

UNIT - 3

Statistical Technique I :-

overview of measures of central tendency, moments, skewness, kurtosis, curve fitting, Method of least squares, fitting of straight lines, fitting of second degree parabola, exponential curves, correlation and Rank Correlation, Regression Analysis : Regression lines of y on x and x on y ,

Statistics :- Statistics is the science which deals with methods of collecting, classifying, presenting, comparing and Interpreting numerical data collected to throw light on any sphere of enquiry.

Variable (variate) :-

A quantity which can vary from one individual to another is called a variable / variate.
e.g. heights, weights, age, rainfall records of cities.

Quantities which can take any numerical value within a certain range are called continuous variables e.g. as the child grows, his/her height takes all possible values from 50 cm to 100 cm.

Quantities which are incapable of taking all possible values are called discrete or discontinuous variables.

(1)

Measure of Central Tendency: →

OR Averages.

A figure which is used to represent a whole (data) series should neither have the lowest value nor the highest in the series, but a value somewhere between these two limits, possibly in the centre, where most of the items of the series cluster. Such figures are called 'Measure of Central Tendency' (or averages).

"Averages are statistical constants which enable us to comprehend in a single effort the significance of the whole."

There are five types of averages

Average

Mathematical Average

1. Arithmetic mean / mean
2. Geometric mean
3. Harmonic mean

Positional Average

1. Median
2. Mode.

Arithmetic mean can be found for Individual series (i.e. where frequency is not given)
For discrete series
For continuous series.

Arithmatic Mean :-

I) For Individual Series / ungrouped data.

if x_1, x_2, \dots, x_n are n variables Then

1) Direct Method :-

$$\text{Mean } (\bar{x}) = \frac{x_1 + x_2 + \dots + x_n}{n} = \frac{\sum x}{n}$$

2) Short cut method :- This is shift of origin
shifting the origin to an arbitrary point a ,
the formula $\bar{x} = \frac{\sum d}{n}$ becomes

$$\bar{x}-a = \frac{1}{n} \sum (x-a)$$

$$\bar{x} = a + \frac{1}{n} \sum (x-a)$$

$$\bar{x} = a + \frac{1}{n} \sum d_n, \text{ where } d_n = x - a$$

a is called assumed mean:

↓
deviation of x from

and n is the no. of observations.

3.) Step deviation:-

$$\bar{x} = a + \frac{c \sum d'_n}{n}, \quad d'_n = \frac{x-a}{c}$$

$= a + \frac{\sum d'_n}{n} \times c$
where c is common value.

Ex:- Daily income of 10 families is given as follows:

S.No	1	2	3	4	5	6	7	8	9	10
Daily Income (\bar{x})	100	120	80	85	95	130	200	250	225	275

find Average daily Income.

S.No	Daily income (in Rs) (x)	$d = x - a$	$d_x' = \frac{d_x}{C}$
1	100	20	2
2	120	40	4
3	80	0	0
4	85	5	0.5
5	95	15	1.5
6	130	50	5
7	200	120	12
8	250	170	17
9	225	145	14.5
10	275	195	19.5
$\sum x = 1560$		$\sum d_x = 76$	$\sum d_x' = 7.6$

By using Direct Method -

Here $n = 10$

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

$$= \frac{1560}{10}$$

$$= 156$$

Hence average daily income is Rs. 156

By using short cut method -

Let assumed mean $a = 80$

$$\bar{x} = a + \frac{\sum d_x}{n} = 80 + \frac{76}{10}$$

$$= 80 + 7.6 = 156$$

Step deviation :-

$$\text{let } c = 10$$

Then

$$\text{Mean } \bar{x} = a + \frac{\sum f d_u'}{n} \times c$$

$$= 80 + \frac{76}{10} \times 10$$

$$= 80 + 76$$

$$= \underline{\underline{156}}$$

Arithmetic Mean for Discrete Series :-

1) Direct Method : - If the frequency distribution is

$$x: x_1, x_2, \dots, x_n$$

$$f: f_1, f_2, \dots, f_n \text{ Then}$$

$$\text{Mean } \bar{x} = \frac{f_1 x_1 + f_2 x_2 + \dots + f_n x_n}{f_1 + f_2 + \dots + f_n}$$

$$\boxed{\bar{x} = \frac{\sum f x}{\sum f}} = \frac{\sum f x}{N} \quad N = \text{sum of frequencies}$$

$$= \sum f$$

2) Short cut method :-

$$\bar{x} = a + \frac{1}{N} \sum f d_x$$

$$d_x = x - a$$

a = assumed mean

$$N = \sum f$$

3) Step deviation method :-

$$\boxed{\bar{x} = a + \frac{\sum f d_u'}{N} \times c}$$

$$d_u' = \frac{x - a}{c}$$

c = step deviation

Ex :- Calculate the mean of the following frequency distribution of marks in a test in Mathematics:

marks	10	20	30	40	50	60	70	80
No. of Students	3	6	10	12	9	6	2	2

Sol:-

Marks <u>x</u>	No. of Students <u>f</u>	$\sum f_x$	$d_x = x - a$	$\sum fd_x$	$d'_x = \frac{x-a}{h}$	$\sum fd'_x$
10	3	30	-30	-90	-3	-9
20	6	120	-20	-120	-2	-12
30	10	300	-10	-100	-1	-10
40	12	480	0	0	0	0
50	9	450	10	90	1	9
60	6	360	20	120	2	12
70	2	140	30	60	3	6
80	2	160	40	80	4	8
$\sum f = 50$		$\sum f_x = 2040$		$\sum fd_x = 40$	$\sum d'_x$	$\sum fd'_x = 4$

By Direct Method:-

$$\text{A. Mean } \bar{x} = \frac{\sum f_x}{\sum f} = \frac{2040}{50} = \frac{204}{5} = 40.8$$

By Shortcut Method :- $a = 40$

$$\begin{aligned}\bar{x} &= a + \frac{\sum fd_x}{\sum f} = 40 + \frac{40}{50} \\ &= 40 + 0.8 = 40.8\end{aligned}$$

By step deviation method:-

$$\text{deviation } h = 10$$

$$\bar{x} = a + \frac{\sum f d_x'}{\sum f} \times h$$

$$= 40 + \frac{4 \times 10}{50} = 40 + 0.8$$

$$= \underline{\underline{40.8}}$$

Arithmetical Mean for Continuous:-

1.) Direct Method:-

$$\bar{x} = \frac{\sum fm}{N} = \frac{\sum fm}{\sum f}$$

2.) Short cut Method:-

$$\bar{x} = a + \frac{\sum f d_x}{N}$$

$d_x = x - a$

3.) Step deviation Method:-

$$\bar{x} = a + \frac{\sum f d_x'}{N} \times h$$

$$d_x' = \frac{x - a}{h}$$

where a is assumed value mean

m is the mid pt of various class

h is step deviation.

Ex:- find the Arithmetic Mean for the following data.

Marks	No. of students
0-10	5
10-20	10
20-30	40
30-40	20
40-50	25

Sol:-

Marks	Mid-values m	No. of students f	$\sum fm$	d	$\sum fd'$	d'	$f d'$
0-10	5	5	25	-20	-100	-2	-10
10-20	15	10	150	-10	-100	-1	-10
20-30	25	40	1000	0	0	0	0
30-40	35	20	700	10	200	1	20
40-50	45	25	1125	20	500	2	50

$N = \sum f = 100$ $\sum fm = 3000$ $\sum fd' = 500$ $\sum f d' = 50$

By D.M

$$\therefore \bar{x} = \frac{\sum fm}{N} = \frac{\sum fm}{\sum f}$$

$$= \frac{3000}{100}$$

$$= 30$$

Short cut Method :- $a = 25$

$$\bar{x} = a + \frac{\sum fd'}{N}$$

$$= 25 + \frac{500}{100} = 30$$

Step

Mat

7. The
the

Proof:-

Then

Step deviation method :-

here $h = 10$

Then

$$\bar{x} = a + \frac{\sum fd'}{N} \times h$$

$$= 25 + \frac{50}{100} \times 10$$

$$= 25 + 5$$

$$= 30$$

d' $\sum fd'$

-2 -10

-1 -10

0 0

1 20

2 50

$\sum fd' = 50$

Mathematical Properties of Arithmetic Mean:-

I. The Total of the deviations of the items from the mean is equal to zero.

Proof:-

Let $x_1, x_2, \dots, x_n \rightarrow x_i^o, i=1, 2, \dots, n$.

then

$$\sum_{i=1}^n (x_i^o - \bar{x}) = \sum_{i=1}^n x_i^o - \sum_{i=1}^n \bar{x}$$

$$= \sum_{i=1}^n x_i^o - \bar{x} \sum_{i=1}^n 1 \quad \because \sum_{i=1}^n 1 = n$$

$$= \sum_{i=1}^n x_i^o - \bar{x} (n)$$

$$= \sum_{i=1}^n x_i^o - n\bar{x} = \sum_{i=1}^n x_i^o - n \frac{\sum_{i=1}^n x_i}{n}$$

$$= 0$$

2. combined Mean of two groups:-

Let \bar{X}_1 = mean of first series

\bar{X}_2 = mean of second series

N_1 = No. of observation in first group

N_2 = " " in second group

Then combined mean of two group

$$\bar{X}_{1,2} = \frac{N_1 \bar{X}_1 + N_2 \bar{X}_2}{N_1 + N_2}$$

Ex :- A cooperative bank has two branches employing 50 & 70 workers respectively. The Average salaries paid by two respective branches are Rs. 360 and Rs 390 per month calculate the mean of the salaries of all the employees.

Sol:-

Mean of the salaries of all the employees

$$\bar{X}_{1,2} = \frac{N_1 \bar{X}_1 + N_2 \bar{X}_2}{N_1 + N_2}$$

$$= \frac{50 \times 360 + 70 \times 390}{50 + 70}$$

$$= \frac{45300}{120}$$

$$= 377.5$$

Ques. 1 The mean of 200 items was 50. Later on it was discovered that two items were misread as 92 and 8 instead of 192 and 88. Find correct mean.

Sol:-

Here incorrect value of $\bar{x} = 50$, $n = 200$

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

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

Using incorrect value of \bar{x}

$$\sum x = 200 \times 50 = 10,000$$

$$\begin{aligned}\therefore \text{corrected value of } \sum x &= 10000 - (92+8) + \\ &\quad (192+88) \\ &= 10000 - 100 + 280 \\ &= 10180\end{aligned}$$

$$\begin{aligned}\text{Corrected Mean} &= \frac{\text{corrected } \sum x}{n} \\ &= \frac{10180}{200} \\ &= 50.9\end{aligned}$$



Weighted Arithmetic Mean :-

If the variate-values are not of equal importance, we may attach to them 'weight' w_1, w_2, \dots, w_n as measures of their importance. The weighted mean \bar{x}_w is defined as

$$\bar{x}_w = \frac{w_1x_1 + w_2x_2 + \dots + w_nx_n}{w_1 + w_2 + \dots + w_n} = \frac{\sum w_i x_i}{\sum w_i}$$

Ex:- calculate weighted mean by weighting each price by the quantity consumed.

8867

Articles of food	Quantity Consumed in kg	Price in Rs per kg)
Flour	11.50	5.8
Ghee	5.60	58.4
Sugar	0.28	8.2
Potato	0.16	2.5
Oil	0.35	20.0

Soln:-

$$\sum w = 17.89$$

<u>wx</u>
66.700
327.04
2.296
0.400
7.000

$$\sum wx = 403.436$$

$$\begin{aligned} \bar{x}_w &= \frac{\sum wx}{\sum w} \\ &= \frac{403.436}{17.89} \\ &= 22.55 \end{aligned}$$

weighted Mean price.

Ques 2 compute the arithmetic mean for the following data :

class Height (in cm)	219	216	213	210	207	204	201	198	195
No. of Persons	2	4	6	10	11	7	5	3	1

Seri:- Height in(cm) x	No. of Persons f	fx
219	2	438
216	4	864
213	6	1278
210	10	2100
207	11	2277
204	7	1428
201	5	1005
198	4	792
195	1	195
$\sum f = 40$		$\sum fx = 10,377$

Arithmetic Mean By Direct Method

$$\bar{x} = \frac{\sum fx}{\sum f}$$

$$= \frac{10,377}{40}$$

$$= \underline{259.425}$$

Geometric Mean :-

(a) Geometric mean for Individual series →

Geometric mean (G.M) of n individual observations x_1, x_2, \dots, x_n ($x_i \neq 0$) is the n^{th} root of their product.

Thus

$$G = (x_1 \cdot x_2 \cdot \dots \cdot x_n)^{1/n}$$

after taking log on both sides, we get

$$G = \text{antilog} \left[\frac{1}{n} \sum_{i=1}^n \log x_i \right]$$

(b) Geometric mean for discrete series →

If x_1, x_2, \dots, x_n occur f_1, f_2, \dots, f_n times resp. and $N = \sum_{i=1}^n f_i$, then G.M is given by

$$\text{G.M}, G = (x_1^{f_1} \cdot x_2^{f_2} \cdot \dots \cdot x_n^{f_n})^{1/N}$$

taking logarithm on both sides

$$\log G = \log (x_1^{f_1} \cdot x_2^{f_2} \cdot \dots \cdot x_n^{f_n})^{1/N}$$

$$= \frac{1}{N} [f_1 \log x_1 + f_2 \log x_2 + \dots + f_n \log x_n]$$

$$G = \text{antilog} \left[\frac{1}{N} \sum_{i=1}^n f_i \log x_i \right]$$

c) G.M for continuous series →

The Geometric mean for continuous series can be obtained by finding out the mid-value of the interval and by using the concept of G.M for discrete series.

Ex-1 Calculate the Geometric mean for the following data:

x	12	13	14	15	16	17	
f	5	4	4	3	2	2	

data:-	x	f	$\log x$	$f \log x$
12	5	1.0792	5.3960	
13	4	1.1139	4.4556	
14	4	1.1461	4.5844	
15	3	1.1761	3.5283	
16	2	1.2041	2.4092	
17	2	1.2304	1.2304	
		$\sum f = 20$	$\sum f \log x = 21.6029$	

we have

$$G.M = \text{Antilog} \left[\frac{1}{N} \sum_{i=1}^n f_i \log x_i \right]$$

$$N = \sum f_i$$

$$= \text{Antilog} \left[\frac{21.6029}{19} \right]$$

$$= \text{Antilog} (1.137)$$

$$= 10^{1.1373} = 13.71$$

Ex:-2 find the Geometric mean for the following data:

marks	0 - 10	10 - 20	20 - 30	30 - 40	40 - 50
No. of students	4	8	10	6	7

marks α	No. of students f	mid values x	$\log x$	$f \log x$
0 - 10	4	5	0.6990	2.7960
10 - 20	8	15	1.1761	9.4088
20 - 30	10	25	1.3979	13.9790
30 - 40	6	35	1.5441	9.2646
40 - 50	7	45	1.6532	11.5724
$\sum f = N = 35$				$\sum f \log x = 47.0208$

$$\text{So } G.M = \text{Antilog} \left[\frac{1}{N} \sum_{i=1}^n f_i \log x_i \right]$$

$$= \text{Antilog} \left[\frac{47.0208}{35} \right]$$

$$= \text{Antilog} [1.3435]$$

$$= 10^{1.3435} = 22.06$$



Harmonic Mean: \rightarrow

for Individual series -

Harmonic mean of a no. of observations is the reciprocal of the Arithmetic mean of the reciprocals of the given values. Thus, H.M H of n observations x_1, x_2, \dots, x_n is

$$H = \frac{1}{\frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}} = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

for discrete series.

if x_1, x_2, \dots, x_n ($\neq 0$) have the frequencies f_1, f_2, \dots, f_n resp. then harmonic mean

$$H = \frac{1}{\frac{1}{N} \sum_{i=1}^n \frac{f_i}{x_i}} = \frac{N}{\frac{f_1}{x_1} + \frac{f_2}{x_2} + \dots + \frac{f_n}{x_n}}, N = \sum_{i=1}^n f_i$$

for continuous series (class-intervals)

x_i is taken to be the mid-value of class-int.

Ex:- An Aeroplane flies along the four sides of a square at speeds of 100, 200, 300 and 400 km/hr resp., what is the average speed of the aeroplane in its flight around the square?

Ans:-

When equal distances are covered with unequal speeds, the Harmonic Mean is the proper average.

\therefore Average speed = Harmonic Mean

$$= \frac{4}{\frac{1}{100} + \frac{1}{200} + \frac{1}{300} + \frac{1}{400}}$$

$$= 192 \text{ Km/hr.}$$

Ex:-2 Find out the Harmonic mean of the following data:

marks (out of 150)	No. of students
10	2
20	3
40	6
60	5
120	4

marks x f	No. of students f	$\frac{1}{x}$	$f \times \frac{1}{x}$
10	2	0.100	0.200
20	3	0.050	0.150
40	6	0.025	0.150
60	5	0.017	0.085
120	4	0.008	0.032
$\sum f = N = 20$			$\sum \frac{f}{x} = 0.617$

Harmonic mean

$$= \frac{N}{\sum \frac{f}{x}}$$

$$= \frac{20}{0.617} = 32.4$$

Median : →

median is the central value of the variables when the values are arranged in ascending or descending order of magnitude.

1. for ungrouped / Individual frequency distribution →

If the n values of the variate are arranged in ascending or descending order of magnitude.

a.) if n is odd Then

$$\text{median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ value} = \text{middle value}$$

if n is even Then

$$\text{median} = \frac{\left(\frac{n}{2} \right)^{\text{th}} \text{ value} + \left(\frac{n}{2} + 1 \right)^{\text{th}} \text{ value}}{2}$$

2. for discrete frequency distribution →

median is obtained by cumulative frequencies

Find $\frac{N+1}{2}$, where $N = \sum f_i$

Then

$$\text{Cof just } \geq \frac{N+1}{2}$$

The corresponding value of x is median.

3. For a Grouped frequency distribution →

$$\text{Median} = l + \frac{f}{2} \left(\frac{N}{2} - c \right)$$

Median class is the class corresponding to cumulative frequency $\geq \frac{N+1}{2}$ or $\frac{N}{2}$

where, l = lower limit of median class

h = width of Median class

f = frequency of Median class

$$N = \sum f$$

$C = C.o.f.$ of the class preceding the median class.

Ex:-1 Find the median of 6, 8, 9, 10, 11, 12, 13.

Soln:-

Arranging given values in ascending order

6, 8, 9, 10, 11, 12, 13

here $n = 7$ (odd)

So

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ value}$$

$$= \left(\frac{7+1}{2} \right)^{\text{th}} \text{ value} = 4^{\text{th}} \text{ value}$$

$$= 10 \quad \text{Ans.}$$

Ex:-2. obtain the median for the following frequency distribution -

$x : 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 9$

$f : 8 \ 10 \ 11 \ 16 \ 20 \ 25 \ 15 \ 9 \ 6$

Soln:-

The given frequency distribution is discrete then cumulative frequency distribution table -

<u>x</u>	<u>f</u>	<u>c.f</u>	
1	8	8	Here $N = 120$
2	10	18	
3	11	29	$\therefore \frac{N+1}{2} = \frac{120+1}{2}$
4	16	45	$= 60.5$
5	20	65	
6	25	90	
7	15	105	$c.f$ just $\geq \frac{N+1}{2} = 60.5$
8	9	114	
9	6	120	is 65
$N = \sum f = 120$			

and the value of x corresponding to
~~c.f~~ 65 is 5
Hence median is 5.

Ex: 13. find the median for the following data:

marks	No. of students	marks	No. of Students
Below 10	15	Below 50	94
Below 20	35	Below 60	127
Below 30	60	Below 70	198
Below 40	84	Below 80	249

Sol:- The cumulative freq. table with class Intervals

marks	No. of students (f)	c.f
0-10	15	15
10-20	20	35
20-30	25	60
30-40	24	84
40-50	10	94
50-60	33	127
60-70	71	198
70-80	51	249

$$\sum f = N = 249$$

$$\frac{N}{2} = 124.5$$

\therefore median class in $50-60$ $l = 50$
 $h = 10, f = 33, c = 94$

$$\text{Median} = l + \frac{h}{f} \left(\frac{N}{2} - c \right)$$

$$= 50 + \frac{10}{33} (124.5 - 94)$$

$$= 59.24 \text{ marks.}$$

#

Mode: \rightarrow

Mode is the value which occurs most frequently in a set of observations and around which the other items of the set cluster densely.

OR mode is that value of the variate for which frequency is maximum.

- a) for discrete frequency distribution
 mode is the value of x corresponding to maximum frequency
 But in case of the following
 i.) maximum freq. is repeated
 ii.) maximum freq. occurs in the very beginning or at the end of the distribution
 iii.) if there is irregularities in the distribution, the value of the mode is determined by method of grouping.

b) for continuous frequency distribution

$$\text{Mode} = l + \frac{f_m - f_1}{2f_m - f_1 - f_2} \times h$$

where l is the lower limit

h is the width

f_m is the frequency of modal class
 f_1, f_2 are the frequencies of the classes preceding and succeeding the modal class respectively.

for a symmetric distribution, mean, median, mode coincide.

* if method of grouping fails then

$$\boxed{\text{mode} = 3 \text{median} - 2 \text{mean}}$$

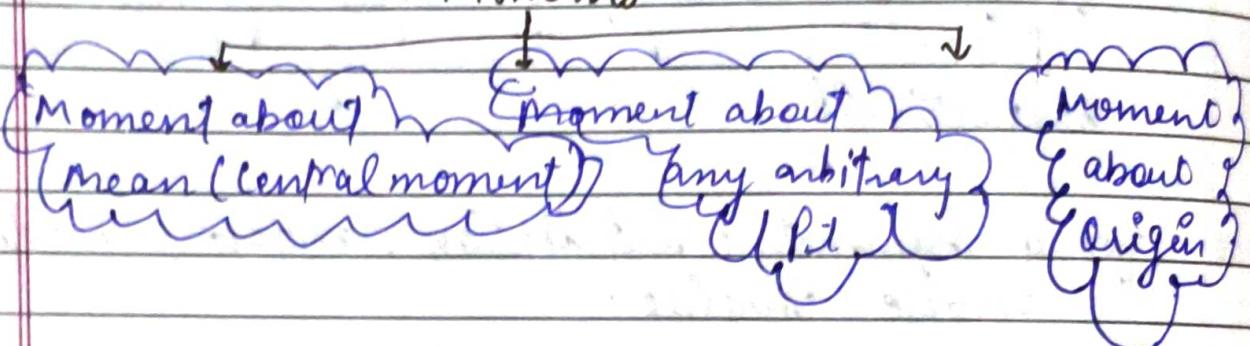
called empirical mode.

Imp Moments : →

Moments are statistical tools, used in statistical investigations.

Moments of a distribution are the arithmetic means of various powers of the deviation of items from some given number.

Moments



7. Moments about Mean (Central Moments)

(a) for an Individual Series →

If x_1, x_2, \dots, x_n are the values of the variable under consideration,

The r th moment M_r about mean \bar{x} is

$$M_r = \frac{\sum_{i=1}^n (x_i - \bar{x})^r}{n}; \quad r = 0, 1, 2, \dots$$

(b) for a frequency distribution →

if the frequency distribution is

$x: x_1, x_2, \dots, x_n$

$f: f_1, f_2, \dots, f_n$ then the r th moment M_r about the mean \bar{x} is defined as

$$\mu_r = \frac{\sum_{i=1}^n f_i (x_i - \bar{x})^r}{N} ; r = 0, 1, 2, 3, \dots$$

$$N = \sum_{i=1}^n f_i$$

for $r=0$,

$$\begin{aligned}\mu_0 &= \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^0 \\ &= \frac{1}{N} \sum_{i=1}^n f_i = \frac{N}{N} = 1\end{aligned}$$

$$\mu_0 = 1$$

\therefore for any distribution,

$$\boxed{\mu_0 = 1}$$

for $r=1$,

$$\begin{aligned}\mu_1 &= \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x}) \\ &= \frac{1}{N} \sum_{i=1}^n f_i x_i - \frac{1}{N} \sum_{i=1}^n f_i \bar{x} \\ &= \frac{1}{N} \sum_{i=1}^n f_i x_i - \bar{x} \frac{1}{N} \sum_{i=1}^n f_i = \frac{1}{N} \sum_{i=1}^n f_i \\ &= \frac{1}{N} \sum_{i=1}^n f_i x_i - \bar{x} \cdot 1\end{aligned}$$

$$= \bar{x} - \bar{x} = 0$$

\therefore for any distribution

$$\boxed{\mu_1 = 0}$$

for $r=2$

$$\mu_2 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^2 = \text{variance} = (S_o)^2$$

\therefore for any distribution, $\mu_2 = \text{variance of the distribution}$.

$S \cdot D = \text{Standard deviation}$

Similarly,

$$\mu_3 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^3$$

f

$$\mu_4 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^4$$

& soon

Ex:- find the first four moments for the following individual series:

x	3	6	8	10	18
-----	---	---	---	----	----

Soln:-

s.no.	x	$x - \bar{x}$	$(x - \bar{x})^2$	$(x - \bar{x})^3$	$(x - \bar{x})^4$
1	3	-6	36	-216	1296
2	6	-3	9	-27	81
3	8	-1	1	-1	1
4	10	1	1	1	1
5	18	9	81	729	6561
n=5	$\sum x = 45$	$\sum (x - \bar{x}) = 0$	$\sum (x - \bar{x})^2 = 128$	$\sum (x - \bar{x})^3 = 486$	$\sum (x - \bar{x})^4 = 7940$

$$\text{mean } \bar{x} = \frac{\sum x}{n} = \frac{45}{5} = 9$$

first four moments

$$\mu_1 = \frac{\sum (x - \bar{x})}{n} = 0$$

$$\mu_2 = \frac{\sum (x - \bar{x})^2}{n} = \frac{128}{5} = 25.6$$

$$\mu_3 = \frac{\sum (x - \bar{x})^3}{n} = \frac{486}{5} = 97.2$$

$$\mu_4 = \frac{\sum (x - \bar{x})^4}{n} = \frac{7940}{5} = 1588$$

Ex-2. calculate μ_1 , μ_2 , μ_3 , σ^2 for the following frequency distribution

marks	5-15	15-25	25-35	35-45	45-55	55-65
No. of students	10	20	25	20	15	10

Soln:-

Marks	No. of students	mid pt (x)	$f x$	$x - \bar{x}$	$f(x - \bar{x})$
5-15	10	10	100	-24	-240
15-25	20	20	400	-14	-280
25-35	25	30	750	-4	-100
35-45	20	40	800	6	120
45-55	15	50	750	16	240
55-65	10	60	600	26	260
$N = \sum f = 100$			$\sum f x = 3400$		$\sum f(x - \bar{x}) = 0$

$$\text{mean } \bar{x} = \frac{\sum f x}{\sum f} = \frac{3400}{100} = 34.$$

$f(x - \bar{x})^2$	$f(x - \bar{x})^3$	$f(x - \bar{x})^4$
5760	-138240	3377760
3920	-54880	768320
400	-1600	6400
720	4320	25920
8840	61440	983040
6760	175760	4569760
$\sum f(x - \bar{x})^2 = 21400$	$\sum f(x - \bar{x})^3 = 46800$	$\sum f(x - \bar{x})^4 = 9671200$

$$\mu_1 = \frac{\sum f(x - \bar{x})}{N} = \frac{0}{100} = 0$$

$$\mu_2 = \frac{\sum f(x - \bar{x})^3}{N} = \frac{21400}{100} = 214$$

$$M_3 = \frac{\sum f(x - \bar{x})^3}{N} = \frac{46800}{100} = 468$$

$$M_4 = \frac{\sum f(x - \bar{x})^4}{N} = \frac{9671200}{100} = 96712$$

2: Moments about an arbitrary number: →
(Raw moments).

(a) For ungrouped data or Individual series →

If x_1, x_2, \dots, x_n are the values of the variable x , then moments about any point A is denoted by μ'_r and is defined as

$$\mu'_r = \frac{1}{n} \sum_{i=1}^n (x_i - A)^r, \quad r=0, 1, 2, 3, \dots$$

For grouped data or for frequency distribution (Discrete or continuous): —

If x_1, x_2, \dots, x_n are the values of a variable x with the corresponding frequencies f_1, f_2, \dots, f_n respectively, then the r th moment about any point A is denoted by μ'_r and is defined as

$$\mu'_r = \frac{1}{N} \sum_{i=1}^n f_i (x_i - A)^r \quad r=0, 1, 2, 3$$

$N = \sum f_i$

$$\text{for } r=0, \quad \mu'_0 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - A)^0 = 1 \quad [\mu'_0 = 1]$$

for $r=1$,

$$\begin{aligned} \mu'_1 &= \frac{1}{N} \sum_{i=1}^n f_i (x_i - A) \\ &= \frac{1}{N} \sum_{i=1}^n f_i x_i - \frac{1}{N} \sum_{i=1}^n f_i A \\ &= \bar{x} - \frac{A}{N} \sum_{i=1}^n f_i \\ &= \bar{x} - \frac{A}{N} \times N \end{aligned}$$

$$\boxed{\mu'_1 = \bar{x} - A}$$

for $r=2$,

$$\mu'_2 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - A)^2$$

for $r=3$,

$$\mu'_3 = \frac{1}{N} \sum_{i=1}^n f_i (x_i - A)^3 \text{ and so on.}$$

3. Moments about the origin \Rightarrow

If x_1, x_2, \dots, x_n be the values of a variable x with corresponding frequencies f_1, f_2, \dots, f_n respectively then r th moment about the origin ν_r is defined as

$$\boxed{\nu_r = \frac{1}{N} \sum_{i=1}^n f_i x_i^r}, \quad r=0, 1, 2, \dots$$

$N = \sum_{i=1}^n f_i$

$$\text{For } r=0, \quad \nu_0 = \frac{1}{N} \sum_{i=1}^n f_i x_i^0 = \frac{N}{N} = 1 \quad \Rightarrow \boxed{\nu_0 = 1}$$

$$\text{For } r=1, \quad \nu_1 = \frac{1}{N} \sum_{i=1}^n f_i x_i^1 = \bar{x} \quad \boxed{\nu_1 = \bar{x}}$$

Relation between μ_x and μ_x' : \rightarrow

We know that

$$\mu_x = \frac{\sum_{i=1}^n f_i (x_i - \bar{x})^r}{N} \quad \dots \quad (1)$$

and

$$\mu_x' = \frac{\sum_{i=1}^n f_i (x_i - A)^r}{N} \quad \dots \quad (2)$$

from (1),

$$\begin{aligned}
 \mu_x &= \frac{1}{N} \sum_{i=1}^n f_i [(x_i - A) - (\bar{x} - A)]^r \\
 &= \frac{1}{N} \sum_{i=1}^n f_i [(x_i - A) - \mu_1']^r \quad \therefore \mu_1' = \bar{x} - A \\
 &= \frac{1}{N} \sum_{i=1}^n f_i [{}^r C_0 (x_i - A) {}^r C_0 (\mu_1')^0 - {}^r C_1 (x_i - A) {}^r C_1 (\mu_1')^1 \\
 &\quad + {}^r C_2 (x_i - A) {}^r C_2 (\mu_1')^2 - \dots + (-1) {}^r C_r (\mu_1')^r] \\
 &= \frac{1}{N} \sum_{i=1}^n f_i (x_i - A)^r - \frac{{}^r C_1}{N} \sum_{i=1}^n f_i (x_i - A)^{r-1} (\mu_1')^1 \\
 &\quad + \frac{{}^r C_2}{N} \sum_{i=1}^n f_i (x_i - A)^{r-2} (\mu_1')^2 + \dots + (-1) \frac{{}^r C_r}{N} \sum_{i=1}^n f_i (\mu_1')^r \\
 &= \mu_1' - {}^r C_1 \mu_2' + {}^r C_2 \mu_3' - {}^r C_3 \mu_4' + \dots + (-1) {}^r C_r (\mu_1')^r
 \end{aligned}$$

Binomial Theorem.

we know that $\boxed{\mu_1' = 0}$

Putting $r=2, 3, 4$, we get

$$\mu_2 = \mu_2' - 2\mu_1'^2 + \mu_0' \cdot \mu_1'^2$$

$$= \mu_2' - 2\mu_1'^2 + \mu_1'^2 \quad \mu_0' = 1$$

$$\boxed{\mu_2 = \mu_2' - \mu_1'^2}$$

$\gamma = 3,$

$$M_3 = \mu'_3 - {}^3C_1 \mu'_2 \mu'_1 + {}^3C_2 \mu'_1 \mu'^2 - {}^3C_3 \mu'_0 \mu'^3$$

$$= \mu'_3 - 3\mu'_2 \mu'_1 + 3\mu'^2 - 1 \cdot 1 \mu'^3 \quad : \mu'_0 = 1$$

$$\boxed{M_3 = \mu'_3 - 3\mu'_2 \mu'_1 + 2\mu'^3}$$

$\gamma = 4,$

$$M_4 = \mu'_4 - {}^4C_1 \mu'_3 \mu'_1 + {}^4C_2 \mu'_2 \mu'^2 - {}^4C_3 \mu'_1 \mu'^3 + {}^4C_4 \mu'_0 \mu'^4$$

$$\boxed{M_4 = \mu'_4 - 4\mu'_3 \mu'_1 + 6\mu'_2 \mu'^2 - 3\mu'^4}$$

Relation between V_r and M_r : \rightarrow

We know that

$$\begin{aligned} V_r &= \frac{1}{N} \sum_{i=1}^n f_i^o x_i^r, \quad r = 0, 1, 2, \dots \\ &= \frac{1}{N} \sum_{i=1}^n f_i^o (x_i^o - \bar{x} + \bar{x})^r \\ &= \frac{1}{N} \sum_{i=1}^n f_i^o [{}^r C_0 (x_i^o - \bar{x})^r \bar{x}^0 + {}^r C_1 (x_i^o - \bar{x})^{r-1} \bar{x}^1 \\ &\quad + {}^r C_2 (x_i^o - \bar{x})^{r-2} \bar{x}^2 + \dots + \bar{x}^r] \end{aligned}$$

$$\begin{aligned} V_r &= \frac{1}{N} \sum_{i=1}^n f_i^o (x_i^o - \bar{x})^r + \frac{{}^r C_1}{N} \sum_{i=1}^n f_i^o (x_i^o - \bar{x})^{r-1} \bar{x}^1 \\ &\quad + \frac{{}^r C_2}{N} \sum_{i=1}^n f_i^o (x_i^o - \bar{x})^{r-2} \bar{x}^2 + \dots + \frac{1}{N} \sum_{i=1}^n f_i^o \bar{x}^r \\ &= M_r + {}^r C_1 M_{r-1} \bar{x} + {}^r C_2 M_{r-2} \bar{x}^2 + {}^r C_3 M_{r-3} \bar{x}^3 \\ &\quad + \dots + \bar{x}^r \end{aligned}$$

Putting $r = 1, 2, 3, 4, \dots$

$$V_1 = \mu_1 + \mu_0 \bar{x} = \bar{x} \quad [\mu_1 = 0, \mu_0 = 1]$$

$$\boxed{V_1 = \bar{x}}$$

$\tau = 2,$

$$V_2 = \mu_2 + 2C_1\mu_1\bar{x} + 2C_2\mu_0\bar{x}^2$$

$$= \mu_2 + 2 \cdot 0 \cdot \bar{x} + 1 \cdot 1 \cdot \bar{x}^2$$

$$V_2 = \mu_2 + \bar{x}^2$$

 $\tau = 3,$

$$V_3 = \mu_3 + 3C_1\mu_2\bar{x} + 3C_2\mu_1\bar{x}^2 + 3C_3\mu_0\bar{x}^3$$

$$V_3 = \mu_3 + 3\mu_2\bar{x} + \bar{x}^3$$

$$\begin{cases} \because \mu_1 = 0 \\ \mu_0 = 1 \end{cases}$$

 $\tau = 4,$

$$V_4 = \mu_4 + 4C_1\mu_3\bar{x} + 4C_2\mu_2\bar{x}^2 + 4C_3\mu_1\bar{x}^3 + 4C_4\mu_0\bar{x}^4$$

$$V_4 = \mu_4 + 4\mu_3\bar{x} + 6\mu_2\bar{x}^2 + \bar{x}^4$$

$$V_4 = \mu_4 + 4\mu_3\bar{x} + 6\mu_2\bar{x}^2 + \bar{x}^4$$

(7)

Karl Pearson's β and γ coefficients:-

These coefficients are based upon the first four moments of a frequency distribution about its mean:

$$\beta_1 = \frac{\mu_3^2}{\mu_2^3}, \quad \beta_2 = \frac{\mu_4}{\mu_2^2}$$

$$\gamma_1 = +\sqrt{\beta_1}, \quad \gamma_2 = \beta_2 - 3$$

The practical use of these coefficients is to measure the skewness & kurtosis of a frequency distribution.

Ex21. The first four moments of a distribution about the value y of a variable are $-1.5, 17, -30$ & 108 . Find the moments about mean, about origin. Also find the moments about the point $x=2$.
Also β_1 & β_2

Sol:-

$$\text{Here } A = y$$

$$\mu'_1 = -1.5$$

$$\mu'_2 = 17$$

$$\mu'_3 = -30$$

$$\mu'_4 = 108$$

Moments about Mean \rightarrow

$$\mu_1 = 0$$

$$\mu_2 = \mu'_2 - \mu'_1^2 = 17 - (-1.5)^2 = 14.75$$

$$\begin{aligned}\mu_3 &= \mu'_3 - 3\mu'_2\mu'_1 + 2\mu'_1^3 \\ &= -30 - 3 \times 17 \times (-1.5) + 2 \times (-1.5)^3 \\ &= 39.75\end{aligned}$$

$$\begin{aligned}\mu_4 &= \mu'_4 - 4\mu'_3\mu'_1 + 6\mu'_2\mu'_1^2 - 3\mu'_1^4 \\ &= 142.3125\end{aligned}$$

$$\text{Also } \mu'_1 = \bar{x} - A$$

$$\text{So } \bar{x} = \mu'_1 + A$$

$$= -1.5 + y$$

$$\boxed{\bar{x} = 2.5}$$

~~6.25
14.75
1.00~~

Moment about origin \rightarrow

$$\nu_1 = \bar{x} = 2.5$$

$$\nu_2 = \mu_2 + \bar{x}^2 = 14.75 + 2.5^2 = 21$$

$$\begin{aligned}\nu_3 &= \mu_3 + 3\mu_2\bar{x} + \bar{x}^3 \\ &= 166\end{aligned}$$

$$\nu_4 = \mu_4 + 4\mu_3\bar{x} + 6\mu_2\bar{x}^2 + \bar{x}^4 = 1132$$

Now Karl Pearson's coef.

$$\beta_1 = \frac{m_3^2}{m_2^3} = \frac{(39.75)^3}{(14.75)^3} = 0.492377$$

$$\beta_2 = \frac{m_4}{m_2^2} = \frac{(141.3195)}{(14.75)^2} = 0.654122$$

$$= 0.654122.$$

$$\gamma_1 = \sqrt{\beta_1} = \sqrt{0.492377}$$

$$\gamma_2 = \beta_2 - 3 = 0.654122 - 3$$

$$=$$

Moments about point $x=2$

$$m'_1 = \bar{x} - A = 2.5 - 2 = 0.5$$

$$\therefore m'_2 = m_2 - m'^2_1$$

$$m'_2 = m_2 + m'^2_1 = 14.75 + 0.5^2$$

$$= 15$$

$$m'_3 = m'_3 - 3m'_2 m'_1 + 2m'^3_1$$

$$m'_3 = m_3 + 3m'_2 m'_1 - 2m'^3_1$$

$$= 39.75 + 3 \times 15 \times 0.5 - 2 \times 0.5^3$$

$$= 62$$

$$m'_4 = m'_4 - 4m'_3 m'_1 + 6m'_2 m'^2_1 - 3m'^4_1$$

$$m'_4 = m_4 + 4m'_3 m'_1 - 6m'_2 m'^2_1 + 3m'^4_1$$

$$= 244$$

Ex:-2

The first three moments of a distribution about the value \bar{x} of the variable are 1, 16, -40 show that the mean is 3, variance is 15, + $\mu_3 = -8$.

$$A = 2,$$

$$\mu'_1 = 1, \mu'_2 = 16, \mu'_3 = -40$$

we have,

$$\mu'_1 = \bar{x} - A$$

$$\therefore \bar{x} = \mu'_1 + A = 1 + 2 = 3$$

$$\therefore \mu'_2 = \mu'_2 - \bar{x}^2 = 16 - 1^2 = 15$$

\therefore Variance $\mu_2 = 15$

$$\mu'_3 = \mu'_3 - 3\mu'_2\bar{x} + 2\bar{x}^3$$

$$= -40 - 3 \times 16 \times 1 + 2 \times 1^3$$

$$= -86.$$

Ex:-3

The first four moments of a distribution about $x=2$ are 1, 2.5, 5.5 and 16. Calculate the first four moments about the mean and about origin.

Ex:-

$$\text{we have } A = 2, \mu'_1 = 1, \mu'_2 = 2.5, \mu'_3 = 5.5$$

$$\mu'_4 = 16$$

Moments about mean.

$$\mu_1 = 0$$

$$\mu_2 = \mu'_2 - \mu'_1^2 = (2.5) - 1^2 = 2.5 - 1 = 1.5$$

$$\mu_3 = \mu'_3 - 3\mu'_2\mu'_1 + 2\mu'_1^3$$

$$= 5.5 - 3 \times 2.5 \times 1 + 2 \times 1^3$$

$$= 5.5 - 7.5 + 2 = 0$$

$$\mu_4 = \mu'_4 - 4\mu'_3\mu'_1 + 6\mu'_2\mu'_1^2 - 3\mu'_1^4$$

$$= 16 - 4 \cdot (5.5) \cdot (1) + 6 \times 2.5 \times 1^2 - 3 \times 1^4$$

$$= 6$$

Moments about Origin \rightarrow

$$\Rightarrow \frac{v_1}{\bar{x}} = \frac{\bar{x} - A}{v_1}$$

$$\therefore \mu'_1 = \bar{x} - A \\ \Rightarrow \bar{x} = \mu'_1 + A = 1 + 2 = 3$$

$$\text{So } v_1 = \bar{x} = 3$$

$$v_2 = \mu_2 + \bar{x}^2 \\ = 1.5 + 3^2 \\ = 1.5 + 9 \\ = 10.5$$

$$v_3 = \mu_3 + 3\mu_2\bar{x} + \bar{x}^3 \\ = 0 + 3 \times 1.5 \times 3 + 3^3 \\ = 40.5$$

$$v_4 = \mu_4 + 4\mu_3\bar{x} + 6\mu_2\bar{x}^2 + \bar{x}^4 \\ = 6 + 4(0)(3) + 6(1.5)(3)^2 + (3)^4 \\ = 168.$$

Ex-4.

for a distribution, the mean is 10, variance is 16, $\gamma_1 = 1$ & β_2 is 4. Find the first four moments about the origin.

Sol:-

Here $\bar{x} = 10, \mu_2 = 16, \gamma_1 = 1, \beta_2 = 4$

$$\beta_2 = \frac{\mu_4}{\mu_2^2} = 4$$

$$\Rightarrow \mu_4 = 4 \times 16^2 = 4 \times 256 = 1024$$

$$\text{also } \gamma_1 = 1 \Rightarrow +\sqrt{\beta_1} = 1 \Rightarrow \beta_1 = 1$$

$$\Rightarrow \frac{\mu_3^2}{\mu_2^3} = 1$$

$$\Rightarrow \mu_3^2 = 1 \times 16^3 \Rightarrow \mu_3^2 = (64)^2$$

$$(M_3 = 64)$$

Moments about the origin \rightarrow

$$V_1 = \bar{x} = 10, V_2 = 116, V_3 = 1544, V_4 = 22184$$

Skewness :-

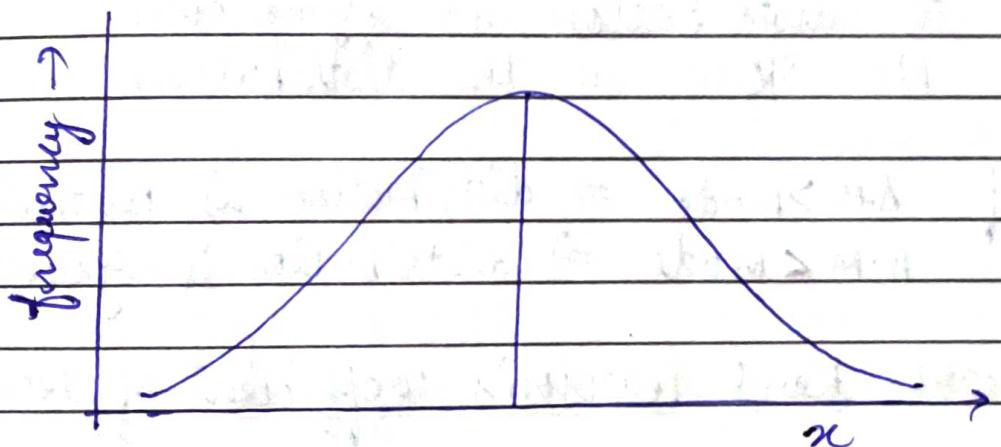
Symmetric distribution

For a symmetrical distribution, the frequencies are symmetrically distributed about the mean i.e. variates equidistant from the mean have equal frequencies.

In symmetric distribution,

$$\text{Mean} = \text{Mode} = \text{Median}$$

and median lies half-way between the two quartiles.



$$M = M_o = M_d$$

Skewness :- lack of symmetry in a frequency distribution is called skewness

Skewness indicates whether the curve is turned more to one side than to other side. i.e. whether the curve has a longer tail on one side.

$$(M_3 = 64)$$

Moments about the origin →

$$V_1 = \bar{x} = 10, V_2 = 116, V_3 = 1544, V_4 = 22184$$

Skewness :-

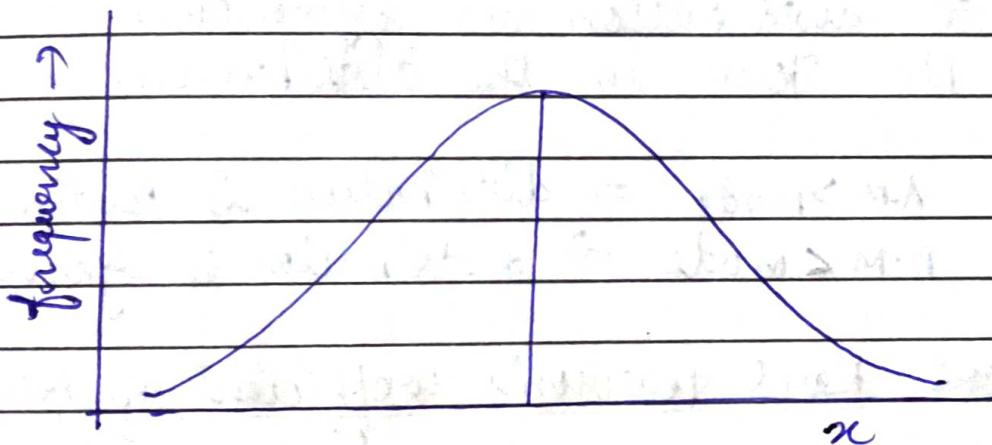
symmetric distribution

For a symmetrical distribution, the frequencies are symmetrically distributed about the mean i.e. variates equidistant from the mean have equal frequencies.

In symmetric distribution,

$$\text{Mean} = \text{Mode} = \text{Median}$$

and median lies half-way between the two quartiles.

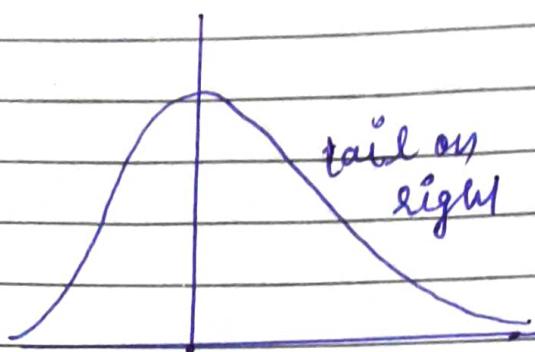


$$M = M_o = M_d$$

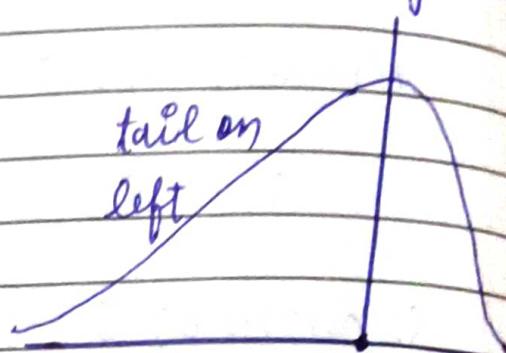
Skewness :- lack of symmetry in a frequency distribution is called skewness

Skewness indicates whether the curve is turned more to one side than to other side. i.e. whether the curve has a longer tail on one side.

Skewness can be positive as well as negative



Positively skewed distribution



Negatively skewed distribution

- * Skewness is positive if the longer tail of the distribution lies towards the right and negative if it lies towards the left.
- * If $A.M = \text{Mode} = \text{Median}$
 \Rightarrow distribution is symmetric
 \Rightarrow No skew in the distribution.
- * If $A.M > \text{Mode} \Rightarrow$ distribution is positively skewed.
 $A.M < \text{Mode} \Rightarrow$ distribution is negatively skewed
- * Karl Pearson's coefficient of skewness
 measure of skewness are called coef. of skewness
Karl Pearson's coefficient of skewness

$$SK_p = \frac{A.M - \text{Mode}}{S.D}$$

As we know that,

$$\text{Mode} = 3 \text{Median} - 2 \text{A.M.}$$

$$\therefore \text{coeff. of skewness} = \frac{\text{A.M} - \text{Mode}}{\text{S.D}}$$

$$= \frac{\text{A.M} - (3 \text{Median} - 2 \text{A.M})}{\text{S.D}}$$

$$= \frac{3 \text{A.M} - 3 \text{Median}}{\text{S.D}}$$

$$\therefore \text{Karl Pearson's coefficient of skewness} = \frac{3(\text{A.M} - \text{Median})}{\text{S.D}}$$

Note:- i) The value of SK_p lies between -1 and 1.

ii) for symmetric distribution. $SK_p = 0$

iii) if $SK_p > 0$, Then distribution is positively skewed.

iv) if $SK_p < 0$, Then distribution is negatively skewed.

Ex:-1. Karl Pearson's coefficient of skewness of a distribution is 0.32. Its standard deviation is 6.5, mean is 29.6 find the mode of the distribution.

soln:- we have $SK_p = 0.32$, $S.D = 6.5$
mean $\bar{x} = 29.6$

$$SK_p = \frac{\text{A.M} - \text{Mode}}{\text{S.D}}$$

$$0.32 = \frac{29.6 - \text{Mode}}{6.5}$$

$$29.6 - \text{Mode} = 0.32 \times 6.5$$

$$\text{Mode} = 27.52$$

Ex:-2. The sum of 20 observations is 300 and sum of their squares is 5000. The median is 15. find the Karl Pearson's coefficient of skewness-

Sol:-

we have,

$$n = 20, \quad \sum x = 300, \quad x \text{ is a variable}$$

$$\sum x^2 = 5000$$

$$\text{Median} = 15$$

$$\text{Mean } \bar{x} = \frac{\sum x}{n} = \frac{300}{20} = 15$$

$$\begin{aligned}\text{Variance} &= \mu_2 = \bar{x}^2 - \frac{\sum x^2}{n} \\ &= \bar{x}^2 - \frac{5000}{20} \\ &= 15^2 - 250 \\ &= 225 - 25 \\ &= 200\end{aligned}$$

$$\begin{aligned}\text{So } S.D. &= \sqrt{\mu_2} \\ &= \sqrt{200} \\ &= 10\sqrt{2}\end{aligned}$$

now Karl Pearson's coeff. of skewness

$$\begin{aligned}SK_p &= \frac{3(A.M - \text{median})}{S.D.} \\ &= \frac{3(15 - 15)}{10\sqrt{2}} \\ &= \frac{3 \times 0}{10\sqrt{2}} \\ &= 0.\end{aligned}$$

∴ distribution is symmetric.

Ex:- 3 find the coefficient of skewness by Karl Pearson's method for the following data:

value	6	12	18	24	30	36	42
Frequency	4	7	9	18	15	10	3

Soln:- Calculation of \bar{x} , S.D

value x	f	fx	fx^2
6	4	24	144
12	7	84	1008
18	9	162	2916
24	18	432	10368
30	15	450	13500
36	10	360	12960
42	3	126	5292
$N = 66$		$\sum fx = 1638$	

$$\text{A.M } \bar{x} = \frac{\sum fx}{\sum f} = \frac{\sum fy}{N} = \frac{1638}{66} \\ = 24.82$$

$$\begin{aligned} \text{Standard deviation, S.D} &= \sqrt{\mu_2} \\ &= \sqrt{\nu_2 - \bar{x}^2} \quad \therefore \nu_2 = \mu_2 + \bar{x}^2 \\ &= \sqrt{\frac{\sum fx^2}{N} - \left(\frac{\sum fy}{N}\right)^2} \\ &= \sqrt{699.82 - 24.82^2} \\ &= \sqrt{83.7876} \\ &= 9.15 \end{aligned}$$

Mode = value of x corresponding to maximum freq. 18
 $= 24$

So coefficient of skewness

$$SK_p = \frac{A - M - \text{Mode}}{S.D}$$

$$= \frac{24 - 82 - 24}{9.15}$$

$$= 0.0896$$

Method of Moments \rightarrow

In this Method, Second and third central moments of the distribution are used. This measure of skewness is called the moment coefficient of skewness, denoted by SK_m or r_1 .

Moment coefficient of skewness = $\frac{\mu_3}{\sqrt{\mu_2^3}}$

$$SK_m = r_1 = \sqrt{\beta_1} = \sqrt{\frac{\mu_3^2}{\mu_2^3}}$$

so

$$SK_m = \frac{\mu_3}{\sqrt{\mu_2^3}}$$

Ex:- The first three central moments of a distribution are 0, 15, -31. find the moment coef. of skewness.

we have $\mu_1 = 0$, $\mu_2 = 15$, $\mu_3 = -31$

$$\text{so moment coef. of skewness} = \frac{\mu_3}{\sqrt{\mu_2^2}} = \frac{-31}{\sqrt{15^3}} = \frac{-31}{58.09} = -0.53$$

Ex:-2. The first four moments of a distribution about the value 2 of the variable are 2, 20, 40 and 50, calculate the moment coeff. of skewness.

Soln:-

We have,

$$A = 2, \mu'_1 = 2, \mu'_2 = 20, \mu'_3 = 40, \mu'_4 = 50$$

$$\text{Now, } \mu_2 = \mu'_2 - \mu'_1^2 \\ = 20 - 2^2 \\ = 16$$

$$\mu_3 M_3 = \mu'_3 - 3\mu'_2 \mu'_1 + 2\mu'_1^3 \\ = 40 + 3 \times 20 \times 2 + 2 \times 2^3 \\ = 40 - 120 + 16 \\ = 56 - 120 \\ = -64$$

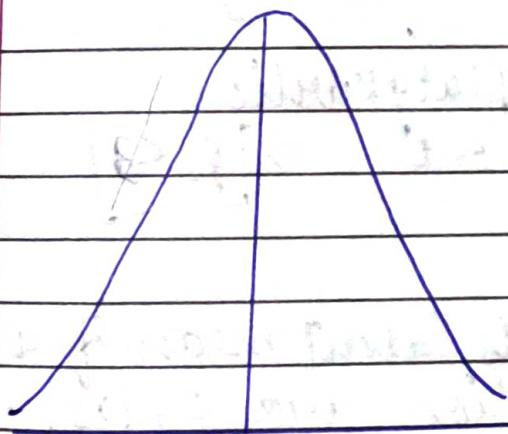
Moment coefficient of skewness

$$SK_M = \frac{\mu_3}{\sqrt{\mu_2^3}} \\ = \frac{-64}{\sqrt{16^3}} \\ = \frac{-64}{64} \\ = -1$$

Kurtosis :- \rightarrow A frequency curve may be symmetrical but it may not be equally flat topped with the normal curve. The bulginess or the relative flatness of the curve of a freq. distribution is called Kurtosis.

Hence. The sharpness of the peak of a frequency-distribution curve is called Kurtosis.

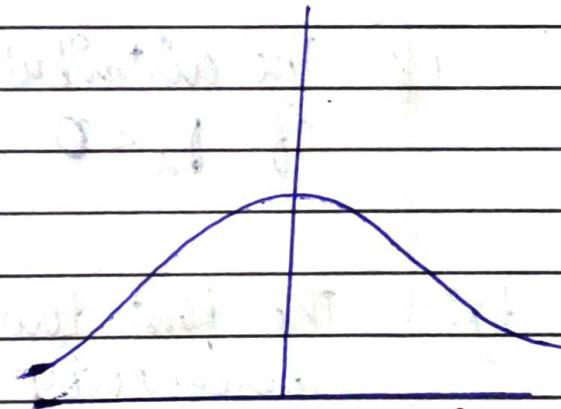
- ① Curves which are neither flat nor sharply peaked are called normal curves or mesokurtic curves.
- ② Curves which are flatter than the normal curve are called Platykurtic curves.
- ③ Curves which are more sharply peaked than the normal curve are called Leptokurtic curves.



Leptokurtic

$$\beta_2 > 3$$

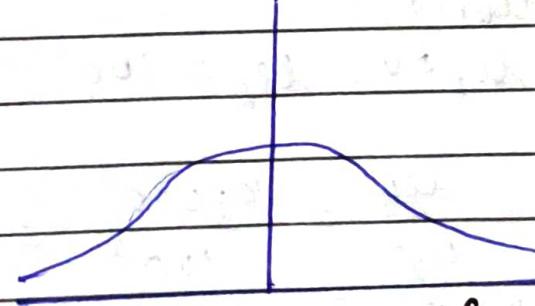
$$\gamma_2 > 0$$



Mesokurtic

$$\beta_2 = 3$$

$$\gamma_2 = 0$$



Platykurtic.

$$\beta_2 < 3$$

$$\gamma_2 \leq 0$$

Measure of Kurtosis :-

Measure of kurtosis is denoted by β_2 and is defined by

$$\beta_2 = \frac{\mu_4}{\mu_2^2}$$

μ_2 and μ_4 are second & fourth Central Moment.

The kurtosis of a distribution is also measured by $\gamma_2 = \beta_2 - 3$.

Note :- ① The distribution is leptokurtic

$$\text{if } \gamma_2 > 0 \Rightarrow \beta_2 - 3 > 0 \Rightarrow [\beta_2 > 3]$$

② The distribution is mesokurtic if ~~$\beta_2 = 3$~~

$$\gamma_2 = 0 \Rightarrow \beta_2 - 3 = 0 \Rightarrow [\beta_2 = 3]$$

③ The distribution is platykurtic

$$\text{if } \gamma_2 < 0 \Rightarrow \beta_2 - 3 < 0 \Rightarrow [\beta_2 < 3]$$

Ex :- The first four moments about mean of a frequency distribution are 0, 100, -7 and 35000. Discuss the kurtosis of the distribution.

Sol :- Given $\mu_1 = 0$, $\mu_2 = 100$, $\mu_3 = -7$, $\mu_4 = 35000$

$$\text{Now } \beta_2 = \frac{\mu_4}{\mu_2^2} = \frac{35000}{100^2} = 3.5 > 3$$

\Rightarrow The distribution is leptokurtic.

Ques. The first four moments of a distribution about $x=4$ are 1, 4, 10 and 45. Obtain the various characteristics of the distribution on the basis of the given information. Comment upon the nature of the distribution. (2018, 2018)

Soln- we have $A = 4$

$$\mu'_1 = 1, \mu'_2 = 4, \mu'_3 = 10, \mu'_4 = 45$$

Moments about mean (Central moments)

$$\mu_1 = 0$$

$$\mu_2 = \mu'_2 - \mu'_1{}^2 = 4 - 1^2 = 3$$

$$\begin{aligned}\mu_3 &= \mu'_3 - 3\mu'_2\mu'_1 + 2\mu'_1{}^3 \\ &= 10 - 3 \times 4 \times 1 + 2 \times 1^3 = 0\end{aligned}$$

$$\begin{aligned}\mu_4 &= \mu'_4 - 4\mu'_3\mu'_1 + 6\mu'_2\mu'_1{}^2 - 3\mu'_1{}^4 \\ &= 45 - 4 \times 10 \times 1 + 6 \times 4 \times 1^2 - 3 \times 1^4 \\ &= 45 - 40 + 24 - 3 \\ &= 26\end{aligned}$$

SKEWNESS:

$$\text{Moment coefficient of skewness} = \frac{\mu_3}{\sqrt{\mu_2^3}}$$

$$= \frac{0}{\sqrt{3^3}} = 0$$

\therefore The distribution is symmetrical

KURTOSIS:

$$\beta_2 = \frac{\mu_4}{\mu_2^2} = \frac{26}{3^2} = \frac{26}{9} = 2.8888$$

\therefore The distribution is platykurtic.

Ex:-3. The following table represents the height of a batch of 100 students. calculate Kurtosis & discuss skewness

height (cm)	59	61	63	65	67	69	71	73	75
No. of students	0	2	6	20	40	20	8	2	2

(2018)

Soln-

Height (in cm)	No. of Students	$u = \frac{x-A}{2}$ $A=67$	fu	fu^2	fu^3	fu^4
59	0	-4	0	0	0	0
61	2	-3	-6	18	-54	162
63	6	-2	-12	24	-48	96
65	20	-1	-20	20	-20	20
67	40	0	0	0	0	0
69	20	1	20	20	20	20
71	8	2	16	32	64	128
73	2	3	6	18	54	162
75	2	4	8	32	128	512
$\sum f = 100$			$\sum fu = 12$	$\sum fu^2 = 164$	$\sum fu^3 = 144$	$\sum fu^4 = 1100$

Moments about the point $A=67 \rightarrow$

$$M_1' = \frac{\sum fu \times h}{N}$$

$$= \frac{12}{100} \times 2 = 0.24$$

$$M_2' = \frac{\sum fu^2 \times h^2}{N}$$

$$= \frac{164}{100} \times 2^2 = 6.56$$

$$M_3' = \frac{\sum fu^3 \times h^3}{N} = \frac{144}{100} \times 2^3 = 11.52$$

$$M_4' = \frac{\sum f u^4 \times h^4}{N} - \frac{1100}{100} \times 2^4 = 176$$

Moments about mean \rightarrow

$$U_1' = 0$$

$$M_2 = U_2' - U_1'^2 = 6.56 - (0.24)^2 = 6.5024$$

$$\begin{aligned} M_3 &= M_3' - 3M_2'U_1' + 2U_1'^3 \\ &= 11.52 - 3 \times 6.56 \times 0.24 + 2 \times 0.24^3 \\ &= 6.8244 \end{aligned}$$

$$\begin{aligned} M_4 &= M_4' - 4M_3'U_1' + 6M_2'U_1'^2 - 3U_1'^4 \\ &= 167.19798 \end{aligned}$$

Skewness :-

$$\text{Moment coef. of Skewness } SK_m = \frac{M_3}{\sqrt{M_2^3}} = \frac{6.8244}{\sqrt{6.5024^3}} = 0.411570$$

\therefore The distribution is positively skewed

Kurtosis :-

$$\beta_2 = \frac{M_4}{M_2^2} = \frac{167.19798}{(6.5024)^2} = 3.9544 > 3$$

$$\Rightarrow \beta_2 > 3$$

Hence the distribution is leptokurtic.

Ques. The first four moments about the value x of a distribution are 0.294, 7.044, 42.409 and 454.98 calculate the moments about mean. Also evaluate β_1 , β_2 & comment upon the skewness & kurtosis of the distribution.

(#)

Curve fitting :-

Consider n - paired observations $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ of two variables x and y .

In order to have an approximate idea about the relationship of these two variables,

we plot these n paired points on a graph, we get a diagram showing the simultaneous variation in values of both the variables called Scatter or dot diagram.

Curve fitting means an exact relationship of two variables by Algebraic equations. In fact this relationship is the eqⁿ of the curve. Therefore, curve fitting means to form an eqⁿ of the curve from the given data.

It is useful in the study of Correlation and Regression.

for this we use Method of least squares

Method of least squares:-

Method of least squares provides a unique set of values to the constants and hence suggests a curve of best fit to the given data.

Suppose we have n - paired observations $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ of x & y .

here $y_i^o \rightarrow$ observed value
or given value for x_i

$y_i^e \rightarrow$ Expected value

$$(y_i^e = f(x_i) \\ i=1, 2, \dots, n)$$

Error or residual

$$E_i^2 = y_i^o - y_i^e \quad (\text{may be +ve / -ve})$$

$$\Rightarrow E^2 = (y_i^o - y_i^e)^2 \quad (\text{for equal weightage})$$

Now introducing $E = \text{Total error}$

$$E = \sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i^o - y_i^e)^2 = \sum_{i=1}^n (y_i^o - f(x_i))^2$$

[if $E=0 \Rightarrow E_i=0 \quad \forall i=1, \dots, n$, then all the points lie on the curve]

The least value of E gives the best fitting curve to the data.

This method is called method of least squares
or principle of least square?

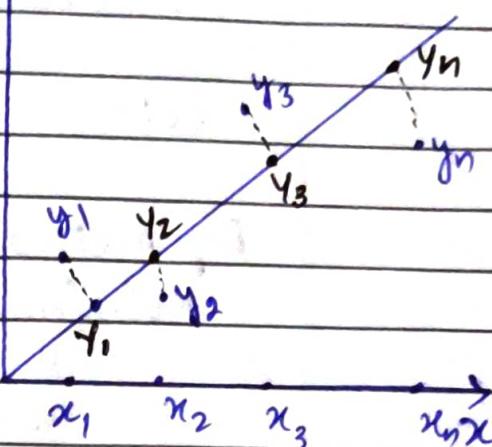
I) fitting of a straight line :-

Let $(x_i, y_i), i=1, 2, \dots, n$ be n sets of observat^ys
of related data and

$$y = a + bx \quad \text{--- ①}$$

be the straight line to be fitted.
now residual / errors

$$E_i^2 = y_i^o - y_i^e$$



$$\Rightarrow E_i = y_i - a - bx_i$$

$$\therefore Y_i = a + bx_i$$

$$\text{Let } U = \sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i - a - bx_i)^2 \quad (2)$$

By method of least squares, U should be minimum.

$$\therefore \frac{\partial U}{\partial a} = 0 \quad \text{and} \quad \frac{\partial U}{\partial b} = 0$$

$$\therefore \frac{\partial U}{\partial a} = 2 \sum_{i=1}^n (y_i - a - bx_i)(-1) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - a \sum_{i=1}^n 1 - b \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \sum y - an - b \sum x = 0$$

$$\Rightarrow \boxed{\sum y = an + b \sum x} \quad (3)$$

$$\text{and } \frac{\partial U}{\partial b} = 2 \sum_{i=1}^n (y_i - a - bx_i)(-x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 = 0$$

$$\Rightarrow \boxed{\sum xy = a \sum x + b \sum x^2} \quad (4)$$

$\therefore U$ is minimum when

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

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

These eqns are called Normal equation for a & b .

on solving (3), (4) we get a & b

but these values of a & b in eqn (1)

$$y = a + bx$$

we get best fit straight line to the given data.

Expt 1 By method of least squares, find the straight line that best fits the following data:

$x :$	1	2	3	4	5
$y :$	14	27	40	55	68

Soln:-

Let the straight line of best fit be

$$y = a + bx \quad \text{--- (1)}$$

Normal eq's are

$$\sum y = na + b \sum x \quad \text{--- (2)}$$

$$\text{and} \quad \sum xy = a \sum x + b \sum x^2 \quad \text{--- (3)}$$

Here $n = 5$.

x	y	xy	x^2
1	14	14	1
2	27	54	4
3	40	120	9
4	55	220	16
5	68	340	25

$$\sum x = 15 \quad \sum y = 204 \quad \sum xy = 748 \quad \sum x^2 = 55$$

Substituting in (2) & (3),

$$204 = 5a + 15b$$

$$748 = 15a + 55b$$

on solving,

$$15a + 45b = 612$$

$$15a + 55b = 748$$

$$-10b = -136$$

$$b = 13.6, \quad a = \frac{204 - 15 \times 13.6}{5} = 0$$

so $a = 0, b = 13.6$
Hence required straight line is

$$y = 13.6x$$

Ex:- 2. fit a straight line to the following data by least square Method.

$x:$	0	1	2	3	y
$y:$	1	1.8	3.3	4.5	6.3

Sol:-

Let the straight line to be fitted in the given data

$$y = a + bx \quad \text{--- } ①$$

Normal eqns are

$$\sum y = an + b \sum x \quad \text{--- } ②$$

$$\sum xy = a \sum x + b \sum x^2 \quad \text{--- } ③$$

Here $n = 5$

x	y	xy	x^2
0	1	0	0
1	1.8	1.8	1
2	3.3	6.6	4
3	4.5	13.5	9
4	6.3	25.2	16
$\sum x = 10$		$\sum y = 16.9$	$\sum xy = 47.1$
			$\sum x^2 = 30$

from ① & ②

$$16.9 = 5a + 10b$$

$$47.1 = 10a + 30b$$

on solving

$$\begin{array}{rcl} 10a + 20b & = & 33.8 \\ 10a + 30b & = & 47.1 \end{array}$$

$$-10b = -13.3$$

$$b = 1.33$$

$$a = 0.72$$

Hence Required straight line is

$$y = 0.72 + 1.33x$$

(ii) fitting a second degree parabola :-

Let (x_i, y_i) , $i=1, 2, \dots, n$ be the given data and

$$y = a + bx + cx^2 \quad \text{--- (1)}$$

be the second degree parabola to be fitted.
Now Errors,

$$E_i = y_i - Y_i \quad [\because Y_i = a + bx_i + cx_i^2]$$

$$E_i = y_i - a - bx_i - cx_i^2 \quad \text{--- (2)}$$

$$\text{let } V = \sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2 \quad \text{--- (3)}$$

By the method of least square, V should be minimum.

For min. value of V

$$\frac{\partial V}{\partial a} = 0, \quad \frac{\partial V}{\partial b} = 0, \quad \frac{\partial V}{\partial c} = 0$$

$$\frac{\partial U}{\partial a} = 2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)(-1) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - a \sum_{i=1}^n 1 - b \sum_{i=1}^n x_i - c \sum_{i=1}^n x_i^2 = 0$$

$$\Rightarrow \boxed{\sum y = a n + b \sum x + c \sum x^2} \quad \textcircled{a}$$

$$\frac{\partial U}{\partial b} = 2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)(-x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 - c \sum_{i=1}^n x_i^3 = 0$$

$$\Rightarrow \boxed{\sum xy = a \sum x + b \sum x^2 + c \sum x^3} \quad \textcircled{b}$$

$$\frac{\partial U}{\partial c} = 2 \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)(-2x_i^2) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i^2 y - a \sum_{i=1}^n x_i^2 - b \sum_{i=1}^n x_i^3 - c \sum_{i=1}^n x_i^4 = 0$$

$$\Rightarrow \boxed{\sum x^2 y = a \sum x^2 + b \sum x^3 + c \sum x^4} \quad \textcircled{c}$$

eqns \textcircled{a} , \textcircled{b} and \textcircled{c} are called Normal eqns. solving these for a , b , c
put these values in eqn $y = a + bx + cx^2$

we get best fitted parabola.

Ex:- fit a second degree parabola in the following data:

x :	0	1	2	3	4
y :	1	4	10	17	30

Sol:- Let $y = a + bx + cx^2$ be second degree parabola to be fitted in the given data, normal eq's are

$$\left. \begin{array}{l} \sum y = a n + b \sum x + c \sum x^2 \\ \sum xy = a \sum x + b \sum x^2 + c \sum x^3 \\ \sum x^2 y = a \sum x^2 + b \sum x^3 + c \sum x^4 \end{array} \right\} -②$$

x	y	x^2	x^3	x^4	xy	x^2y
0	1	0	0	0	0	0
1	4	1	1	1	4	4
2	10	4	8	16	20	40
3	17	9	27	81	51	153
4	30	16	64	256	120	480
$\Sigma x = 10$	$\Sigma y = 62$	$\Sigma x^2 = 30$	$\Sigma x^3 = 100$	$\Sigma x^4 = 354$	$\Sigma xy = 195$	$\Sigma x^2y = 677$

using all these values in ②

$$62 = 5a + 10b + 30c \quad -④$$

$$195 = 10a + 30b + 100c \quad -⑤$$

$$677 = 30a + 100b + 354c \quad -⑥$$

on solving ④, ⑤ & ⑥, we have
 $a = 1.2$, $b = 1.1$, $c = 1.5$

Putting all these values in eqn ①

$y = 1.2 + 1.1x + 1.5x^2$
which is the required eqn of Parabola
to be fitted in the given data.

(iii)

Fitting of the Curve: \rightarrow $y = ax + bx^2$

Let (x_i, y_i) , $i=1, 2, \dots, n$ be the given data
and $y = ax + bx^2$ — (1)
be the curve fitted to the data.

Residual at $x = x_i$:

$$E_i = y_i - y_i$$

$$[\because y_i = ax_i + bx_i^2]$$

$$\Rightarrow E_i = y_i - ax_i - bx_i^2$$

Let

$$U = \sum_{i=1}^n E_i^2 = \sum_{i=1}^n (y_i - ax_i - bx_i^2)^2 — (2)$$

Now by Method of least square, U should be minimum.

for min. value U ,

$$\frac{\partial U}{\partial a} = 0, \quad \frac{\partial U}{\partial b} = 0$$

$$\frac{\partial U}{\partial a} = 0 \Rightarrow 2 \sum_{i=1}^n (y_i - ax_i - bx_i^2)(-x_i) = 0$$

$$\Rightarrow \sum_{i=1}^n (x_i y_i - ax_i^2 - bx_i^3) = 0$$

$$\Rightarrow \sum_{i=1}^n x_i y_i = a \sum_{i=1}^n x_i^2 + b \sum_{i=1}^n x_i^3$$

$$\Rightarrow \boxed{\sum xy = a \sum x^2 + b \sum x^3} — (3)$$

$$\frac{\partial U}{\partial b} = 0 \Rightarrow 2 \sum_{i=1}^n (y_i - ax_i - bx_i^2)(-x_i^2) = 0$$

$$\Rightarrow \sum_{i=1}^n (x_i^2 y_i - ax_i^3 - bx_i^4) = 0$$

$$\Rightarrow \boxed{\sum x^2 y = a \sum x^3 + b \sum x^4} — (4)$$

eqn (3) and (4) are Normal equations.

Ex:- By method of least square, find the curve $y = ax + bx^2$ that best fit the following data:

(2014)

x :	1	2	3	4	5
y :	1.8	5.1	8.9	14.1	19.8

Sol:- Let $y = ax + bx^2$ be the curve that best fits to the given data.

Here Normal equations are

$$\begin{aligned} \sum xy &= a \sum x^2 + b \sum x^3 \\ \sum x^2 y &= a \sum x^3 + b \sum x^4 \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} - \textcircled{2}$$

Here $n = 5$.

x	y	x^2	x^3	x^4	xy	x^2y
1	1.8	1	1	1	1.8	1.8
2	5.1	4	8	16	10.2	20.4
3	8.9	9	27	81	26.7	80.1
4	14.1	16	64	256	56.4	225.6
5	19.8	25	125	625	99	495
		$\sum x^2 = 55$	$\sum x^3 = 225$	$\sum x^4 = 979$	$\sum xy = 194.1$	$\sum x^2y = 822.9$

Substituting these values in equation ~~②~~ we get

$$194.1 = 55a + 225b$$

$$822.9 = 225a + 979b$$

on solving, we get

$$a = \frac{83.85}{55} = 1.52$$

$$b = \frac{317.4}{664} = 0.49$$

Hence Required eqn of curve

$$y = 1.52x + 0.49x^2$$

(IV)

Fitting of an exponential curve $y = ae^{bx}$

$$\text{Given } y = ae^{bx} \quad \dots \text{ (1)}$$

Taking logarithm on both sides of eqⁿ (1)

$$\log_{10} y = \log_{10}(a \cdot e^{bx}) = \log_{10} a + \log_{10} e^{bx}$$

$$\log_{10} y = \log_{10} a + bx \log_{10} e$$

$$\Rightarrow \log_{10} y = \log_{10} a + (b \cdot \log_{10} e) \cdot x$$

$$\boxed{Y = A + BX}$$

where $A = \log_{10} a$, $B = b \cdot \log_{10} e$, $X = x$

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

Normal equations are -

$$\Sigma Y = A \cdot n + B \sum X$$

$$\text{and } \Sigma XY = A \sum X + B \sum X^2$$

Solving on these, we get A and B

$$\text{Then } \boxed{a = \text{Antilog } A}$$

$$\boxed{b = \frac{B}{\log_{10} e}}$$

Ex:- Find the curve of best fit of the type $y = ae^{bx}$ to the following data by the method of least squares.

$$x: 1 \quad 5 \quad 7 \quad 9 \quad 12$$

$$y: 10 \quad 15 \quad 12 \quad 15 \quad 21$$

Soln:- The curve to be fitted is $y = ae^{bx}$ — (1)

or $Y = A + BX$ — (2)

where $Y = \log_{10} y$, $A = \log_{10} a$, $B = b \log_{10} e$, $X = x$

\therefore The normal equations are

$$\sum y = A \cdot n + B \sum x$$

$$\sum xy = A \sum x + B \sum x^2$$

Here $n = 5$

— (3)
— (4)

$x = x$	y	$y = \log_{10} y$	x^2	xy
1	10	1	1	1
5	15	1.1761	25	5.8805
7	12	1.0792	49	7.5544
9	15	1.1761	81	10.5849
12	21	1.3222	144	15.8664
$\sum x = 34$		$\sum y = 5.7536$	$\sum x^2 = 300$	$\sum xy = 40.8862$

Substituting these values in eqn (3) & (4), we get

$$5.7536 = 5A + 34B$$

$$40.8862 = 34A + 300B$$

on solving, we get

$$A = 0.9766, B = 0.02561$$

\therefore

$$a = \text{Antilog}_{10} A, b = \frac{B}{\log_{10} e} = \frac{0.02561}{\log_{10} e}$$

$$\Rightarrow a = 9.4754, b = 0.059.$$

Hence The required curve is $y = 9.4754 e^{0.059x}$

(v)

Fitting of the curve $y = ax^b$.

$$\text{given } y = ax^b$$

Taking logarithm on both sides, we get

$$\log_{10} y = \log_{10} a + b \log_{10} x$$

$$\text{or } y = A + BX$$

$$\text{where } y = \log_{10} y, A = \log_{10} a, B = b, X = \log_{10} x$$

Normal eqns are

$$\sum y = A \cdot n + B \sum x$$

$$\sum xy = A \sum x + B \sum x^2$$

on solving we get A and B.

$$a = \text{Antilog } A, \quad b = B$$

(7)

Fitting of the curve $y = ab^x$

given curve is $y = ab^x$ —①

taking logarithm on both sides

$$\log y = \log_{10}(ab^x)$$

$$\log_{10} y = \log_{10} a + x \log_{10} b$$

or

$$y = A + BX$$

—②

where

$$y = \log_{10} y, \quad A = \log_{10} a, \quad B = \log_{10} b, \quad x = z$$

Normal eqn are

$$\sum y = A \cdot n + B \sum x$$

$$\sum xy = A \sum x + B \sum x^2$$

where n is the no. of pairs of values of x and y.

After solving Normal eqn's we get A and B
and hence

$$a = \text{Antilog}_{10} A, \quad b = \text{Antilog}_{10} B$$

Ex:-

Obtain a relation of the form $y = ab^x$ for the following data by the Method of least squares

x	2	3	4	5	6
y	8.3	15.4	33.1	65.9	127.4

SOLN-

The curve to be fitted is $y = ab^x$ —①
 taking log on both sides

$$\log_{10} y = \log_{10} a + x \log_{10} b$$

$$\text{or } y = A + BX \quad \text{—②}$$

where

$$Y = \log_{10} y, A = \log_{10} a, B = \log_{10} b, X = x$$

Normal equations are -

$$\sum Y = A \cdot n + B \sum X$$

$$\sum XY = A \sum X + B \sum X^2, n=5$$

$X=x$	y	$Y = \log_{10} y$	X^2	XY
2	8.3	0.9191	4	1.8382
3	15.4	1.1872	9	3.5616
4	33.1	1.5198	16	6.0792
5	65.2	1.8142	25	9.0710
6	127.4	2.1052	36	12.6312
$\sum X = 20$		$\sum Y = 7.5455$	$\sum X^2 = 90$	$\sum XY = 33.1812$

Substituting these values in Normal equations

$$7.5455 = 5A + 20B$$

$$33.1812 = 20A + 90B$$

on solving we get,

$$A = 0.31, B = 0.3$$

$$a = \text{antilog } A = \text{antilog } 0.31 = 2.04$$

$$b = \text{antilog } B = \text{antilog } 0.3 = 1.995$$

Hence the required curve is

$$y = 2.04 (1.995)^x$$

III

Fitting of the curve $PV^r = C$:-

Given $PV^r = C \Rightarrow V^r = \frac{C}{P} \Rightarrow P = CV^r$
Taking logarithm on both sides

$$\log P = \log C - r \log V$$

$$\text{or } Y = A + BX$$

where $Y = \log P$, $A = \log C$, $B = -r$, $X = \log V$
Normal eqns are

$$\sum Y = A \cdot n + B \sum X$$

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

Imp

Ex:- The pressure of the gas corresponding to various volumes V is measured given by the following data

$$V (\text{cm}^3) : 50 \quad 60 \quad 70 \quad 90 \quad 100$$

$$P (\text{kg cm}^{-2}) : 64.7 \quad 51.3 \quad 40.5 \quad 25.9 \quad 78$$

Sol:- fit the data to the eqn $PV^r = C$ (2019)

Given curve $PV^r = C$

$$\Rightarrow P = CV^{-r}$$

Taking logarithm on both sides, we get

$$\log P = \log C - r \log V$$

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

where

$$Y = \log P, A = \log C, B = -r, X = \log V$$

Normal eqns are

$$\sum Y = A \cdot n + B \sum X$$

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

Here $n = 5$.

V	P	X = log V	Y = log P	X ²	XY
50	64.7	1.69897	1.81090	2.88650	3.07666
60	51.3	1.777815	1.71012	3.16182	3.04085
70	40.5	1.84510	1.60746	3.40439	2.96592
90	25.9	1.95424	1.41330	3.81905	2.76193
100	78	2	1.89209	4	3.78418
		$\Sigma X = 9.27646$	$\Sigma Y = 8.43387$	$\Sigma XY = 15.62954$	$\Sigma X^2 = 15.62954$
					$\Sigma X^2 = 17.27176$

Putting these values in Normal equations,

$$8.43387 = 5A + 9.27646B$$

$$15.62954 = 9.27646A + 17.27176B$$

Solving these eqns, we get

$$A = 2.22476, \quad B = -0.28997$$

$$\therefore Y = -B = 0.28997$$

$$C = \text{antilog } A = \text{antilog } 2.22476 \\ = 167.78765$$

Hence, the required eqn of curve is

$$PV^{0.28997} = 167.78765$$

Fitting of the curve $y = \frac{C_0}{x} + C_1 \sqrt{x}$

Given curve to be fitted $y = \frac{C_0}{x} + C_1 \sqrt{x}$ ————— (1)

Error at point $x = x_i$ is

$$E_i = y_i - Y_i = y_i - \frac{C_0}{x_i} - C_1 \sqrt{x_i}$$

By method of least squares,

the values of C_0, C_1 are such that

$$O = \sum_{i=1}^n E_i^2 = \sum_{i=1}^n \left(y_i - \frac{C_0}{x_i} - C_1 \sqrt{x_i} \right)^2 \text{ is Minimum}$$

so normal equations are given by

$$\frac{\partial U}{\partial c_0} = 0 \quad \text{and} \quad \frac{\partial U}{\partial c_1} = 0$$

Now $\frac{\partial U}{\partial c_0} = 0$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - \frac{c_0}{x_i} - c_1 \sqrt{x_i}) \left(-\frac{1}{x_i} \right) = 0$$

$$\Rightarrow \sum_{i=1}^n \left(\frac{y_i}{x_i} - \frac{c_0}{x_i^2} - \frac{c_1}{\sqrt{x_i}} \right) = 0$$

$$\Rightarrow \boxed{\sum \frac{y}{x} = c_0 \sum \frac{1}{x^2} + c_1 \sum \frac{1}{\sqrt{x}}}$$

and $\frac{\partial U}{\partial c_1} = 0$

$$\Rightarrow 2 \sum_{i=1}^n (y_i - \frac{c_0}{x_i} - c_1 \sqrt{x_i}) (-\sqrt{x_i}) = 0$$

$$\Rightarrow \sum_{i=1}^n (y_i \sqrt{x_i} - \frac{c_0}{\sqrt{x_i}} - c_1 x_i) = 0$$

$$\Rightarrow \boxed{\sum y \sqrt{x} = c_0 \sum \frac{1}{\sqrt{x}} + c_1 \sum x}$$

so normal eqs are

$$\sum \frac{y}{x} = c_0 \sum \frac{1}{x^2} + c_1 \sum \frac{1}{\sqrt{x}}$$

$$\text{and } \sum y \sqrt{x} = c_0 \sum \frac{1}{\sqrt{x}} + c_1 \sum x$$

Ex:

Use Method of Least squares to fit the curve

$y = \frac{c_0}{x} + c_1 \sqrt{x}$ to the following table of values

$x : 0.1 \quad 0.2 \quad 0.4 \quad 0.5 \quad 1 \quad 2$

$y : 21 \quad 11 \quad 7 \quad 6 \quad 5 \quad 6$

Soln:-

The given curve to be fitted is $y = \frac{c_0}{x} + c_1 \sqrt{x}$ — (1)

The Normal equations are

$$\sum \frac{y}{x} = c_0 \sum \frac{1}{x^2} + c_1 \sum \frac{1}{\sqrt{x}} \quad \text{--- (2)}$$

and $\sum y \sqrt{x} = c_0 \sum \frac{1}{\sqrt{x}} + c_1 \sum x \quad \text{--- (3)}$

x	y	y/x	$y \sqrt{x}$	$1/\sqrt{x}$	y/x^2
0.1	21	210	6.64078	3.16228	100
0.2	11	55	4.91935	2.23607	25
0.4	7	17.5	4.42719	1.58114	6.25
0.5	6	12	4.24264	1.41421	4
1	5	5	5	1	1
2	6	3	8.48528	0.70711	0.25
$\Sigma x = 4.2$		$\sum \frac{y}{x} = 302.5$	$\sum y \sqrt{x} = 33.715$	$\sum \frac{1}{\sqrt{x}} = 10.10081$	$\sum \frac{1}{x^2} = 136.5$

from eqn (2) and (3), we have

$$302.5 = 136.5 c_0 + 10.10081 c_1$$

$$33.71524 = 10.10081 c_0 + 4.2 c_1$$

Solving these, we get

$$c_0 = 1.97327, \quad c_1 = 3.28182$$

Hence the required eqn of curve is

$$y = \frac{1.97327}{x} + 3.28182 x$$

(X)

Fitting of the curve $xy = b + ax = ax + b$

$$xy = b + ax$$

$$y = \frac{b}{x} + a$$

$$y = bx + a, \quad x = \frac{1}{x}, \quad y = y$$

Normal eqns are

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

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

(Y)

Fitting of the curve $y = ax + \frac{b}{x}$

Normal eqns are

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

$$\text{and } \sum \frac{y}{x} = a \cdot n + b \sum \frac{1}{x^2}$$

⑪ Correlation : →

In a Bivariate distribution, if the change in one variable affects a change in other variable, the variables are said to be correlated.

- if the two variables deviate in the same direction i.e. if the increase (or decrease) in one results in a corresponding increase (decrease) in the other, correlation is said to be direct or positive correlation.
e.g. the correlation between income and expenditure is positive.
- if the two variables deviate in opposite direction i.e. if the increase (or decrease) in one results in a corresponding decrease (or increase) in the other, correlation is said to be inverse or Negative correlation.
e.g. The correlation between volume and the pressure of a perfect gas or the correlation b/w price and demand is Negative.
- correlation is said to be Perfect if the deviation in one variable is followed by a corresponding proportional deviation in the other.

(1)

Measure of correlation

Karl Pearson's coefficient of correlation
(Product Moment Correlation Coefficient) :-

Consider two variables x and y then
we know

$$\text{Var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \sigma_x^2$$

$$\text{Var}(y) = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2 = \sigma_y^2$$

σ_x is S.D for x
 σ_y is S.D for y

$$\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

Karl Pearson's coefficient of correlation b/w two variables x and y , denoted by $r(x, y)$ or r_{xy} is a numerical measure of linear relationship b/w them and is defined as

$$\begin{aligned} r(x, y) &= r_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x) \cdot \text{Var}(y)}} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2 \cdot \frac{1}{n} \sum (y_i - \bar{y})^2}} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sigma_x \cdot \sigma_y} \end{aligned}$$

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Alternative form of $\rho(x, y)$:

$$\rho_{xy} = \rho(x, y) = \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

n = No. of pairs of values of x and y .

if $X = x_i - \bar{x}$, $Y = y_i - \bar{y}$ deviation from A.M

$$\rho_{xy} = \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}}$$

Note :- $-1 \leq \rho_{xy} \leq 1$

$$\rightarrow \text{if } u = \frac{x - a}{h}, v = \frac{y - b}{k}$$

$$\text{Then } \rho_{xy} = \rho_{uv}$$

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

Ex:-1 from the data given below, find the no. of items n :

$$\rho_{xy} = 0.5, \sum XY = 120, \sum X^2 = 90, \sigma_y = 8$$

where X and Y are deviations from the arithmetic mean. (2018)

Sol:- Given $\rho_{xy} = 0.5, \sum XY = 120, \sum X^2 = 90, \sigma_y = 8$
we have

$$\rho_{xy} = \frac{\sum XY}{\sqrt{\sum X^2 \cdot \sum Y^2}} \quad \text{where } X = x_i - \bar{x} \\ Y = y_i - \bar{y}$$

$$\text{Now } \sum y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$$

$$S^2 = \frac{1}{n} \sum y^2 = \frac{\sum y^2}{n}$$

$$\Rightarrow \sum y^2 = S^2 n = 64n$$

So

$$r_{xy} = \frac{120}{\sqrt{90 \times 64n}} = 0.5$$

Squaring both sides

$$\frac{14400}{90 \times 64n} = 0.25$$

$$n = \frac{14400}{90 \times 0.25 \times 64}$$

$$\boxed{n = 10}$$

Ex:-2 Calculate the coefficient of correlation b/w the age of husband & wife from the following data:

Age of husband : 35 34 40 43 56 20 38

Age of wife : 32 30 31 32 53 20 33

$$\text{Sol:- } \bar{x} = \frac{\sum x_i}{n} = \frac{266}{7} = 38, \bar{y} = \frac{\sum y_i}{n} = \frac{231}{7} = 33$$

x_i	y_i	$X = x_i - \bar{x}$	$Y = y_i - \bar{y}$	x^2	y^2	xy
35	32	-3	-1	9	1	3
34	30	-4	-3	16	9	12
40	31	2	-2	4	4	-4
43	32	5	-1	25	1	-5
56	33	18	20	324	400	360
20	20	-18	-13	324	169	234
38	33	0	0	0	0	0
$\sum x_i = 266$		$\sum y_i = 231$		$\sum x^2 = 702$		$\sum y^2 = 584$

\therefore Karl Pearson's coefficient of correlation

$$\begin{aligned}
 r_{xy} &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}} \\
 &= \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} \\
 &= \frac{600}{\sqrt{702} \sqrt{584}} = \frac{600}{640.287} \\
 &= 0.937
 \end{aligned}$$

Ex:- 3. find the coefficient of correlation b/w the values of x and y :

$x:$	1	3	5	7	8	10
$y:$	8	12	15	17	18	20

soln:-

x	y	x^2	y^2	xy
1	8	1	64	8
3	12	9	144	36
5	15	25	225	75
7	17	49	289	119
8	18	64	324	144
10	20	100	400	200
$\Sigma x = 34$	$\Sigma y = 90$	$\Sigma x^2 = 248$	$\Sigma y^2 = 1446$	$\Sigma xy = 582$

Karl Pearson's Coefficient of correlation is

$$r_{xy} = \frac{n \sum xy - \sum x \sum y}{\sqrt{n \sum x^2 - (\sum x)^2} \sqrt{n \sum y^2 - (\sum y)^2}}$$

$$\gamma_{xy} = \frac{6 \times 582 - 34 \times 96}{\sqrt{6 \times 248 - 34^2} \sqrt{6 \times 1446 - 96^2}} \\ = 0.9879$$

Ex-4. Find the coefficient of correlation for the following table

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

Sol:-

$$\text{Let } u = \frac{x-22}{4}, v = \frac{y-24}{6}$$

$$a=22, b=4$$

$$d=24$$

$$k=6$$

x	y	$u = \frac{x-22}{4}$	$v = \frac{y-24}{6}$	u^2	v^2	uv
10	18	-3	-1	9	1	3
14	12	-2	-2	4	4	4
18	24	-1	0	1	0	0
22	6	0	-3	0	9	0
26	30	1	1	1	1	1
30	36	2	2	4	4	4
		$\sum u = -3$	$\sum v = -3$	$\sum u^2 = 19$	$\sum v^2 = 19$	$\sum uv = 12$

$$\text{Here } n=6, \bar{u} = \frac{\sum u}{n}, \bar{v} = \frac{\sum v}{n}$$

$$\text{so } \bar{u} = \frac{-3}{6} = -\frac{1}{2}, \bar{v} = \frac{-3}{6} = -\frac{1}{2}$$

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

$$= \frac{6 \times 12 - (-3)(-3)}{\sqrt{6 \times 19 - (-3)^2} \sqrt{6 \times 19 - (-3)^2}}$$

$$= \frac{63}{\sqrt{105} \sqrt{105}} = 0.6$$

$$\gamma_{xy} = \gamma_{uv} = 0.6$$

Rank Correlation: →

Sometimes we have to deal with problems in which data can not be quantitatively measured but qualitative assessment is possible.

Spearman's Coefficient of Rank correlation

Let a group of n individuals be arranged in order of merit or proficiency in possession of two characteristics A and B. The ranks in two characteristics are, in general, different.

Let $x_i, y_i, i=1, 2, \dots, n$ be ranks of n indi. in the group for char. A and B respectively.

Then

Rank correlation coefficient

$$\gamma = 1 - \left[\frac{6 \sum D_i^2}{n(n^2-1)} \right]$$

where n is no. of observations

$D_i = x_i - y_i$ = Difference in ranks.

Note :- $-1 \leq \gamma \leq 1$

Ex: 1 Calculate the Rank correlation coefficient.

Marks in chem: 78 36 98 25 75 82 90 62 65 39

Marks in Maths: 84 51 91 60 68 62 86 58 63 47

Soln:-

X (chem)	Y (Maths)	R _x	R _y	D = R _x - R _y	D ²
78	84	4	3	1	1
36	51	9	9	0	0
98	91	1	1	0	0
25	60	10	7	3	9
75	68	5	4	1	1
82	62	3	6	-3	9
90.	86	2	2	0	0
62	58	7	8	-1	1
65	63	6	5	1	1
39	47	8	10	-2	4
					$\sum D^2 = 26$

80 rank correlation coefficient

$$\gamma = 1 - \frac{6 \sum D^2}{n(n^2-1)} \quad (n=10)$$

$$\begin{aligned} \gamma &= 1 - \frac{6 \times 26}{10(10^2-1)} \\ &= 1 - \frac{156}{990} \\ \gamma &= 0.8424 \end{aligned}$$

Ex:- Ten competitors in a beauty contest were ranked by three judges in the following order

First judge: 1 6 5 10 3 2 4 9 7 8

Second judge: 3 5 8 4 7 10 2 1 6 9

Third judge: 6 4 9 8 1 2 3 10 5 7

Use Method of rank correlation to determine which pair of judges has the nearest approach to

common taste in beauty?

Competitors	R ₁	R ₂	R ₃	D ₁₂ = R ₁ - R ₂	D ₁₃ = R ₁ - R ₃	D ₂₃ = R ₂ - R ₃	D ₁₂ ²	D ₁₃ ²	D ₂₃ ²
A	1	3	6	-2	-5	-3	4	25	9
B	6	8	4	1	2	1	1	4	1
C	5	8	9	-3	-4	-1	9	16	1
D	10	4	8	6	2	-4	36	4	16
E	3	7	1	-4	2	6	16	4	36
F	2	10	2	-8	0	8	64	0	84
G	4	2	3	2	1	-1	4	1	1
H	9	1	10	8	-1	-9	64	1	81
I	7	6	5	1	2	1	1	4	1
J	8	9	7	-1	1	2	1	1	4
n = 10							$\sum D_{12}^2 = 200$	$\sum D_{13}^2 = 60$	$\sum D_{23}^2 = 214$

rank correlation coefficient b/w first and second judges

$$\begin{aligned} r_{12} &= 1 - \frac{6 \sum D_{12}^2}{n(n^2-1)} \\ &= 1 - \frac{6 \times 200}{10(10^2-1)} = 1 - \frac{1200}{990} \\ &= -0.212 \end{aligned}$$

$$\begin{aligned} \text{likewise, } r_{13} &= 1 - \frac{6 \sum D_{13}^2}{n(n^2-1)} \\ &= 1 - \frac{6 \times 60}{10(10^2-1)} = 1 - \frac{360}{990} \\ &= 0.636 \end{aligned}$$

$$\begin{aligned} \text{and } r_{23} &= 1 - \frac{6 \sum D_{23}^2}{n(n^2-1)} = 1 - \frac{6 \times 214}{10(10^2-1)} \\ &= -0.297 \end{aligned}$$

correlation b/w first and second judges is negative i.e. their opinions regarding beauty test are opposite to each other. Similarly, opinions of second and third judges are opposite to each other, but opinions of first and third judges are of similar type as their correlation is positive. It means their likings and dislikings are very much common.

if the ranks are repeated
Then Rank correlation coefficient

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

where

$$E = \frac{1}{12} m_1(m_1^2 - 1) + \frac{m_2(m_2^2 - 1)}{12} + \dots$$

where m_1, m_2, \dots stands for the number of times different items repeats

Ex:- obtain the rank correlation coefficient of 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

sol:-

X	Y	R _x	R _y	D = R _x - R _y	D ²
68	62	4	5	-1	1
64	58	6	7	-1	1
75	68	2.5	3.5	-1	1
50	45	9	10	-1	1
64	81	6	1	5	25
80	60	1	6	-5	25
75	68	2.5	3.5	-1	1
40	48	10	9	1	1
55	50	8	8	0	0
64	70	6	2	4	16
$\sum D^2 = 72$					

Here $n = 10$

80

$$r = 1 - \frac{6 [\sum D^2 + E]}{n(n^2 - 1)}$$

$$E = \frac{m_1(m_1^2 - 1)}{12} + \frac{m_2(m_2^2 - 1)}{12} + \frac{m_3(m_3^2 - 1)}{12}$$

$$m = 2, 3, 2$$

$$= \frac{2(2^2 - 1)}{12} + \frac{3(3^2 - 1)}{12} + \frac{2(2^2 - 1)}{12}$$

$$= \frac{2 \times 3 + 3 \times 8 + 2 \times 3}{12} = \frac{36}{12}$$

$$= 3$$

$$r = 1 - \frac{6 [72 + 3]}{10(10^2 - 1)}$$

$$= 1 - \frac{6 \times 75}{99 \times 10}$$

$$= \frac{6}{11} = 0.545$$



Regression :-

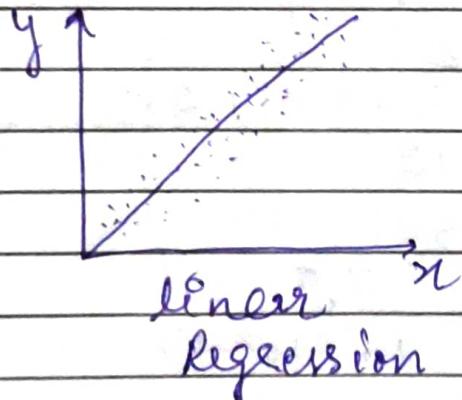
Regression means functional relationship between two or more related variables.

The only fundamental difference, if any, between problems of curve fitting and regression is that in regression, any of the variables may be considered as independent or dependent while in curve fitting, one variable can not be dependent.

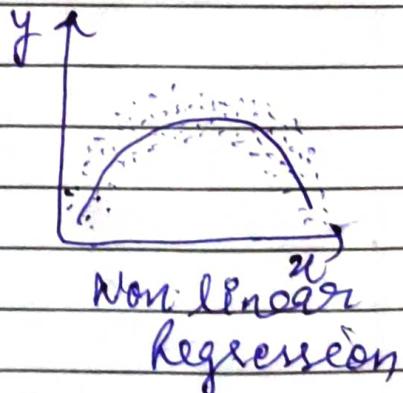
Regression measures the nature and extent of correlation. Regression is the estimation or prediction of unknown values of one variable from known values of another variable.

Linear Regression :-

When points in scatter diagram are concentrated around a straight line then it is linear regression otherwise Non-linear regression.



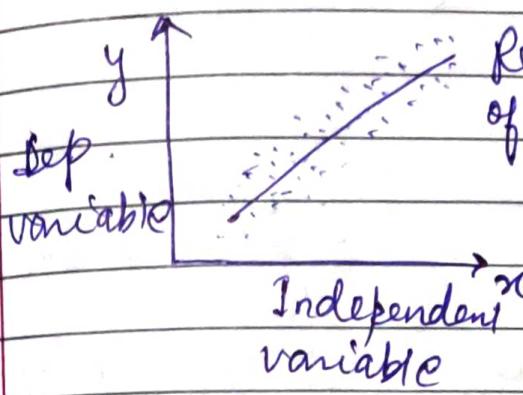
Scatter diagram



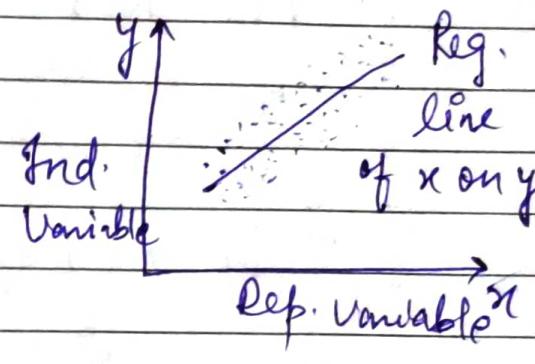
Suppose that the scatter diagram indicates some relationship b/w two variables x & y , the dots of the scatter diagram will be more or less concentrated around the curve.

This curve is called the curve of regression.

When curve is a straight line, it is called a line of regression, & regression is said to be linear.



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



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

① Derivation of line of Regression :-

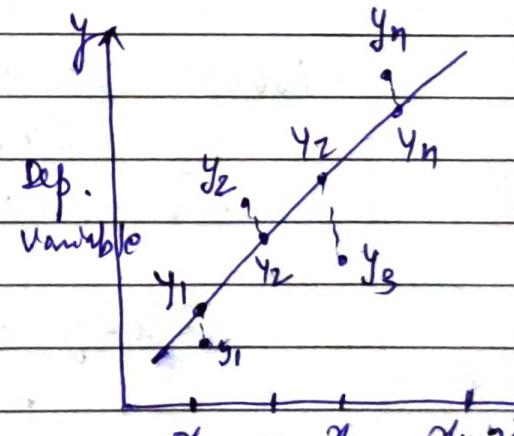
② y on x :-

Let $y = a + bx$ be the eqn of regression line of y on x

given value $\rightarrow y_i$

$y_i \rightarrow$ expected value

$$\text{Error } e_i = y_i - y_i^e$$



$y_i^e = a + bx_i$

$$e_i^2 = (y_i - y_i^e)^2$$

$$= (y_i - a - bx_i)^2 \quad \because y_i^e = a + bx_i$$

Now consider Total Error

$$U = \sum_{i=1}^n E_i^2$$

$$= \sum_{i=1}^n (y_i - a - b x_i)^2 \quad \text{--- (1)}$$

By Method of least squares, the constants a & b are chosen in such a way that the sum of the square of errors / residuals is minimum.

for min. value of U

$$\frac{\partial U}{\partial a} = 0$$

$$2 \sum_{i=1}^n (y_i - a - b x_i) (-1) = 0$$

$$\# \sum_{i=1}^n y_i - a \sum_{i=1}^n 1 - b \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \boxed{\sum y = a \cdot n + b \sum x} \quad \text{--- (2)}$$

$$\frac{\partial U}{\partial b} = 0$$

$$2 \sum_{i=1}^n (y_i - a - b x_i) (-x_i) = 0$$

$$\sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 = 0$$

$$\boxed{\sum xy = a \sum x + b \sum x^2} \quad \text{--- (3)}$$

(2) & (3) are normal eqns of $y = a + bx$

$$\text{Now } \sum x \sum y = n a / \sum x + b (\sum x)^2 \quad \text{from (2)}$$

$$= n \sum xy = n a / \sum x + b n \sum x^2 \quad \text{from (3)}$$

$$\sum x \sum y - n \sum xy = b [(\sum x)^2 - n \sum x^2]$$

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

from ② $\sum y = na + b \sum x$

$$\Rightarrow a = \frac{\sum y}{n} - \frac{b \sum x}{n}$$

$$\Rightarrow a = \bar{y} - b \bar{x}$$

$$\Rightarrow \boxed{\bar{y} = a + b \bar{x}}$$

\Rightarrow line $y = a + bx$ passes through the point (\bar{x}, \bar{y})

Putting $a_y = \bar{y} - b \bar{x}$ in eqⁿ $y = a + bx$

$$y = \bar{y} - b \bar{x} + bx$$

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

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

where

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

called Regression coefficient of y on x

In eqⁿ ③ shifting the origin to (\bar{x}, \bar{y}) , we get

$$\sum (x - \bar{x})(y - \bar{y}) = a \sum (x - \bar{x}) + b \sum (x - \bar{x})^2$$

$$\therefore \sum (x - \bar{x}) = \mu_x - n \bar{x}$$

$$\Rightarrow n \bar{x} \bar{y} = a \cdot 0 + b \cdot n \bar{x}^2$$

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

$$\Rightarrow b = \boxed{b_{yx} = r \frac{\bar{y}}{\bar{x}}}$$

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

called Regression Coef.

Hence line of regression of y on x is

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

where

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

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

ii) y

II. lin of regression of x on y \rightarrow

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

where

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

or

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

Note:-

if $r=0$, the two line of regression become $y = \bar{y}$ and $x = \bar{x}$ which are two straight lines parallel to x and y axes respectively and passing through their means \bar{y} and \bar{x} . They are mutually perpendicular.

if $r=\pm 1$, the two line of regression will coincide.

(H) Properties of Regression :-

1. Correlation coefficient is the geometric mean b/w the regression coefficients.

$$\text{Regression coefficients } b_{yx} = r \frac{\sigma_y}{\sigma_x}, b_{xy} = r \frac{\sigma_x}{\sigma_y}$$

$$\begin{aligned}\text{G.M of } b_{yx} \text{ and } b_{xy} &= \sqrt{b_{yx} \times b_{xy}} \\ &= \sqrt{r \frac{\sigma_y}{\sigma_x} \times r \frac{\sigma_x}{\sigma_y}} \\ &\leftarrow \sqrt{r^2} \\ &= r\end{aligned}$$

2. Arithmetic mean of regression coefficient is greater than the correlation coefficient.

$$\text{we have } \frac{b_{yx} + b_{xy}}{2} > r$$

$$\text{Now } \frac{1}{2} \left(r \frac{\sigma_y}{\sigma_x} + r \frac{\sigma_x}{\sigma_y} \right) > r$$

$$\Rightarrow r \left(\frac{\sigma_y^2 + \sigma_x^2}{2\sigma_x \sigma_y} \right) > r$$

$$\Rightarrow \frac{\sigma_y^2 + \sigma_x^2}{2\sigma_x \sigma_y} > 1$$

$$\Rightarrow \sigma_y^2 + \sigma_x^2 - 2\sigma_x \sigma_y > 0$$

$$\Rightarrow (\sigma_x - \sigma_y)^2 > 0$$

which is true

③ b_{yx} , b_{xy} and r have same sign.

Jump

#

Angle between two lines of regression :-

If θ is the acute angle between the two regression lines for the case of two variable x and y , show that

$$\tan \theta = \frac{1 - r^2}{r} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$$

where r, σ_x, σ_y are correlat^h coef. and regression coefficients

explain the significance of the formula $r=0$ and $r = \pm 1$. (2015, 17)

Proof:- Eqⁿ of line of regression
y on x

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

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

its slope $m_1 = r \frac{\sigma_y}{\sigma_x}$ ————— ①

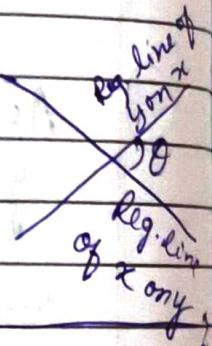
eqⁿ of line of regression of x on y

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

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

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

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



its slope $m_2 = \frac{1}{\gamma} \frac{\sigma_y}{\sigma_x}$ - (2)

for angle $\tan \theta = \pm \frac{m_2 - m_1}{1 + m_1 m_2}$

$$= \pm \frac{\frac{\sigma_y}{\sigma_x} - \gamma \frac{\sigma_y}{\sigma_x}}{1 + \frac{\sigma_y}{\sigma_x} \cdot \gamma \frac{\sigma_y}{\sigma_x}}$$

$$= \pm \frac{\frac{\sigma_y}{\sigma_x}}{\gamma} \cdot \frac{1 - \gamma^2}{\gamma} \cdot \frac{\frac{\sigma_x}{\sigma_x^2 + \sigma_y^2}}{\frac{\sigma_x^2}{\sigma_x^2 + \sigma_y^2}}$$

$$= \pm \frac{1 - \gamma^2}{\gamma} \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$$

Since $-1 \leq \gamma \leq 1 \Rightarrow \gamma^2 \leq 1$

and σ_x, σ_y are positive

\therefore the sign gives the acute angle b/w the lines

Hence

$$\boxed{\tan \theta = \frac{1 - \gamma^2}{\gamma} \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}}$$

(1) when $\gamma = 0 \Rightarrow \tan \theta = \infty \Rightarrow \theta = \pi/2$

\therefore two lines of regression are perpendicular to each other

(2) when $\gamma = \pm 1 \Rightarrow \tan \theta = 0 \Rightarrow \theta = 0 \text{ or } \pi$

Hence the lines of regression coincide & there is perfect correlation b/w x & y.

Ex:1. if the regression coefficients are 0.8 and 0.2 what would be the value of coefficient of correlation?

Sol:- we know that $r = \sqrt{b_{yx} \cdot b_{xy}}$

or $r^2 = b_{yx} \cdot b_{xy}$

$$= 0.8 \times 0.2$$

$$= 0.16$$

Since r has the same sign as both the regression coefficients b_{yx} and b_{xy} .
Hence

$$r = \sqrt{0.16}$$

$$= 0.4$$

Ex:2. calculate linear regression coefficients from the following data:

$x:$	1	2	3	4	5	6	7	8
$y:$	8	7	10	12	14	17	20	24

Sol:-

Linear regression coefficients are given by

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

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

Here $n = 8$

$$b_{yx} = \frac{8 \times 599 - 36 \times 107}{8 \times 204 - (36)^2} = \frac{940}{336} = 2.7776$$

and $b_{xy} = \frac{8 \times 599 - 36 \times 107}{8 \times 1763 - (107)^2} = \frac{940}{2655} = 0.3540$

table :-

x	y	x^2	y^2	xy
1	3	1	9	3
2	7	4	49	14
3	10	9	100	30
4	12	16	144	48
5	14	25	196	70
6	17	36	289	102
7	20	49	400	140
8	24	64	576	192
$\sum x = 36$	$\sum y = 107$	$\sum x^2 = 204$	$\sum y^2 = 1763$	$\sum xy = 599$

20/2/19
Ex:- 3

The following data regarding the heights (y) and the weights (x) of 100 college students are given

$$\sum x = 15000, \sum x^2 = 2272500, \sum y = 6800,$$

$$\sum y^2 = 463025, \sum xy = 1022250$$

find the Regression line of height on weight.

Sol:-

$$\text{given } \sum x = 15000, \sum x^2 = 2272500$$

$$\sum y = 6800, \sum y^2 = 463025$$

$$\sum xy = 1022250, n = 100$$

we have to find Regression line y on x

which is

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

where

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

$$= \frac{100 \times 1022250 - 15000 \times 6800}{100 \times 2272500 - 15000^2}$$

$$b_{yx} = 0.1$$

$$\text{also } \bar{x} = \frac{\sum x}{n} = \frac{15000}{100} = 150$$

$$\text{and } \bar{y} = \frac{\sum y}{n} = \frac{6800}{100} = 68$$

\therefore Regression line of height (y) on weight (x)

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

$$y - 68 = 0.1(x - 150)$$

$$y = 0.1x - 15 + 68$$

$$y = 0.1x + 53$$

Ex: 4. The following table gives age (x) in years of cars and annual maintenance cost (y) in hundred rupees:

$X:$	1	3	5	7	9
$Y:$	15	18	21	23	22

Estimate the maintenance cost for a 4-year-old car after finding the regression eqn.

Sol:-

x	y	xy	x^2
1	15	15	1
3	18	54	9
5	21	105	25
7	23	161	49
9	22	198	81

$$\sum x = 25 \quad \sum y = 99 \quad \sum xy = 533 \quad \sum x^2 = 165$$

$$\text{Here } n = 5$$

$$\bar{x} = \frac{\sum x}{n} = \frac{25}{5} = 5, \quad \bar{y} = \frac{\sum y}{n} = \frac{99}{5} = 19.8$$

Regression coefficient

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

$$= \frac{5 \times 533 - 25 \times 99}{5 \times 165 - (25)^2}$$

$$= 0.95$$

Hence Maintenance cost on age of cars (x) is given by Regression line

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

$$y - 19.8 = 0.95 (x - 5)$$

$$y - 19.8 = 0.95 x - 5 \times 0.95$$

$$y = 0.95x + 15.05$$

when $x = 4$ years

$$y = 0.95 \times 4 + 15.05$$

$$= 18.85 \bar{x}$$

$$\boxed{y = 18.85 \bar{x}}$$

2018

Ex: 5 In a partially destroyed laboratory record of an analysis of a correlation data the following results are only legible:

Variance of $x = 9$

Regression eqn : $8x - 10y + 66 = 0$

$$40x - 18y = 214$$

- What were (a) the mean values of x and y
 (b) the standard deviation of y
 (c) the coeff of correlation b/w x & y .

Soln:- given, variance of $x = \sigma_x^2 = 9$

regression eqn:

$$8x - 10y + 66 = 0 \quad \text{--- (2)}$$

$$40x - 18y - 214 = 0 \quad \text{--- (3)}$$

since regression lines (2) and (3) passes through the central point (\bar{x}, \bar{y}) then

$$8\bar{x} - 10\bar{y} + 66 = 0 \quad \text{--- (4)}$$

$$40\bar{x} - 18\bar{y} - 214 = 0 \quad \text{--- (5)}$$

from (4) & (5),

$$40\bar{x} - 50\bar{y} = -330$$

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

+

$$32\bar{y} = 544$$

$$\bar{y} = \frac{544}{32} = 17$$

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

using in eqn (1)

$$8\bar{x} - 10 \times 17 = -66$$

$$8\bar{x} = -66 + 170$$

$$8\bar{x} = 104$$

$$\boxed{2\bar{x} = 13}$$

eqn (2) and (3) can be written as

$$y = 0.8x + 6.6 \quad (\text{R.L of } y \text{ on } x)$$

\therefore Reg. Coeff. of y on x

$$b_{yx} = r \frac{\sigma_y}{\sigma_x} = 0.8 \quad \text{--- (6)}$$

eqⁿ ③ can be written as

$$40x = 18y + 214$$

Reg. line of x on y

$$x = 0.45y + 5.35$$

\therefore Reg. coef. of x on y

$$b_{xy} = \rho \frac{\sigma_x}{\sigma_y} = 0.45 \quad \text{--- ⑦}$$

Multiplying eqⁿ ⑥ & ⑦.

$$\rho \frac{\sigma_y}{\sigma_x} \times \rho \frac{\sigma_x}{\sigma_y} = 0.8 \times 0.45$$

$$\rho^2 = 0.36$$

Correlatⁿ coef.

$$\boxed{\rho = 0.6}$$

from ⑥,

$$\rho \frac{\sigma_y}{\sigma_x} = 0.8$$

$$\frac{\sigma_y}{\sigma_x} = 0.8$$

$$\sigma_y = \frac{0.8 \times \sigma_x}{\rho}$$

$$= \frac{0.8 \times 3}{0.6}$$

$$= 4$$

from ①

Hence S.D of y , $\boxed{\sigma_y = 4}$

& Variance $\sigma_y^2 = 16$.

2019

Ex:- 6. The equations of two regression lines obtained in a correlation analysis of 60 observations are

$$5x = 6y + 24$$

$$1000y = 768x - 3608$$

What is the correlation coefficient? Show that the ratio of coefficient of variability of x to that of y is 5. What is the ratio of variances of x and y ?

Soln:-

Regression line of x on y

$$5x = 6y + 24$$

$$x = \frac{6}{5}y + \frac{24}{5}$$

\therefore Regression coef. of x on y is

$$b_{xy} = \frac{6}{5} = r \frac{\sigma_x}{\sigma_y} \quad \text{--- (1)}$$

Regression line of y on x

$$1000y = 768x - 3608$$

$$y = 0.768x - 3.608$$

\therefore Regression coef. of y on x is

$$b_{yx} = 0.768 = r \frac{\sigma_y}{\sigma_x} \quad \text{--- (2)}$$

Multiplying eqn (1) & (2),

$$b_{xy} \times b_{yx} = \frac{6}{5} \times 0.768$$

$$r \frac{\sigma_x}{\sigma_y} \times r \frac{\sigma_y}{\sigma_x} = 0.9216$$

$$r^2 = 0.9216$$

correlation coef.

$$r = 0.96$$

Now dividing eq ① and ②

$$\frac{\sigma_x^2}{\sigma_y^2} = \frac{6}{5 \times 0.768} = 1.5625$$

Taking square root on both sides,

$$\frac{\sigma_x}{\sigma_y} = 1.25 = \frac{5}{4}$$

Ration of variances of x and y

$$\frac{\text{Var}(x)}{\text{Var}(y)} = \frac{\sigma_x^2}{\sigma_y^2} = \frac{5^2}{4^2} = \frac{25}{16}$$

now since the regression line passes through the point (\bar{x}, \bar{y})

$$5\bar{x} = 6\bar{y} + 24 \quad \dots \text{③}$$

$$1000\bar{y} = 768\bar{x} - 3608 \quad \dots \text{④}$$

Solving eqs ③ & ④, we get

$$\bar{x} = 6, \bar{y} = 1$$

Coef. of variability of $x = \frac{\sigma_x}{\bar{x}}$

coef. of variability of $y = \frac{\sigma_y}{\bar{y}}$

Ratio of coef. of variability of x & y

$$= \frac{\sigma_x}{\bar{x}} \times \frac{\bar{y}}{\sigma_y}$$

$$= \frac{5}{6} \times \left(\frac{6}{4} \right) = \frac{1}{6} \times \frac{5}{4} = \frac{5}{24}$$

(2)

Polynomial fit: Non-linear Regression

Let $y = a + bx + cx^2$ —①

be a second degree parabolic curve of regression of y on x to be fitted for the data (x_i, y_i) , $i = 1, 2, \dots, n$.

Normal eqns of ① are

$$\sum y = a \cdot n + b \sum x + c \sum x^2$$

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

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

We find a, b and c from these eqns and put in eqn ①

Ex:- fit a second degree parabolic curve of regression of y on x to the following data:

$x :$	1	2	3	4
$y :$	6	11	18	27

Sol:-

Let $y = a + bx + cx^2$ be the 2nd degree parabolic curve to be fitted in the data its normal eqns are

$$\sum y = a \cdot n + b \sum x + c \sum x^2$$

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

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

Here $n = 4$

x	y	x^2	x^3	x^4	xy	x^2y
1	6	1	1	1	6	6
2	11	4	8	16	22	44
3	18	9	27	81	54	162
4	27	16	64	256	108	432
$\Sigma x = 10$	$\Sigma y = 62$	$\Sigma x^2 = 30$	$\Sigma x^3 = 100$	$\Sigma x^4 = 354$	$\Sigma xy = 195$	$\Sigma x^2y = 699$

\therefore Normal eq's are

$$62 = 4a + 10b + 30c$$

$$195 = 10a + 30b + 100c$$

$$644 = 30a + 100b + 354c$$

Solving these eq's we get

$$a = 3, b = 2, c = 1$$

Hence $y = 3 + 2x + x^2$ is
 required regression parabolic curve
 of y on x .