Module IV- CURVE FITTING AND STATISTICAL METHODS

CORRELATION

CORRELATION: Covariation of two independent magnitudes is known as correlation

Ex: 1. Field of crop varies with the amount of rainfall.

2. Price of the commodity increases with the reduction of its supply.

If an increase (or decrease) in the values of one variable corresponds to an increase (or decrease) in the other ,the correlation is said to be positive .

If the increase (or decrease) in one variable corresponds to the decrease (or increase) in the other, the correlation is said to be negative.

If there is no relationship indicated between the variables, they are said to be independent or uncorrelated.

KARL-PEARSON'S COEFFICIENT OF CORRELATION:

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

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

Also r is given by

$$r = \frac{\sum XY}{n\sigma_x\sigma_y}.$$

$$r = \frac{\sigma^2 + \sigma^2 - \sigma^2}{2\sigma_x\sigma_y}.$$

$$r = \frac{n\sum f d_x d_y - (\sum f_x d_x)(\sum f_y d_y)}{\sqrt{[n\sum f_x d_x^2 - (\sum f_x d_x)^2] \times [n\sum f_y d_y^2 - (\sum f_y d_y)^2]}}$$

(step deviation method)

Note:

1. $X = x - \overline{x}$ and $Y = y - \overline{y}$

2.
$$\bar{x} = \frac{\Sigma S}{n} akd \bar{y} = \frac{\Sigma y}{n}$$

3. $\sigma_{S} = \sqrt{\frac{\Sigma(S-\bar{S})^{2}}{n}} = \frac{\Sigma S^{2}}{n} - \bar{x}^{2}$ and $\sigma_{y} = \sqrt{\frac{\Sigma(y-\bar{y})^{2}}{n}} = \frac{\Sigma y^{2}}{n} - \bar{y}^{2}$

EXAMPLES:

Example 1: Psychological tests of intelligence and of engineering ability were applied to 10 students. Here is a record of ungrouped data showing intelligence ratio (IR) and Engineering ratio (ER). Calculate the coefficient of correlation.

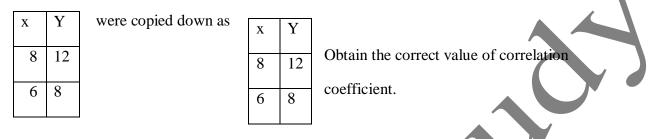
Students	A	В	C	D	E	F	G	Η	I
I.R.	105	104	102	101	100	99	98	96	93 92
E.R.	101	103	100	98	95	96	104	92	97 94
Solution:						(

Student	Intelligen	ce ratio	Enginee	ring ratio	X ²	Y ²	XY
	x	$X = x - \bar{x}$	У	$Y = y - \bar{y}$			
A	105	6	101	3	36	9	18
В	104	5	103	5	25	25	25
С	102	3	100	2	9	4	6
D	101	2	98	0	4	0	0
E	100	5	95	-3	1	9	-3
F	99	0	96	-2	0	4	0
G	98	-1	104	6	1	36	-6
H	96	-3	92	-6	9	36	18
Ι	93	-6	97	-1	36	1	6
J	92	-7	94	-4	49	16	28
Total	$\sum x = 990$	0	$\sum y = 980$	0	$\sum X^2 = 170$	$\sum Y^2 = 140$	$\sum XY = 92$

$$\bar{x} = \frac{\sum S}{n} = \frac{990}{10} = 99 \ akd \ \bar{y} = \frac{\sum y}{n} = \frac{980}{10} = 98$$

Using $r = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} = \frac{92}{\sqrt{170 \times 140}} = 0.59.$

Example 2: While calculating correlation coefficient between two variables x and y from 25 pairs of observations, the following results were obtained: k = 25; $\sum x = 125$; $\sum x^2 = 650$; $\sum y = 100$; $\sum y^2 = 460$; $\sum xy = 508$. Later it was discovered at the time of checking that the pairs of values



Solution: k = 25,

 $\sum x = 125 - 6 - 8 + 8 + 6 = 125$ (6 & 8 were wrongly taken so subtracted; instead 8 & 6 are the correct values so added)

 $\sum y = 100 - 14 - 6 + 12 + 8 = 100$ (14 & 6 were wrongly taken so subtracted; instead 12 & 8 are the correct values so added)

Similarly,

$$\sum x^{2} = 650 - 6^{2} - 8^{2} + 8^{2} + 6^{2} = 650$$

$$\sum y^{2} = 460 - 14^{2} - 6^{2} + 12^{2} + 8^{2} = 436$$

$$\sum xy = 508 - (6 \times 14) - (8 \times 6) + (8 \times 12) + (6 \times 8) = 12$$

$$r = \frac{n \sum Sy - (\sum S)(\sum y)}{\sqrt{\{n \sum S^{2} - (\sum S)^{2}\}\{n \sum y^{2} - (\sum y)^{2}\}}} = \frac{2}{3}$$

Example 3: Establish the formula $r = \frac{\frac{\sigma^2 + \sigma^2 - \sigma^2}{x - y}}{\frac{2\sigma_x \sigma_y}{2\sigma_x \sigma_y}}$. Hence calculate r from the following data

520

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

Solution: Let z = x - y so that $\overline{z} = \overline{x} - \overline{y}$

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

$$(z - \bar{z})^2 = (x - \bar{x})^2 + (y - \bar{y})^2 - 2(x - \bar{x})(y - \bar{y})$$

Summing up for n terms, and dividing by n

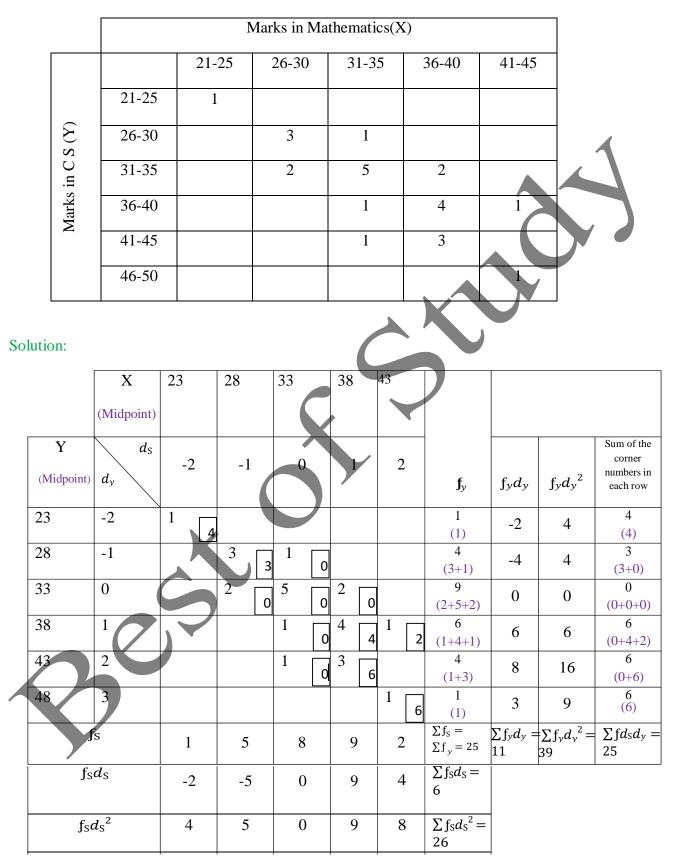
$$\frac{\Sigma(Z-Z)^2}{n} = \frac{\Sigma(S-S)^2}{n} + \frac{\Sigma(y-\bar{y}^2)}{n} - 2\frac{\Sigma(S-S)(y-\bar{y})}{n}$$
$$\sigma_Z^2 = \sigma_S^2 + \sigma_y^2 - 2\sigma_S \sigma_y r$$
$$\therefore r = \frac{\sigma_S^2 + \sigma_y^2 - \sigma_z^2}{2\sigma_x \sigma_y}.$$

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

x	X = x - 54	X ²	У	Y = y - 100	Y ²	z = x - y	<i>z</i> ²
21	-33	1089	60	-40	1600	-39	1521
23	-31	969	71	-29	841	-48	2304
30	-24	576	72	-28	784	-42	1764
54	0	0	83	-17	289	-29	841
57	3	9	110	10	100	-53	2809
58	4	16	84	-16	256	-26	676
72	18	324	100	0	0	-28	784
78	24	576	92	-8	64	-14	196
87	33	1089	113	13	169	-26	676
90	36	1296	135	35	1225	-45	2025
	$\sum X = 30$	$\sum X^2 = 5936$	5	$\sum Y = -80$	$\sum X^2 = 5328$	$3\Sigma X = -350$	$\sum z^2 = 13596$

$$\sigma_{s}^{2} = \sum_{n}^{\Sigma x^{2}} \sum_{n} \sum_{n}^{2} \sum_{n}^{2}$$

Example 4: The following table shows the bivariate frequency distribution of marks obtained by 25 students in maths(X) and Computer science(Y). Determine the co-efficient of correlation (r).



Sum of the corner numbers in each row	4	3	0	10	8	$\sum_{z \in S} f d_S d_y =$	
---------------------------------------	---	---	---	----	---	------------------------------	--

Here n = total frequency =25

Substituting the values from the above table in the formula for r i.e

$$r = \frac{n \sum \mathbf{f} d_x d_y - (\sum \mathbf{f}_x d_x) (\sum \mathbf{f}_y d_y)}{\sqrt{[n \sum \mathbf{f}_x d_x^2 - (\sum \mathbf{f}_x d_x)^2] \times [n \sum \mathbf{f}_y d_y^2 - (\sum \mathbf{f}_y d_y)^2]}} = 0.772$$

Note:

 $d_{\rm S} = \frac{{\rm S}-a}{h_1} \& d_y = \frac{y-b}{h_2}$, a & b are the arbitrary values taken as class marks(assumed values usually the middle or near to middle value in the midpoint), $h_1 akd h_2$ are the size of class intervals respectively.

Purple coloured are just for detailed information, not to be included in the table.

Try these:

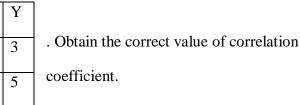
Example 1: Calculate the coefficient of correlation for the following ages of husbands(x) and wives(y)

Х	23	27	28	28	29	30	31	33	35	36			
Y	18	20	22	27	21	29	27	29	28	29			
	0.00												

Answer: 0.82

Example 2: From 10 pairs of observations for x and y the following data is obtained n=10, $\sum x = 66$, $\sum y = 69$, $\sum x^2 = 476$, $\sum y^2 = 521$, $\sum xy = 485$. It was later found that the two pairs of (correct) values





Answer: 0.55

Example 3: The correlation table below shows that the ages of husband and wife of 53 married couples living together om the census night of 1991. Calculate the coefficient of correlation for the following ages of husband and wife.

				Age o	of wife			Total
		15-25	25-35	35-45	45-55	55-65	65-75	
	15-25	1	1					2
pu	25-35	2	12	1				15
Age of husband	35-45		4	10	1			15
ge of l	45-55			3	6	1		10
Υβ	55-65				2	4	2	8
	65-75					1	2	3
Т	otal	3	17	14	9	6	4	53

Answer: 0.91



A group of n 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 order corresponding to two characteristics A & B, the correlation between these n pairs of ranks is called the **Rank correlation** in the characteristics A and B for the group of individuals.

Let x_i , y_i be the ranks of the ith individuals in A and B respectively, then the rank correlation coefficient between these variables is denoted by ρ and is given by

 $\rho = 1 - \sum_{n=1}^{\infty} d_i^2$, where $d_i = x - y_i$ EXAMPLES:

Example . 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 ρ .

Solution:

X:	1	6	5	10	3	2	4	9	7	8
y:	6	4	9	8	1	2	3	10	5	7
$d_{\rm i} = x_{\rm i} - y_{\rm i}$	-5	2	-4	2	2	0	1	-1	2	1
d ²	25	4	16	4	4	0	1	1	4	1

k = 10

 $\sum d_1^2 = 60$

$$\therefore \ \rho = 1 - \frac{6\sum d_{i}^{2}}{n^{3} - n} = 1 - \frac{6 \times 60}{990} = 0.6(kearly).$$

Try these:

Example 1:

The marks scored by recruits in the selection test(X) in the proficiency test(Y) are given below

Serial No			3	4 5 6 7 8	9
Х	10	15	12	17 13 16 24 14	22
Y	30	42	45	46 33 34 40 35	39

Answer: 0.4 (hint: take X as x and Y as y)

Example 2:

Three judges A.B.C give the following ranks. Find which pair of judges have common approach?

A	1	6	5	10	3	2	4	9	7	8
В	3	5	8	4	7	10	2	1	6	9
С	6	4	9	8	1	2	3	10	5	7

Answer: The pair of judges A and C have the nearest common approach(:: ρ_{Z-S} is maximum).

(hint: make table as below

A = x	B = y	C = z	$d_1 = x - y$	$d_2 = y - z$	$d_3=z-x$	d_1^2	d_2^2	d_3^2
								1

$$\rho_{S-y} = -0.2, \quad \rho_{y-z} = -0.3, \quad \rho_{Z-S} = 0.6.$$



REGRESSION ANALYSIS

Lines of regression

Introduction: 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 gives reasonably a good relation between x & y. The best fitting straight line of the form y = ax + b is called the regression line of $y \circ k x$ and x = ay + b is called the regression line of $x \circ k y$.

Equations of the regression lines

The equation (i) $y - y = r \frac{\sigma_y}{\sigma_x} (x - x)$ or $Y = \frac{\sum XY}{\sum X^2} (X)$ is called the regression line of y ok x

The equation (ii) $x - \overline{x} = r_{\sigma_y}^{\mathcal{A}} (y - \overline{y})$ or $X = \frac{\Sigma XY}{\Sigma Y^2}(Y)$ is called the regression line of x ok y

where $X = x - \overline{x}$ and $Y = y - \overline{y}$.

The coefficient of x in (i) and the coefficient of y in (ii) respectively given by $r \frac{\sigma_y}{\sigma_x}$ and $r \frac{\sigma_x}{\sigma_y}$ are known as the regression coefficients, their product is equal to r^2 and hence

$$r = \sqrt{(\text{coeffof } x)(c \text{ oeffof } y)}$$
 is the coefficient of correlation
OR $r = \frac{\sum XY}{\sqrt{\sum X^2}\sqrt{\sum Y^2}}$ OR $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{X-Y}^2}{2\sigma_x\sigma_y}$

where S.Ds can be found by applying the formula $\sigma^2 = \frac{\sum S^2}{n} - (\pi)^2$

EXAMPLES:

Example 1: Obtain the lines of regression and hence find the coefficient of correlation for the data

Solution: Let $x = \frac{\sum S}{n} = \frac{70}{10} = 7$ and $y = \frac{\sum y}{n} = \frac{150}{10} = 15$ we note that X = x - x and Y = y - y.

The relevant table is as follows

x	y	$X = x - \overline{x}$	$Y = y - \overline{y}$	X ²	Y^2	XY
1	8	-6	-7	36	49	42
Š	6	-4	-9	16	81	36
4	10	-3	-5	9	25	15
2	8	-5	-7	25	49	35
5	12	-2	-3	4	9	6
8	16	1	1	1	1	1

9	16	2	1	4	1	2
10	10	3	-5	9	25	-15
13	32	6	17	36	289	102
15	32	8	17	64	289	136
$\Sigma = 70$	$\sum y=150$			$\Sigma X^2 = 204$	$\Sigma Y^2 = 818$	∑XY=360

The equation of regression lines be $Y = \frac{\sum XY}{\sum X^2}(X)$ and $X = \frac{\sum XY}{\sum Y^2}(Y)$ or $y - y = r \frac{\sigma_y}{\sigma_x}(x - x)$ and $x - \overline{x} = r \frac{\sigma_x}{\sigma_y}(y - \overline{y})$ that is $y - 15 = \frac{360}{204}(x - 7)$ and $x - 7 = \frac{360}{818}(y - 15)$

 $\therefore y = 1.76x + 2.68$ and x = 0.44y + 0.4

Coefficient of correlation $r = \pm \sqrt{(\text{coeffof } x)(\text{coeffof } y)} = \sqrt{(1.76)(0.44)} = 0.88$

Example 2. Compute the coefficient of correlation and the equation of regression lines for the data

х	1	2	3	4	5	6	7	
у	9	8	10	12	11	13	14	

Solution:

First we shall prepare the following table of values

	x	у	x - y = z	<i>x</i> 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		
	$\sum x = 28$	$\Sigma y = 77$	$\sum z = -49$	$\sum x^2 = 140$	$\sum y^2 = 875$	$\sum z^2 = 347$		
No	$x x = \sum_{k=1}^{\infty} x^{k}$	$=\frac{28}{-}=4$ ake	$d y = \frac{\Sigma y}{2} = 7$	$\frac{7}{2}$ = 11 akd z	$= \Sigma Z = -49$	-7		
	n	7	n	7	n 7			
Now $x = \frac{\Sigma S}{n} = \frac{28}{7} = 4$ and $y = \frac{\Sigma y}{n} = \frac{77}{7} = 11$ and $z = \frac{\Sigma Z}{n} = \frac{-49}{7} = -7$ $\sigma_{S}^{2} = \frac{\Sigma S^{2}}{n} - (x)^{2} = \frac{140}{7} - (4)^{2} = 4$, $\sigma_{y}^{2} = \frac{\Sigma y^{2}}{n} - (y)^{2} = \frac{875}{7} - (11)^{2} = 4$								
5	n	7		y n	7			
σ	$^{2}=^{\Sigma Z^{2}}-($	$(\underline{z})^2 = {}^{347} -$	$(-7)^2 = 0.5$	57				
Ζ	n	7						

The coefficient of correlation is given by $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y^2}}{2\sigma_x\sigma_y} = \frac{4+4-0.57}{2\sqrt{4}\sqrt{4}} = 0.93$

The regression lines are given by

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

$$\therefore y - 11 = \frac{(0.93)(2)}{2} (x - 4) akd x - 4 = \frac{(0.93)(2)}{2} (y - 11)$$

ie $y - 11 = 0.93(x - 4)$; $x - 4 = 0.93(y - 11)$
Thus $y = 0.93x + 7.28$; $x = 0.93y - 6.23$ are the regression lines

Example 3. if θ is the angle between the two regression lines, show that tan θ

Explain the significance when r = 0 akd $r = \pm 1$

Solution: The equations to the line of regression of y ok x akd x ok y are

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

: Their slopes are $m_1 = r^{\mathcal{G}}_{\alpha_x} akd m_2 = \frac{\sigma_y}{ra_x}$

Thus $\tan \theta = \frac{m_2 - m_1}{1 + m_1 m_2} = \frac{\sigma_y / r \sigma_x - r \sigma_y / \sigma_x}{1 + \sigma_y^2 / \sigma_x^2} = \frac{1 - r^2}{r} \times \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$

When r = 0, $\tan \theta \to \infty$ or $\theta = \pi/2$ i. e. when the variables are independent, the two lines of regression are perpendicular to each other.

When $r = \pm 1$, $\tan \theta = 0$ i.e., $\theta = 0$ or π . Thus the lines of regression coincide i.e., there is perfect correlation between the two variables.

Example 3. In a partially destroyed laboratory record, only the lines of regression of y ok x akd x ok y are available as 4x - 5y + 33 = 0 akd 20x - 9y = 107 respectively. Calculate \overline{x} , \overline{y} and the coefficient of correlation between x akd y.

Solution: since the regression lines pass through $(\overline{x}, \overline{y})$, therefore,

$$4\overline{x} - 5\overline{y} + 33 = 0, \ 20\overline{x} - 9\overline{y} = 107.$$

Solving these equations, we get $\overline{x} = 13$, $\overline{y} = 17$.

Rewriting the line of regression of y ok x as $y = \frac{4}{5}x + \frac{33}{5}$, we get

$$b_{yS} = r \frac{\sigma_y}{\sigma_x} = \frac{4}{5}$$
 (i)

Rewriting the line of regression of x ok y as $x = \frac{9}{20}y + \frac{107}{9}$ we get

$$b_{Sy} = r \frac{\sigma_x}{\sigma_y} = \frac{9}{20}$$
(ii)

Multiplying (i) and (ii), we get $r^2 = \frac{4}{5} \times \frac{9}{20} = 0.36$ \therefore r = 0.6

Hence r = 0.6, the positive sign being taken as b_{yS} and b_{Sy} both are positive.

Example 4. 8x - 10y + 66 = 0 akd 40x - 18y = 214 are the two regression lines. Find the mean of x s , y s and the coefficient of correlation. Find σ_y if $\sigma_s = 3$

Solution: Given that the regression lines: $8x - 10y + 66 = 0 \Rightarrow 8x - 10y = -66$ and 40x - 18y = 214 we know that the regression lines passes through x = akd y

Therefore $8\overline{x} - 10\overline{y} = -66$ ----- (i)

$$40\overline{x} - 18\overline{y} = 214$$
 ---- (ii

By solving we get $\overline{x} = 13$ akd $\overline{y} = 17$ are means of x s akd y s.

We shall now rewrite the equations of the regression lines to find the correlation coefficient

10y = 8x + 66 or y = 0.8x + 6.6-----(iii), 40x = 18y + 214 or x = 0.45y + 5.35--- (iv)

Correlation coefficient $r = \sqrt{(coeffof x)(coeffof y)} = \sqrt{(0.8)(0.45)} = \pm 0.6$ Also $\sigma_{s} = 3$ by the data and we have $r^{\frac{\sigma_{y}}{2}} = 0.8$ i.e., $(0.6)\sigma_{y} = 2.4$ thus $\sigma_{y} = 4$

	x series	y series
Mean	18	100
S.D	14	20

and r = 0.8, write down the equation of the regression lines and hence find the most probable value of y when x = 70.

Solution: By data
$$\overline{x} = 18 \ akd \ \overline{y} = 100$$
, $\sigma_{\rm S} = 14 \ akd \ \sigma_{\rm y} = 20$

We have the equations of the regression lines

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

i. e., $y - 100 = 0.8 \times \frac{20}{14} (x - 18) akd x - 18 = 0.8 \times \frac{14}{20} (y - 100)$

i.e., y - 100 = 1.14(x - 18); x - 18 = 0.56(y - 100)

Thus y = 1.14x + 79.48; x = 0.56y - 38 these are the regression lines.

When x = 70, we obtain from the first equation y = 1.14(70) + 79.48 = 159.28

Try these

Example 1. Find the coefficient of correlation for the following data

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

Example 2. Compute \overline{x} , \overline{y} and r from the following equation of the regression lines:

2x + 3y + 1 = 0; x + 6y - 4 = 0

Example3. The two regression equations of variables x & y are x = 19.13 - 0.87y akd y = 11.64 - 0.50x. Find a) mean of x's b) mean of y's c) the correlation coefficient betweek x & y.

Curve Fitting

The process of determining a curve of best fit is called curve fitting. This method is called as the method of least square sense. It is generally used for curve fitting.

I) Fitting of a straight line y = a + bx:

Consider a set of n – given values (x, y) for fitting the straight line y = a + bx, where a and b are parameters to be determined. They can be determined by using the following normal equations for fitting the straight line y = a + bx in the least squares sense.

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

$$(\sum x) a + (\sum x^2) b = \sum x y$$
(1)
(2)

Examples:

1) Fit a straight line y = a + bx in the least square sense for the data

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

Answer: Here n = 8 values of (x, y) are given.

We know the normal equations for straight line y = a + bx.

$$n a + (\sum x)b = \sum y \tag{1}$$

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

We prepare a table with respect to equations (1) and (2)

X	у	<i>x</i> ²	<i>x y</i>
1	1	1	1
3	2	9	6
4	4	16	16
6	4	36	24
8	5	64	40
9	7	81	63
11	8	121	88
14	9	196	126
$\sum x = 56$	$\sum y = 40$	$\sum x^2 = 524$	$\sum x y = 364$

Put these table values in equations (1) and (2), we have

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

 $56 a + 524 b = 364$

Solving these equations simultaneously or by calculator, we get

$$a = 0.5454, \quad b = 0.6363$$

Put these values in y = a + bx, we get

$$y = 0.5454 + 0.6363 x$$

2. Find the equation of the best fitting straight line y = ax + b for the following data

ſ	х	5	10	15	20	25
	у	16	19	23	26	30

Answer: Here n = 5 values of (x, y) are given.

We know the normal equations for straight line y = ax + bBefore to write normal equations, we write straight line as y = b + a xNow normal equations are

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

$$\sum x b + (\sum x^2) a = \sum x y$$
(1)
(2)

We prepare a table with respect to equations (1) and (2)

x	У	x ²	x y
5	16	25	80
10	19	100	190
15	23	225	345
20	26	400	520
25	30	625	750
$\sum x = 75$	$\sum y = 114$	$\sum x^2 = 1375$	$\sum x y = 1885$

Put these table values in equations (1) and (2), we have

$$5 b + 75 a = 114,$$

 $75 b + 1375 a = 1885$

Solving these equations simultaneously or by calculator, we get

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

Put these values in y = ax + b, we get y = 0.7 x + 12.3.

3. 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

Р	12	15	21	25
W	50	70	100	120

Where *P* and W are taken in *kg-wt*. Compute *P* when W = 150kg.wt. Answer: Here n = 4 values of (W, P) are given.

We know the normal equations for straight line P = mW + c or P = c + mW.

$$n c + (\Sigma W) m = \Sigma P$$

 $(\Sigma W) c + (\Sigma W^2) m = \Sigma WP$

We prepare the table with respect to equations (1) and (2), we have

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

Put these values in equations (1) and (2), we have

$$4 c + 340 m = 73$$

 $340 c + 3180 m = 6750$

After solving these equations, we get c = 2.2785 and m = 0.1879Put *c*, *m* in equation P = mW + c, we get

$$P = 0.1879W + 2.2785.$$
(3)

Now, we find P by putting W = 150 kg in equation (3), we get

$$P = 0.1879 (150) + 2.2785 \qquad \therefore P = 30.4635$$

II Fitting of a second degree parabola $y = a + bx + cx^2$:

Consider a set of n given values (x, y) for fitting a second degree parabola $y = a + bx + c x^2$, where a, b and c are parameters to be determined. They can be determined by using the following normal equations for fitting the parabola $y = a + bx + c x^2$ in the least squares sense.

$$n a + (\sum x) b + (\sum x^2) c = \sum y$$
⁽¹⁾

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

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

Solving these equations by simultaneously or by calculator, we get *a*, *b*. *c*. Put these in above equations, we get required best fitting parabola

 $y = a + bx + c x^2.$

Examples:

1.	Fit a paral	bola of 2 nd	degree par	abola y =	a + bx + c	x ² for the d	lata
	x	0	1	2	3	4	

У	1	1.8	1.3	2.5	2.3

Answer: Here n = 5 values of (x, y) are given.

We know the normal equations for parabola $y = a + bx + cx^2$

$$n a + (\sum x) b + (\sum x^2) c = \sum y$$
⁽¹⁾

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

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

We prepare the table with respect to equations (1), (2), and (3)

x	У	xy	<i>x</i> ²	x^2y	<i>x</i> ³	x ⁴
0	1	0	0	0	0	0
1	1.8	1.8	1	1.8		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
$\sum x = 10$	$\sum y =$	$\sum x y =$	$\sum x^2 =$	$\sum x^2 y =$	$\sum x^3 =$	$\sum x^4 =$
	8.9	21.1	30	66.3	100	354

Put these in equations (1), (2), and (3) we have

$$5 a + 10 b + 30 c = 8.9$$

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

$$30 a + 100 b + 354 c = 66.3$$

After solving these equations, we get, a = 1.078, b = 0.414, c = -0.021.

Now, put *a*, *b*, *c*, in
$$y = a + bx + cx^2$$

We get, $y = 1.078 + 0.414 \times -0.021 \times^2$.

2. Fit a parabola of second degree parabola $y = a + bx + c x^2$ for the data

х		2	3	4	5	6	7	8	9
у	2	6	7	8	10	11	11	10	9

Estimate y when x = 4.5

Answer: Here n = 9 values of (x, y) are given.

We know the normal equations for parabola $y = a + bx + cx^2$

$$n a + (\sum x) b + (\sum x^2) c = \sum y$$
(1)

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

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

x	У	xy	x^2	x^2y	<i>x</i> ³	<i>x</i> ⁴
1	2	2	1	2	1	1
2	6	12	4	24	8	16
3	7	21	9	63	27	81
4	8	32	16	128	64	256
5	10	50	25	250	125	625
6	11	66	36	396	216	1296
7	11	77	49	539	343	2401
8	10	80	64	640	512	4096
9	9	81	81	729	729	6561
$\sum x =$	$\sum y =$	$\sum x y =$	$\sum x^2 =$	$\sum x^2 y =$	$\sum x^3 =$	$\sum x^4 =$
$\sum_{\substack{X=\\45}} x =$	74	421	285	2771	2025	15333

We prepare the table with respect to equations (1), (2), and (3)

Put these in equations (1), (2), and (3) we have

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

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

$$285a + 2025b + 15333c = 2771$$
After solving these equations, we get,
 $a = -0.9286, b = 3.5232, c = -0.2673.$
Now, put *a*, *b*, *c*, in $y = a + bx + cx^2$, we get
 $y = -0.9286 + 3.5232 \times -0.2673 \times^2.$
(3)

We find y by substituting x = 4.5 in equation (3), we have

$$y = -0.9286 + 3.5232 (4.5) - 0.2673 (4.5)^{2}$$
$$\therefore y = 9.5125$$

III Fitting of an exponential curve $y = ax^b$:

Consider
$$y = ax^b$$
, (1)

Taking *log* on both sides, we have

Take
$$logy = loga + b logx$$
,
 $Take logy = Y$, $loga = A$, $b = B$, $logx = X$, we get
 $Y = A + BX$
(2)

This is the equation of straight line, the normal equation for straight are

$$nA + (\sum X)B = \sum Y \tag{3}$$

$$(\sum X) A + (\sum X^2) B = \sum XY$$
(4)

Solving these equations (3) and (4), we will get A and B,

But, we have to calculate *a* and *b*.

Here, we have taken loga = A \therefore $a = e^A$ and b = BFinally, substitute *a* and *b* in $y = ax^b$. This gives the curve of best fit.

Examples:

1. Fit a least square geometric curve $y = ax^b$ from the following data

	1 8		<i>.</i>		8				
x	1	2	3	4	5				
У	0.5	2	4.5	8	12.5				
Answer:	Here $n = 5$, v	values of (x, y)	y) are given						
Consider	$y = ax^b$								
Taking <i>lo</i>	og on both side	es, we have							
	logy = loga	+ b logx,							
Take <i>logy</i>	v = Y, loga =	A, b = B,	log x = X, w	e get					
$Y = A + BX \tag{2}$									
This is the equation of straight line, the normal equation for straight are									
$nA + (\sum X)B = \sum Y $ (3)									
		· _		1 (1)	(4)				
we prepare ti	he table with re	espective equ							
<i>x</i>	X = logx	у у	Y = logy	X ²	XY				
1	0	0.5	-0.6931	0	0				
2	0.6931	2	0.6931	0.4804	0.4804				
3	1.0986	4.5	1.5041	1.2069	1.6524				
4									
5 1.6094 12.5 2.5257 2.5902 4.									
	$\Sigma X = \Sigma Y = \Sigma X^2 = \Sigma XY =$								
	4.7874		6.1092	6. 1993	9.0804				
out these valu	ues in equation	s $\overline{(3)}$ and $\overline{(4)}$, we have						
			D (100)	h					

5A + 4.7874B = 6.1092

4.7874 A + 6.1993 B = 9.0804

After solving these equations, we get A = -0.6929 and B = 1.9998.

Here, we have taken, loga = A $\therefore a = e^A$

 $a = e^{-0.9629}$ $a = 0.5001 \approx 0.5$ and b = B $b = 1.998 \approx 2$,

Put a and b in $y = ax^b$, we get

$$y = 0.5x^2$$

This is the curve of best fit.

2. An experiment gave the following values:

v(ft/min)	<i>v</i> (<i>ft</i> / <i>min</i>) 350		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.

Answer: Here n = 4, values of (t, v) are given

Consider $v = at^b$

Taking *log* on both sides, we have

logv = loga + b logt,

Take logv = Y, loga = A, b = B, logt = X, we get Y = A + BX

This is the equation of straight line, the normal equation for straight are

(3) (4)

$$nA + (\sum X)B = \sum Y$$

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

We prepare the table with respective equations (3) and (4)

t	X = logt	ν	Y = logv	X ²	XY
61	1.7853	350	2.5441	3.187	4.542
26	1.4150	400	2.6021	2.002	3.682
7	0.8451	500	2.6990	0.714	2.281
2.6	0.4150	600	2.7782	0.172	1.153
	$\sum X =$		$\sum Y =$	$\sum X^2 =$	$\sum XY =$
	4.4604		10.6234	6.075	11.658

Put these data in equations (3) and (4), we have

4A + 4.4604B = 10.6234

$$4.4604 A + 6.075 B = 11.658$$

After solving these equations, we get A = 2.845 and B = -0.1697. Here, we have taken, loga = A $\therefore a = e^A$

 $a = e^{2.845}$ a = 699.8and b = B b = -0.1697,

Practice Problems:

1.	1. Predict y at $x = 3.75$, by fitting a power curve $y = ax^b$ to the given data:									
	x	1	2		3		4		5	6
	у	2.98	4.2	26	5.21		6.10		6.80	7.50
-	Answer:	<i>y</i> = 2.97	8x ^{0.}	⁵¹⁴³ C	ınd y	= 5.	8769	at x	= 3.75	
2.	Fit a para	bola $y = a$	+ k	bx + cx	2 to the	ne fol	lowing	data	:	
	x	2		4			6		8	10
	у	3.07		12.8	35	31	.47	4	57.38	91.29
	Answer:	y = 0.34	- 0). 78x +	- 0.99	9x ²				
3.	3. Fit a second degree parabola									
	x 0 1 2 3 4									
	у	1		1.8	3]	.3		2.5	6.3
	Answer:	y = 1.42	- 1	.07x +	- 0.5	$5x^2$				
4.	Fit a seco	nd degree p	bara	bola						
	x	1.5		2.0	2	.5	3.0		3.5	4.0
	у	1.3		1.6	2.	.0	2.7		3.4	4.1
	Answer:	y = 1.04	- 0). 198x	+ 0.2	244x ²	2		-	
5.	The resul	ts of meas	ure	ment of	f elec	tric r	esistanc	e R	of a cop	oper bar at
	various te	mperature	t°С	are liste	ed bel	ow:				

various (temperatu	re t°C are	listed belo	ow:			
+	10	25	30	36	40	15	50

R 76 77 79 80 82 83 85	t	- /	19	25	30	36	40	45	50
	ŀ	2	76	77	79	80	82	83	85

Find a relation R = a + bt where a and b are constants to be determined. Answer: R = 70.052 + 0.292t.

6. Find the best possible curve of the form y = a + bx, using the method of least squares for the data:

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

Answer: y = 0.545 + 0.636x