 0.93x + 7.28$ & $x = 0.93y - 6.23$ are the

(2) Obtain the lines of regression & hence find the coefficient of correlation for the data.

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

Soln:- Here $\bar{x} = 4$, $\bar{y} = 11$

$$\therefore x = x - \bar{x}, y = y - \bar{y}$$

$$= x - 4, y = y - 11$$

x	y	x	y	xy	x^2	y^2
1	9	-3	-2	6	9	4
2	8	-2	-3	6	4	9
3	10	-1	-1	1	1	1
4	12	0	1	0	0	1
5	11	1	0	0	1	0
6	13	2	2	4	4	4
7	14	3	3	9	9	9

$$\sum xy = 26, \sum x^2 = 28, \sum y^2 = 28$$

We shall consider regression lines in the form.

$$y = \frac{\sum xy}{\sum x^2} \cdot x \quad \text{and} \quad x = \frac{\sum xy}{\sum y^2} \cdot y$$

$$\text{i.e. } y - 11 = \frac{26}{28} (x - 4), \quad x - 4 = \frac{26}{28} (y - 11)$$

$$y - 11 = 0.93(x - 4), \quad x - 4 = 0.93(y - 11)$$

$$y = 0.93x + 7.28, \quad x = 0.93y - 6.23$$

These are the regression lines and we compute ' γ ' as the geometric mean of the regression coefficients.

$$\text{i.e. } \gamma = \sqrt{(\text{coeff of } x)(\text{coeff of } y)} = \sqrt{(0.93)(0.93)}$$

$$\gamma = 0.93$$

The sign of ' γ ' must be +ve since both the regression coefficients are +ve.

$$\text{Thus } \gamma = 0.93.$$

③ Find the correlation coefficient & the eqn of the liner of regression for the following values of x & y .

$x = 1, 2, 3, 4, 5$

$y = 2, 5, 3, 8, 7$

Soln: $n = 5$

x	y	$z = x - y$	x^2	y^2	z^2
1	2	-1	1	4	1
2	5	-3	4	25	9
3	3	0	9	9	0
4	8	-4	16	64	16
5	7	-2	25	49	4
$\sum x = 15$		$\sum y = 25$	$\sum z = -10$	$\sum x^2 = 55$	$\sum y^2 = 151$
					$\sum z^2 = 30$

$$\bar{x} = \frac{\sum x}{n}, \bar{y} = \frac{\sum y}{n}$$

$$\bar{z} = \frac{\sum z}{n}$$

$$\bar{x} = \frac{15}{5} = 3, \bar{y} = \frac{25}{5} = 5,$$

$$\bar{z} = \frac{-10}{5} = -2$$

$$\sigma_x^2 = \frac{\sum x^2}{n} - (\bar{x})^2 = \frac{55}{5} - (3)^2 = 2$$

$$\sigma_y^2 = \frac{\sum y^2}{n} - (\bar{y})^2 = \frac{151}{5} - (5)^2 = 5.2$$

$$\sigma_z^2 = \frac{\sum z^2}{n} - (\bar{z})^2 = \frac{30}{5} - (-2)^2 = 2$$

we have, $\sigma_x^2 = 2, \sigma_y^2 = 5.2, \sigma_z^2 = 2$

$$\text{Now, } \gamma = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_z^2}{2\sigma_x \sigma_y}$$

$$\gamma = \frac{2 + 5.2 - 2}{2\sqrt{2} \sqrt{5.2}} = 0.8062 \approx 0.81$$

Thus $\boxed{\gamma = 0.81}$

The eqn of the regression liner are as follows.

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

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

$$y - 5 = (0.81) \cdot \frac{\sqrt{5.2}}{\sqrt{2}} (x - 3),$$

$$x - 3 = 0.81 \cdot \frac{\sqrt{2}}{\sqrt{5.2}} (y - 5)$$

$$y = 5 + 1.306(x-3), \quad x-3 = 0.502(y-5)$$

Thus $y = 1.306x + 1.082$ & $x = 0.502y + 0.49$
These are the lines of regression.

- ④ Find the correlation coefficient b/w x & y for the following data. Also obtain the regression lines.

x	1	2	3	4	5	6	7	8	9	10
y	10	12	16	28	25	36	41	49	40	50

Soln:- Here $n = 10$

x	y	$z = x-y$	x^2	y^2	z^2
1	10	-9	1	100	81
2	12	-10	4	144	100
3	16	-13	9	256	169
4	28	-24	16	784	576
5	25	-20	25	625	400
6	36	-30	36	1296	900
7	41	-34	49	1681	1156
8	49	-41	64	2401	1681
9	40	-31	81	1600	961
10	50	-40	100	2500	1600
$\sum x = 55$		$\sum y = 307$	$\sum z = -252$	$\sum x^2 = 385$	$\sum y^2 = 11387$
					$\sum z^2 = 7624$

$$\bar{x} = \frac{\sum x}{n} = \frac{55}{10} = 5.5$$

$$\bar{y} = \frac{\sum y}{n} = \frac{307}{10} = 30.7$$

$$\bar{z} = \frac{\sum z}{n} = \frac{-252}{10} = -25.2$$

$$\sigma_x^2 = \frac{\sum x^2}{n} - (\bar{x})^2 = \frac{385}{10} - (5.5)^2 = 8.25, \sigma_x = 2.87$$

$$\sigma_y^2 = \frac{\sum y^2}{n} - (\bar{y})^2 = 1138.7 - (30.7)^2 = 196.21, \sigma_y = 14.01$$

$$\sigma_{xy}^2 = \sigma_{x-y}^2 = \frac{\sum xy}{n} - (\bar{x})(\bar{y}) = 762.4 - (5.5)(30.7) = 127.36$$

we have, $\gamma = \frac{\sigma_{xy}^2 + \sigma_y^2 - \sigma_x^2}{2\sigma_x\sigma_y}$

$$= \frac{8.25 + 196.21 - 127.36}{2 \times 2.87 \times 14.01} = 0.96$$

Thus $\boxed{\gamma = 0.96}$

Equation of the lines of regression are

$$y - \bar{y} = \frac{\gamma \sigma_y}{\sigma_x} (x - \bar{x}), \quad x - \bar{x} = \gamma \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

On substituting & simplifying we get,

$$y - 30.7 = 0.96 \times \frac{14.01}{2.87} (x - 5.5)$$

$$\boxed{y = 4.686x + 4.927}$$

$$x - 5.5 = 0.96 \times \frac{2.87}{14.01} (y - 30.7)$$

$$\boxed{x = 0.197y - 0.548}$$

These are the lines of regression.

- * ⑤ find the coefficient of correlation for the following data.

x	10	14	18	22	26	30
y	18	12	24	6	30	36

Soln:- we have $\bar{x} = \frac{\sum x}{n} = \frac{120}{6} = 20$

$$\bar{y} = \frac{\sum y}{n} = \frac{126}{6} = 21$$

Let $x = x - \bar{x}$, & $y = y - \bar{y}$

$$x = x - 20, \quad y = y - 21 \quad \&$$

x	y	x	y	x^2	y^2	xy
10	18	-10	-3	100	9	30
14	12	-6	-9	36	81	54
18	24	-2	3	4	9	-6
22	6	2	-15	4	225	-30
26	30	6	9	36	81	54
30	36	10	15	100	225	150
				$\sum x^2 = 280$	$\sum y^2 = 630$	$\sum xy = 252$

$$\rho = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{252}{\sqrt{280} \times \sqrt{630}} = 0.6$$

- * ⑥ Find the correlation coefficient and the eqⁿ of the line of regression for the following.

x	1	2	3	4	5
y	2	5	3	8	7

Soln:- Let $n=5$

x	y	$z = x - y$	x^2	y^2	z^2
1	2	-1	1	4	1
2	5	-3	4	25	9
3	3	0	9	9	0
4	8	-4	16	64	16
5	7	-2	25	49	4
$\sum x = 15$	$\sum y = 25$	$\sum z = -10$	$\sum x^2 = 55$	$\sum y^2 = 151$	$\sum z^2 = 30$

$$\bar{x} = \frac{\sum x}{n} = \frac{15}{5} = 3, \bar{y} = \frac{\sum y}{n} = \frac{25}{5} = 5, \bar{z} = \frac{\sum z}{n} = \frac{-10}{5} = -2$$

$$S_x^2 = \frac{\sum x^2}{n} - (\bar{x})^2 = \frac{55}{5} - (3)^2 = 2$$

$$S_y^2 = \frac{\sum y^2}{n} - (\bar{y})^2 = \frac{151}{5} - (5)^2 = 5.2$$

$$S_z^2 = S_x^2 - S_y^2 = \frac{\sum z^2}{n} - (\bar{z})^2 = \frac{30}{5} - (-2)^2 = 2$$

$$Now = \frac{S_x^2 + S_y^2 - S_z^2}{2 S_x S_y}$$

(12)

$$\gamma = \frac{2+5.2-2}{2\sqrt{2}\sqrt{5.2}} = 0.8062 \approx 0.81$$

$$\text{Thus } \gamma = 0.81$$

The eqns of the regression lines are as follows.

$$y - \bar{y} = \gamma \frac{\sigma_y}{\sigma_x} (x - \bar{x}), \quad x - \bar{x} = \gamma \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$y - 5 = (0.81) \frac{\sqrt{5.2}}{\sqrt{2}} (x - 3), \quad x - 3 = (0.81) \frac{\sqrt{2}}{\sqrt{5.2}} (y - 5)$$

$$y - 5 = 1.306 (x - 3), \quad x - 3 = 0.502 (y - 5)$$

$$\text{Thus } y = 1.306x + 1.182 \text{ & } x = 0.502y + 0.49$$

These are the lines of regression.

- * 7) Find the correlation coefficient b/w x & y for the following data. Also obtain the regression lines.

x	1	2	3	4	5	6	7	8	9	10
y	10	12	16	28	25	36	41	49	40	50

<u>Soln:-</u>		<u>Here</u>	<u>$n=10$</u>	<u>x^2</u>	<u>y^2</u>	<u>z^2</u>
x	y	$z = x - y$				
1	10	-9		1	100	81
2	12	-10		4	144	100
3	16	-13		9	256	169
4	28	-24		16	784	576
5	25	-20		25	625	900
6	36	-30		36	1296	1156
7	41	-34		49	1681	1681
8	49	-41		64	2401	961
9	40	-31		81	1600	1600
10	50	-40		100	8500	
					$\sum y^2 = 11387$	$\sum z^2 = 7624$
				$\sum x^2 = 385$		
				$\sum z = -252$		

$$\sum x = 55 \quad \sum y = 307$$

$$\hat{x} = \frac{\sum x}{n} = \frac{55}{10} = 5.5, \quad \hat{y} = \frac{\sum y}{n} = \frac{307}{10} = 30.7, \quad \hat{z} = \frac{\sum z}{n} = \frac{-252}{10} = -25.2$$

$$\sum x^2 = \frac{\sum x^2}{n} - (\bar{x})^2 = \frac{385}{10} - (5.5)^2 = 8.25, \bar{x} = 2.87$$

$$\sum y^2 = \frac{\sum y^2}{n} - (\bar{y})^2 = \frac{11387}{10} - (30.7)^2 = 196.21, \bar{y} = 14.01$$

$$\sum z^2 = \sum x^2 + \sum y^2 - (\sum x^2 + \sum y^2) = \frac{7624}{10} - (-25.2)^2 = 127.36$$

We have $\gamma = \frac{\sum x^2 + \sum y^2 - \sum z^2}{2 \sum x \sum y}$

$$\gamma = \frac{8.25 + 196.21 - 127.36}{2 \times 2.87 \times 14.01} = 0.96$$

Thus $\gamma = 0.96$

Eqs of the line of regression are

$$y - \bar{y} = \gamma \frac{\sum y}{\sum x} (x - \bar{x}), x - \bar{x} = \gamma \frac{\sum x}{\sum y} (y - \bar{y})$$

$$y - 30.7 = 0.96 \times \frac{14.01}{2.87} (x - 5.5), (x - 5.5) = 0.96 \times \frac{2.87}{14.01} (y - 30.7)$$

$$y = 4.686x + 4.927 \text{ and } x = 0.197y - 0.548$$

There are lines of regression.

Q) Find the regression line of 'y' on 'x' for the following data.

x	1	3	4	6	8	9	11	14
y	1	2	4	4	5	7	8	9

$$\text{Soln: } \bar{x} = \frac{\sum x}{n} = \frac{56}{8} = 7, \bar{y} = \frac{\sum y}{n} = \frac{40}{8} = 5$$

We denote $x = x - \bar{x}$ and $y = y - \bar{y}$
 $x = x - 7, y = y - 5$

We have lines of regression in the form

$$y = \frac{\sum xy}{\sum x^2} x \quad x = \frac{\sum xy}{\sum y^2} \cdot y.$$

(13)

x	y	$x = x - \bar{x}$	$y = y - \bar{y}$	x^2	y^2	xy
1	1	-6	-4	36	16	24
3	2	-4	-3	16	9	12
4	4	-2	-1	4	1	2
6	4	-1	-1	1	1	1
8	5	1	0	1	0	0
9	7	2	2	4	4	4
11	8	4	3	16	9	12
14	9	7	4	49	16	28

$$\sum x = 56 \quad \sum y = 40$$

$$\sum x^2 = 132 \quad \sum y^2 = 56 \quad \sum xy = 84$$

i.e $y - \bar{y} = \frac{\sum xy}{\sum x^2} (x - \bar{x})$, $x - \bar{x} = \frac{\sum xy}{\sum y^2} (y - \bar{y})$

$$y - 5 = \frac{84}{132} (x - 7) \quad x - 7 = \frac{84}{56} (y - 5)$$

$$y - 5 = 0.63 (x - 7) \quad x - 7 = 1.4285 (y - 5)$$

$$y - 5 = 0.63x - 4.47 \quad x - 7 = 1.4285y - 7.1425$$

$$y = 0.63x - 4.47 + 5 \quad x = 1.4285y - 7.1425 + 7$$

$$y = 0.63x + 0.58 \quad x = 1.4285y - 0.1425$$

These are the lines of regression.

- * (9) calculate the karl pearson's coefficient of correlation for 10 students who have obtained the following % of marks in mathematics & Electronics.

Roll No	1	2	3	4	5	6	7	8	9	10
Mark in mathematics	78	36	98	25	75	82	90	62	65	39
Mark in Electronics	84	51	91	60	68	62	86	58	53	47

Soln:- we have $\bar{x} = \frac{\sum x}{n} = \frac{650}{10} = 65$

$$\bar{y} = \frac{\sum y}{n} = \frac{660}{10} = 66$$

$$\text{Let } x = x - \bar{x} \quad \text{&} \quad y = y - \bar{y}$$

$$x = x - 65 \quad y = y - 66$$

x	y	$x = x - \bar{x}$	$y = y - \bar{y}$	x^2	y^2	xy
78	84	13	18	169	324	234
36	51	-29	-15	841	225	435
98	91	33	25	1089	625	825
85	60	-10	-6	1600	36	240
75	68	10	2	100	4	20
82	62	17	-4	289	16	68
90	86	25	20	625	400	500
62	58	-3	-8	9	64	94
65	53	0	-13	0	169	0
39	47	-26	-19	676	361	494

$$\sum x^2 = 5398 \quad \sum y^2 = 2224 \quad \sum xy = 2840$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{2840}{\sqrt{5398} \sqrt{2224}} = \frac{2840}{\sqrt{5398} \sqrt{2224}} = \frac{2840}{3464.1105}$$

$$r = 0.819 \approx 0.8$$

- ★ ⑩ In a partially destroyed lab record, only the lines of regression of y on x & x on y are available as $4x - 5y + 33 = 0$ & $20x - 9y = 107$ respectively. calculate \bar{x} , \bar{y} and coefficient of correlation b/w x & y .

Soln: we know that regression lines passes through \bar{x} & \bar{y} .

$$4\bar{x} - 5\bar{y} = -33$$

$$20\bar{x} - 9\bar{y} = 107$$

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

we shall now rewrite the eqn of the regression line to find the regression coefficients.

$$5y = 4x + 33 \quad \text{or} \quad y = 0.8x + 6.6 \rightarrow ①$$

$$20x = 9y + 107 \quad \text{or} \quad x = 0.45y + 5.35 \rightarrow ②$$

From ① & ②

$$\gamma \cdot \frac{\bar{y}}{\bar{x}} = 0.8 , \quad \gamma \cdot \frac{\bar{x}}{\bar{y}} = 0.45$$

correlation coefficient $\gamma = \sqrt{0.8 \times 0.45} = \pm 0.6$

Thus $\gamma = 0.6$

* Show that if θ is the angle b/w the lines of regression, then

$$\tan \theta = \frac{\bar{x} \bar{y}}{\bar{x}^2 + \bar{y}^2} \left(\frac{1 - \gamma^2}{\gamma} \right)$$

Sol'n: If θ is the angle b/w the lines $y = m_1 x + c_1$ and $y = m_2 x + c_2$ is given by.

$$\tan \theta = \frac{m_2 - m_1}{1 + m_1 m_2}$$

We have the lines of regression,

$$y - \bar{y} = \gamma \frac{\bar{y}}{\bar{x}} (x - \bar{x}) \quad ① \quad x - \bar{x} = \gamma \frac{\bar{x}}{\bar{y}} (y - \bar{y}).$$

We write the second of the eqn as

$$y - \bar{y} = \frac{\bar{y}}{\gamma \bar{x}} (x - \bar{x}) \rightarrow ②$$

Slopes of ① & ② are respectively given by

$$m_1 = \gamma \frac{\bar{y}}{\bar{x}} \quad \text{and} \quad m_2 = \frac{\bar{y}}{\gamma \bar{x}}$$

Substituting these in the formula for $\tan \theta$
we have,

$$\tan \theta = \frac{\frac{\bar{y}}{\gamma \bar{x}} - \frac{\bar{y}}{\bar{x}}}{1 + \frac{\bar{y}}{\gamma \bar{x}} \cdot \frac{\bar{y}}{\bar{x}}} = \frac{\frac{\bar{y}}{\bar{x}} \left(\frac{1}{\gamma} - 1 \right)}{1 + \frac{\bar{y}^2}{\bar{x}^2}}$$

$$\text{Thus } \tan \theta = \frac{\frac{\bar{y}}{\bar{x}} \left(\frac{1 - \gamma^2}{\gamma} \right)}{\frac{\bar{x}^2 + \bar{y}^2}{\bar{x}^2}} = \frac{\bar{x} \bar{y}}{\bar{x}^2 + \bar{y}^2} \left(\frac{1 - \gamma^2}{\gamma} \right)$$

② Find the rank correlation for the data

x	8	5	11	13	10	5	18	15	2	8
y	56	44	79	72	70	54	94	85	33	65

Solu:-

x	y	Rank (x)	Rank (y)	Adjust rank Rx	Adjust rank Ry	$d = R_x - R_y$
8	56	6	7	$6+7=6.5$	7	-0.5
15	44	8	9	$8+9=8.5$	9	-0.5
11	79	4	3	4	3	1
13	72	3	4	3	4	-1
10	70	5	5	5	5	0
5	54	8	8	8.5	8	0.5
18	94	1	1	1	1	0
15	85	2	2	2	2	0
2	33	10	10	10	10	0
8	65	6	6	6.5	6	0.5

d^2	0.25	0.25	1	1	0	0.25	0	0	0	0.25
-------	------	------	---	---	---	------	---	---	---	------

$$\sum d^2 = 3$$

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

$$n = 10$$

$$i = 2$$

$$C.F = \frac{1}{12} \sum_{i=1}^K m_i (m_i^2 - 1) = \frac{1}{12} \left(\sum_{i=1}^2 m_i (m_i^2 - 1) \right)$$

$$m_1 = 2$$

$$m_2 = 2$$

$$= \frac{1}{12} \left[m_1 (m_1^2 - 1) + m_2 (m_2^2 - 1) \right]$$

$$= \frac{1}{12} \left[2(2^2 - 1) + 2(2^2 - 1) \right] = 1$$

$$\therefore \rho = 1 - \frac{6[3+1]}{10(10^2-1)} = 1 - \frac{24}{990} = 0.976$$

③ Find the rank correlation coefficient for the given data.

x	50	33	40	10	15	15	65	24	15	57
y	12	12	24	6	15	4	20	9	6	18

Solu:

x	y	Rank x	Rank y	Adjust Rank Rx	Adjust Rank Ry	d = Rx - Ry	d^2
50	12	3	5	3	$\frac{5+6}{2} = 5.5$	-2.5	6.25
33	12	5	5	5	5.5	-0.5	0.25
40	24	4	1.5	4.5	8	3	9
10	6	10	8	10	$\frac{8+9}{2} = 8.5$	1.5	2.25
15	15	7	4	$\frac{7+8+9}{3} = 8$	4	4	16
15	4	8	10	8	10	-2	4
65	20	1	2	1	2	-1	1
24	9	6	7	6	7	-1	1
15	6	7	8	8	8.5	-0.5	0.25
57	18	2	3	2	3	-1	1

$$i=3, m_1=3, m_2=2, m_3=2$$

$$n=10$$

$$\rho = 1 - \frac{6(\sum d^2 + C.F)}{n(n^2-1)}$$

$$C.F = \frac{1}{12} \sum_{i=1}^3 m_i(m_i^2 - 1) = \frac{1}{12} \left[m_1(m_1^2 - 1) + m_2(m_2^2 - 1) + m_3(m_3^2 - 1) \right]$$

$$\begin{aligned}
 C.F &= \frac{1}{12} [3(3^2 - 1) + 2(2^2 - 1) + 2(2^2 - 1)] \\
 &= \frac{1}{12} [24 + 6 + 6] = \frac{36}{12} = 3 \\
 \therefore \rho &= 1 - \frac{6[41+3]}{10(10^2 - 1)} = 1 - \frac{264}{990} \\
 &\boxed{\rho = 0.7333}
 \end{aligned}$$

③ Find the rank correlation for the data

x	68	64	75	50	64	80	75	40	55	64
y	74	58	68	45	61	60	68	48	50	74