

UNIT - I :

Measures of Central Tendency:

Class:

Each stated interval such as 0-10, 10-20, ... etc., is called a class.

Class limit:

There are two limits of every class, they are lower limit and upper limit

Class Interval / Class width:

The difference between the boundaries is known as Class Interval i.e., $(U-L)$ or (L_2-L_1) , i.e., upper limit - lower limit

Mid values:

The values lying in the middle of the two class limits of a group or a class is known as mid value.

$$\text{i.e., } \frac{L_1 + L_2}{2}$$

Class frequency:

The number of observations to a particular class is known as class frequency.

Types of Series:

There are 3 types of series:

1) Individual series :

Every observation is independent or separate is called Individual series.

Ex: 1, 2, 3, 4, ...

2) Discrete series:

The series dealing with the discrete

Discrete series means where frequencies are given but variable is without class interval.

Ex:	Marks (X)	No. of students (f)
	10	2
	20	3
	30	4
	40	5

3) Continuous series:

There are 5 types of continuous series:

i) Inclusive series:

In this Inclusive series first class interval of its upper limit is included and immediate next class interval of its lower limit is also included.

Ex:	Marks (X)	No. of students (f)
	10 - 19	2
	20 - 29	3
	30 - 39	4
	40 - 49	5

ii) Exclusive series:

In this exclusive series first class interval of its upper limit is excluded and immediate next class interval of its lower limit is included.

Ex:	Marks (X)	No. of students (f)
	10 - 20	2
	20 - 30	3
	30 - 40	4
	40 - 50	5

iii) Open-end series:

In this open-end series first class interval of its lower limit is not given and last class interval (highest class interval) of its upper limit is not given.

Ex:

Marks (X)	No. of students (f)
below 20	2
20-30	3
30-40	4
Above 40	5

iv) Mid-value series:

In this series, difference between each mid point and dividing with two, therefore subtracting and adding to each midpoint we get the interval.

Ex:

midpoint	frequency	class interval
6 $\frac{10-6}{2}=2$	2	(6-2 - 6+2) 4 - 8
10 $\frac{14-10}{2}=2$	7	8 - 12
14 $\frac{18-14}{2}=2$	18	12 - 16
18	29	16 - 20
22	17	20 - 24

v) Cumulative frequency series:

It is of 2 types:

i) Less than cumulative frequency series:

(Lcf)

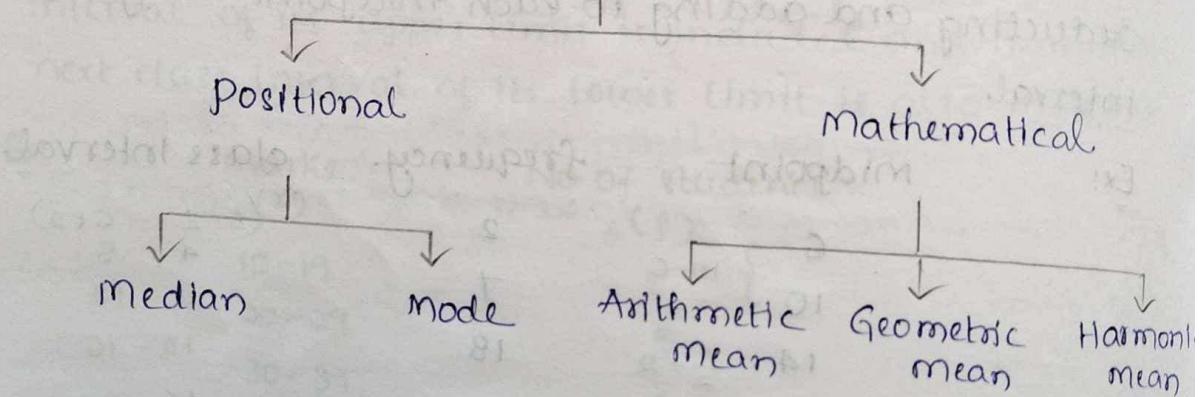
Ex:	C-I	f	Less than C.F	C-I
	below 40	3	3	30-40
"	50	11	11-3=8	40-50
"	60	34	34-11=23	50-60
"	70	59	59-34=25	60-70
"	80	72	72-59=13	70-80
"	90	80	80-72=8	80-90

2) More than Cumulative frequency Series:

<u>Ex:</u>	<u>C.I</u>	<u>f</u>	<u>more than C.F</u>	<u>C.I</u>
Above 40		80	80-71 = 3	
" 50		77	77-69 = 8	
" 60		69	69-46 = 23	
" 70		46	46-21 = 25	
" 80		21	21-8 = 13	
" 90		8	8-0 = 8	
				$N = 80$

Types of Averages:

There are 2 types of averages:



Arithmetic mean:

It is denoted by \bar{x}

$$\therefore \bar{x} = \frac{\text{Sum of observations}}{\text{No. of observations}}$$

• Individual series:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}; n \rightarrow \text{no. of observations}$$

• Discrete series:

$$\bar{x} = \frac{\sum f x}{N}; N \rightarrow \text{sum of frequencies}$$

problems:

1) find mean for the following data:

30, 41, 47, 54, 23, 34, 37, 51, 53, 47

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{417}{10} = 41.7$$

2) calculate arithmetic mean from the following data:

Marks (x)	No. of students (f)	fx
4	6	24
8	12	96
12	18	216
16	15	240
20	9	180
	<u>N=60</u>	<u>$\sum fx = 756$</u>

$$\bar{x} = \frac{\sum fx}{N} = \frac{756}{60} = 12.6$$

• continuous series:

$$\bar{x} = \frac{\sum fm}{N}; m \rightarrow \text{mid value}, N \rightarrow \text{sum of frequencies}$$

3) find average income of:

Exclusive series	Income	No. of persons (f)	m	fm
	10 - 20	4	15	60
	20 - 30	7	25	175
	30 - 40	16	35	560
	40 - 50	20	45	900
	50 - 60	15	55	825
	60 - 70	8	65	520
		<u>N=70</u>		<u>$\sum fm = 3040$</u>

$$\bar{x} = \frac{\sum fm}{N} = \frac{3040}{70} = 43.42$$

4) Find arithmetic mean for the following data:

C.I	frequency	C.I	m	fxm
4-6	1	3.5-6.5	5	5
7-9	3	6.5-9.5	8	24
$\frac{1-6}{2} = 0.5$	10-12	9.5-12.5	11	77
$\frac{10-9}{2} = 0.5$	13-15	12.5-15.5	14	210
16-18	11	15.5-18.5	17	187
19-21	3	18.5-21.5	20	60
22-24	2	21.5-24.5	23	46
<u>$N = 42$</u>				<u>$\sum fm = 609$</u>

$$\bar{x} = \frac{\sum fm}{N} = \frac{609}{42} = 14.5$$

Median:

Median is a positional measure. The position which is exactly in the centre, equal number of terms lie on either side of it, when terms are arranged in ascending order / descending order.

• Individual series:

case(i): when no. of obs'ns are odd

Median = size of $(\frac{n+1}{2})^{\text{th}}$ observation

case(ii): when no. of obs'ns are even

$$\begin{aligned} \text{Median} &= \text{Avg. of size of } (\frac{n}{2})^{\text{th}} \text{ & } (\frac{n}{2} + 1)^{\text{th}} \text{ obs} \\ &= \frac{\text{size of } (\frac{n}{2})^{\text{th}} + (\frac{n}{2} + 1)^{\text{th}} \text{ obs}}{2} \end{aligned}$$

Problems:

1) Find median of 4, 8, 16, 20, 2, 7, 13

2, 4, 7, 8, 13, 16, 20

$$n = 7$$

Median = size of $(\frac{n+1}{2})^{\text{th}}$ obs

2) Find median of 4, 8, 16, 20, 2, 7

2, 4, 7, 8, 16, 20 ; n = 6

$$\text{median} = \frac{\text{size of } \left(\frac{6}{2}\right)^{\text{th}} + \left(\frac{6}{2}+1\right)^{\text{th}} \text{ obs}}{2} = \frac{\text{size of } 3^{\text{rd}} + 4^{\text{th}} \text{ obs}}{2}$$

$$= \frac{7+8}{2}$$

$$= 7.5$$

- Discrete Series:

Median = size of $\left(\frac{N+1}{2}\right)^{\text{th}}$ observation

- 3) Compute median from the following data:

X	f	Cf	
10	4	4	Median = size of $\left(\frac{N+1}{2}\right)^{\text{th}}$ obs.
20	7	11	
30	21	32	$= \left(\frac{106+1}{2}\right)^{\text{th}}$ obs.
median $\leftarrow \frac{40}{50}$	34	66 \rightarrow Median class	$= (53.5)^{\text{th}}$ obs.
	25	91	
60	12	103	\therefore median = 40
70	3	106	
<u>N = 106</u>			

- Continuous Series:

$$\text{Median} = L + \frac{\frac{N}{2} - Cf}{f} \times i;$$

where, L \rightarrow lower limit of the median class

cf \rightarrow cumulative frequency preceding the median class

f \rightarrow frequency corresponding to the median class

N \rightarrow Total of frequency.

i \rightarrow common factor of median class.
(upper limit - lower limit)

4) Find median from the following data:

Mark	No. of students (f)	C.f
0-10	20	20
10-20	45	65
20-30	85	150
30-40	160	310
40-50	70	380
50-60	55	435
60-70	35	470
70-80	30	500
<u>N = 500</u>		

$$\frac{N}{2} = 250 ;$$

$$cf = 150 ; f = 160 ; L = 30 ; i = 10$$

$$\text{median} = 30 + \frac{250 - 150}{160} \times 10 = 30 + \frac{100}{16} = 36.25$$

Mode:

Mode is the most repeated frequency or value.

• Individual series;

Problems:

1) Find mode : 1, 2, 3, 4, 5, 6, 1, 2, 3, 4, 5, 5

$$\text{Mode} = 5$$

• Discrete series:

2) Find mode from the following data:

X	f
4	3
7	9
11	14

$$\text{mode} = \frac{16}{25} \rightarrow \text{Max. frequency}$$

$$\therefore \text{Mode} = 16$$

3) find mode from the following data:

	X	I	(2)	II	III	IV	V	VI	$\frac{1}{3}$	$\frac{2}{3}$
Grouping Table	X	f								
5	1	4			Leave 1st					
10	3	7				8				
15	4	13					16			
20	9		20					24		
25	11		23			32				
30	max 12			15			26			17
35	3	5								
40	2		4		7					
45	2			Add 2			Add 3			

Analysis Table

X	I	II	III	IV	V	VI	Total
5							
10							
15					✓		1
20			✓	✓		✓	3
25		✓	✓	✓	✓	✓	5
30	✓	✓		✓	✓		4
35					✓		1
40							
45							

$$\therefore \text{Mode} = 25$$

• Continuous series:

$$\text{Mode} = L + \frac{f_1 - f_0}{2f_1 - f_0 - f_2} \times i$$

where, $L \rightarrow$ lower limit corresponding to the modal class

$f_1 \rightarrow$ modal class frequency

$f_0 \rightarrow$ preceding the modal class "

$f_2 \rightarrow$ succeeding the modal class "

$i \rightarrow$ common factor corresponding to the modal class

4) Calculate mode for the following data:

Grouping table X	f	(I)	Add 2 (1/2)	Add 3 (1/3)	(2/3)	X
		(II)	(III)	(IV)	(V)	(VI)
0-10	4			Leave 1st	Leave 1st	Leave 1st & add 3
10-20	13	17		38		
20-30	21		34		78	
30-40	44	65			98	
40-50	33		77	99		
50-60	22	55				
60-70	7		29	62		

Analysis table

X	I	II	III	IV	V	VI	Total
0-10							
10-20					✓		1
20-30		✓			✓	✓	3
Modal class 30-40	✓	✓	✓	✓	✓	✓	6
40-50			✓	✓		✓	3
50-60				✓			1
60-70							

$$\text{Mode} = L + \frac{f_1 - f_0}{2f_1 - f_0 - f_2} \times i$$

$$= 30 + \frac{44 - 21}{2 \times 44 - 21 - 33} \times 10$$

$$= 30 + \frac{230}{34}$$

$$= 30 + 6.76$$

$$\therefore z = \text{Mode} = 36.76$$

Standard deviation: (σ)

{Introduced by Karl Pearson in 1893}

The idea of standard deviation was introduced by Karl Pearson in 1893. This measure is widely used for studying dispersion.

standard deviation means the square root of the mean of the squared deviations from the actual mean

• individual series:

$$\sigma = \sqrt{\frac{\sum x_i^2}{n}} ; x_i = x - \bar{x} ; \bar{x} \rightarrow \text{sample mean}$$

(direct method)

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

$$\sigma = \sqrt{\frac{\sum d^2}{n} - \left(\frac{\sum d}{n}\right)^2} ; d = x - A ; A \rightarrow \text{Assumed mean}$$

(short cut method)

Problems:

1. Find standard deviation by using direct and short cut methods:

25, 27, 31, 32, 35

direct method

$$\bar{x} = \frac{\sum x_i}{n} = \frac{25+27+31+32+35}{5} = \frac{150}{5} = 30$$

$$x \quad x_i - \bar{x} \rightarrow 30 \quad (x_i - \bar{x})^2$$

$$25 \quad -5 \quad 25$$

$$27 \quad -3 \quad 9$$

$$31 \quad 1 \quad 1$$

$$32 \quad 2 \quad 4$$

$$35 \quad 5 \quad 25$$

$$\sigma = \sqrt{\frac{64}{5}}$$

$$\therefore \sigma = 3.56$$

$$\sum (x_i - \bar{x})^2 = 64$$

Short cut method

$$\begin{array}{ccc}
 X & d = x - A & d^2 \\
 25 & -6 & 36 \\
 27 & -4 & 16 \\
 31 & 0 & 0 \\
 32 & 1 & 1 \\
 35 & 4 & 16 \\
 \hline
 \sum d = -5 & \sum d^2 = 69
 \end{array}$$

$$\begin{aligned}
 \sigma &= \sqrt{\frac{\sum d^2}{n} - \left(\frac{\sum d}{n}\right)^2} \\
 &= \sqrt{\frac{69}{5} - \left(-\frac{5}{5}\right)^2} \\
 &= \sqrt{13.8 - 1} \\
 &= \sqrt{12.8} \\
 \sigma &= 3.56
 \end{aligned}$$

• Discrete Series:

Direct Method:

$$\sigma = \sqrt{\frac{\sum f x^2}{N}} ; x = x - \bar{x} \text{ where } \bar{x} = \frac{\sum f x}{N}$$

Shortcut Method:

$$\sigma = \sqrt{\frac{\sum f d^2}{N} - \left(\frac{\sum f d}{N}\right)^2} ; d = x - A \text{ where } A \rightarrow \text{Assumed mean}$$

2. Calculate standard deviation from the following data:

<i>Direct Method</i>	X	f	fx	$x - \bar{x} \rightarrow 8.47$	$(x - \bar{x})^2$	$f(x - \bar{x})^2$
	2	5	10	-6.47	41.8609	209.3045
	4	15	60	-4.47	19.9809	299.7135
	6	20	120	-2.47	6.1009	122.018
	8	25	200	-0.47	0.2209	5.5225
	10	25	250	1.53	2.3409	58.5225
	12	20	240	3.53	12.4609	249.218
	15	8	120	6.53	42.6409	341.1272
	$\sum N = 118$		$\sum f x = 1000$			

$$\sum f(x - \bar{x})^2 = 1285.4262$$

$$\bar{x} = \frac{\sum f x}{N} = \frac{1000}{118} = 8.47$$

$$\sigma = \sqrt{\frac{\sum f x^2}{N}} = \sqrt{\frac{1285.4262}{118}} = \sqrt{10.8934} = 3.3005$$

<i>X</i>	<i>f</i>	$d = X - A$	d^2	fd	fd^2
2	5	-6	36	-30	180
4	15	-4	16	-60	240
6	20	-2	4	-40	80
(8) A	25	0	0	0	0
10	25	2	4	50	100
12	20	4	16	80	320
15	8	7	49	56	392
<hr/> $N = 118$			<hr/> $\sum fd = 56$		
<hr/> $\sum fd^2 = 1312$			<hr/>		

$$\sigma = \sqrt{\frac{\sum fd^2}{N} - \left(\frac{\sum fd}{N}\right)^2}$$

$$= \sqrt{\frac{1312}{118} - \left(\frac{56}{118}\right)^2}$$

$$= \sqrt{11.11 - (0.47)^2}$$

$$= \sqrt{11.11 - 0.2209}$$

$$= \sqrt{10.8891}$$

$$= 3.2998$$

Variance:

Variance is the square of standard deviation. If a phenomenon is affected by a number of variables, variance helps in isolating the effects of different factors.

This term was used by R.A. Fisher in 1913

$$\text{Variance} = (\sigma)^2 = (S.D.)^2$$

• Continuous series:

Direct Method:

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

Step-deviation method: (short cut method)

$$(if class interval = 1) \quad \sigma = \sqrt{\frac{\sum f d_1^2}{N} - \left(\frac{\sum f d_1}{N}\right)^2} \quad m-A$$

$$(\text{if class intervals are not equal}) \quad \sigma = \sqrt{\frac{\sum fd'^2}{N} - \left(\frac{\sum fd'}{N}\right)^2}$$

3. Compute standard deviation and variance from the following data:

Short cut method C.I	f	m	$d' = \frac{m-A}{C}$	d'^2	fd'	fd'^2		
			Assumed mean $\uparrow = 5$					
0-10	1	5	0	0	0	0		
10-20	4	15	1	1	4	4		
20-30	17	25	2	4	34	68		
30-40	45	35	3	9	135	405		
40-50	26	45	4	16	104	416		
50-60	5	55	5	25	25	125		
60-70	2	65	6	36	12	72		
$\sum f = 100$			$\sum fd' = 314$			$\sum fd'^2 = 1090$		

$$\sigma = \sqrt{\frac{\sum fd'^2}{N} - \left(\frac{\sum fd'}{N}\right)^2} \times C$$

$$= \sqrt{\frac{1090}{100} - \left(\frac{314}{100}\right)^2} \times 10$$

$$= \sqrt{10.9 - 9.85} \times 10$$

$$= \sqrt{1.05} \times 10$$

$$= 1.0246 \times 10$$

$$\therefore \sigma = 10.246 \quad \& \quad \sigma^2 = 104.04$$

$$\overline{x} = A + \frac{\sum fd'}{N} \times C$$

$$= 5 + \frac{314}{100} \times 10$$

$$= 36.4$$

Diced
method

C.I	f	m	fm	$m-\bar{x}$	$(m-\bar{x})^2$	$f(m-\bar{x})^2$
0-10	1	5	5	-31.4	985.96	985.96
10-20	4	15	60	-21.4	457.96	1831.84
20-30	17	25	425	-11.4	129.96	2209.32
30-40	45	35	1575	-1.4	1.96	88.2
40-50	26	45	1170	8.6	73.96	1922.96
50-60	5	55	275	18.6	345.96	1729.8
60-70	2	65	130	28.6	817.96	1635.92

$$\underline{N = 100}$$

$$\underline{\sum fm = 3640}$$

$$\underline{\sum f(m-\bar{x})^2 = 10404}$$

$$\bar{x} = \frac{\sum fm}{N} = \frac{3640}{100} = 36.4$$

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

$$= \sqrt{\frac{10404}{100}}$$

$$= \sqrt{104.04}$$

$$\therefore \sigma = 10.2$$

Moments:

- There are 2 types of Moments:
- Non-central Moments / Raw Moments / Moments about a point
 - Central Moments / Moments about mean

Non-central Moments:

These moments are calculated from any value 'A' (or) any point 'A'. The r^{th} non-central moments is denoted by μ_r' .

$$\text{i.e., } \mu_r' = E(x-A)^r \text{ (or) } \frac{1}{N} \sum f_i (x_i - A)^r \quad [\text{Discrete series}]$$

$$\mu_r' = E(x-A)^r = \frac{1}{n} \sum (x_i - A)^r \quad [\text{Individual series}]$$

If $r=0$,

$$\mu_0' = \frac{1}{N} \sum f_i \underbrace{(x_i - A)}_1^0$$

$$\therefore \mu_0' = 1$$

If $r=1$,

$$\begin{aligned} \mu_1' &= \frac{1}{N} \sum f_i (x_i - A) \\ &= \frac{\sum f_i x_i}{N} - \frac{1}{N} \sum f_i A \\ &= \bar{x} - A \end{aligned}$$

$$\therefore \mu_1' = \bar{x} - A$$

By, $r=2, 3, 4 \dots$ we get 1st 4 Non-central Moments

$$\mu_2' = \frac{1}{N} \sum f_i (x_i - A)^2 \quad (\text{or}) \quad \mu_2' = \frac{\sum f_i d_i^2}{N}, \quad d_i = x_i - A$$

$$\mu_3' = \frac{1}{N} \sum f_i (x_i - A)^3 \quad (\text{or}) \quad \mu_3' = \frac{\sum f_i d_i^3}{N}$$

$$\mu_4' = \frac{1}{N} \sum f_i (x_i - A)^4$$

Central Moments:

These moments are calculated from actual mean \bar{x} . The r^{th} central moments is defined as

i.e., $M_r = E(X - \bar{X})^r$ (or) $\frac{1}{N} \sum f_i (x_i - \bar{X})^r$ [Discrete series]

$M_r = E(X - \bar{X})^r = \frac{1}{n} \sum f_i (x_i - \bar{X})^r$ [Individual series]

If $r=0$,

$$M_0 = \frac{1}{N} \sum f_i (\underbrace{x_i - \bar{X}}_0)^r$$

$$\therefore M_0 = 1$$

If $r=1$,

$$M_1 = \frac{1}{N} \sum f_i (x_i - \bar{X})^1$$

$$= \frac{\sum f_i x_i}{N} - \frac{1}{N} \sum f_i \bar{X}$$

$$= \bar{X} - \bar{X}$$

$$\therefore M_1 = 0$$

If $r=2, 3, 4$ we get 1st 4 central moments

$$M_2 = \frac{1}{N} \sum f_i (x_i - \bar{X})^2 \text{ (or)} M_2 = \frac{\sum f_i d_i^2}{N}; d_i = x_i - \bar{X}$$

$$M_3 = \frac{1}{N} \sum f_i (x_i - \bar{X})^3 \text{ (or)} M_3 = \frac{\sum f_i d_i^3}{N}$$

$$M_4 = \frac{1}{N} \sum f_i (x_i - \bar{X})^4$$

Properties of Moments (Central Moments):

1. The first moment about mean is always zero

$$\text{i.e., } M_1 = 0$$

2. The second moment about mean indicates variance

$$\text{i.e., } M_2 = \sigma^2$$

3. The third moment about mean measures of skewness of a given distribution.

4. The fourth moment about mean measures kurtosis of a frequency distribution.

Relation between central and Non-central Moments:

Central Moments:

$$\mu_0 = 1$$

$$\mu_1 = 0$$

$$\mu_2 = \mu_2' - (\mu_1')^2$$

$$\mu_3 = \mu_3' - 3\mu_2'\mu_1' + 2(\mu_1')^3$$

$$\mu_4 = \mu_4' - 4\mu_3'\mu_1' + 6\mu_2'(\mu_1')^2 - 3(\mu_1')^4$$

Non-central Moments:

$$\mu_0' = 1$$

$$\mu_1' = \bar{x} - A$$

$$\mu_2' = \mu_2 + (\mu_1')^2$$

$$\mu_3' = \mu_3 + 3\mu_2(\mu_1') + (\mu_1')^3$$

$$\mu_4' = \mu_4 + 4\mu_3\mu_1' + 6\mu_2(\mu_1')^2 + (\mu_1')^4$$

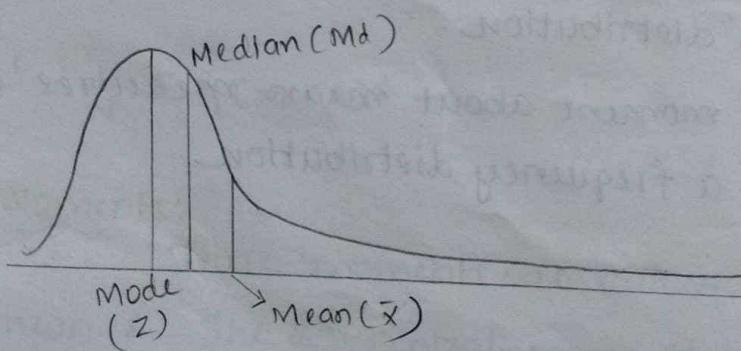
Skewness:

Skewness is a measure of symmetry or more precisely the lack of symmetry.

We study skewness to have an idea about the shape of the curve which we can draw with the help of given data.

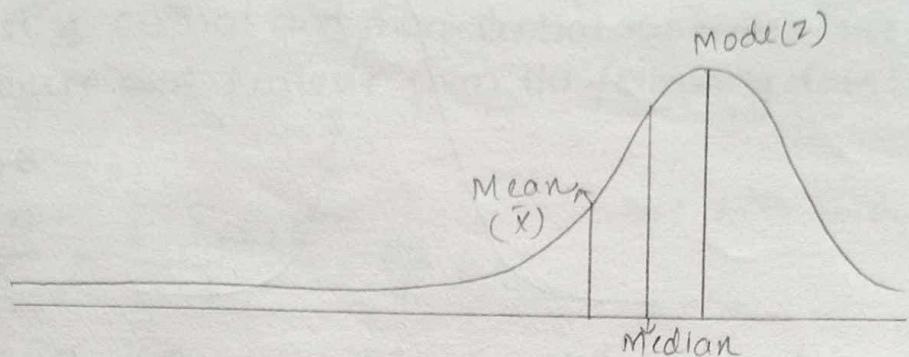
The third moment about mean measure of skewness of given frequency distribution

1. If $\mu_3 > 0$, that is if the frequency curve has longer tail to right, then the distribution is positively skewed.

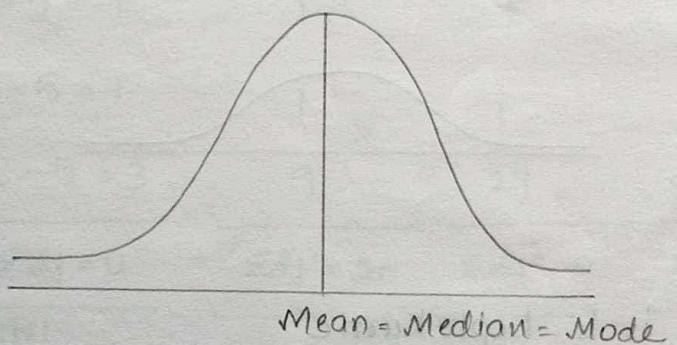


$$B_1 = \frac{\mu_3^2}{\mu_2^3}$$

2. If $\mu_3 < 0$, i.e., if the frequency curve has longer tail to left, then the distribution is negatively skewed.



3. If $\mu_3 = 0$, then the distribution is symmetrical.



Kurtosis:

Kurtosis measures the relative peakedness or flatness of the distribution.

Coefficient of kurtosis based on second and fourth central moments, and it is denoted by B_2 .

$$\therefore B_2 = \frac{\mu_4}{\mu_2^2} \quad (\text{or}) \quad V_2 = B_2 - 3$$

limits are -3 to +3

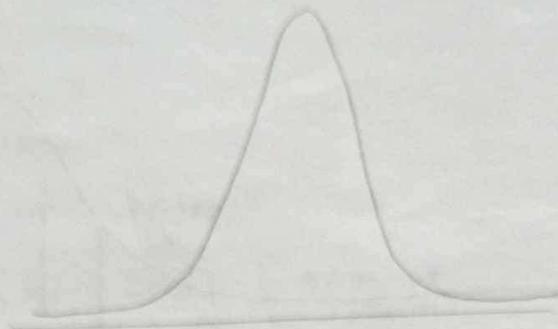
$\left(\text{gamma}_2\right)$

Meso Kurtic (or) Normal curve:

The curve is Normal or Symmetrical is called Meso Kurtic.

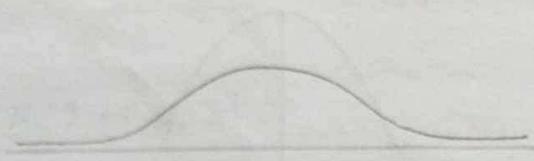
Lepto Kurtic Curve:

If a curve which is more peaked than the normal (or) Meso Kurtic curve, then it is called Lepto Kurtic curve.



Platy Kurtic Curve:

The curve which is less peaked than the normal (or) Meso Kurtic curve, is called Platy Kurtic Curve.



Note:

- 1) If $B_2 > 3$, it is Lepto Kurtic
- 2) If $B_2 < 3$, it is Platy Kurtic
- 3) If $B_2 = 3$, it is Meso Kurtic / Normal curve

Coefficient of Skewness:

The coefficient of skewness based on second and third central moments, and it is denoted by β_1 .

$$\therefore \beta_1 = \frac{\mu_3^2}{\mu_2^3} \text{ and } \gamma_1 = \pm \sqrt{\beta_1}$$

- The sign of skewness depends upon the sign of μ_3

Note:

- 1) If $\beta_1 > 0$ / $\gamma_1 > 0$, then the distribution is positively skewed.

3) If $\beta_1 = 0 / \gamma_1 = 0$, then the distribution is normal & symmetrical

problems

1. find first 4 central and non-central moments and also find Skewness and Kurtosis from the following data:

2, 4, 6, 8

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{2+4+6+8}{4} = \frac{20}{4} = 5$$

x	$d_i = x - \bar{x}$	d_i^2	d_i^3	d_i^4
2	$2-5 = -3$	9	-27	81
4	$4-5 = -1$	1	-1	1
6	$6-5 = 1$	1	1	1
8	$8-5 = 3$	9	27	81
	$\sum d_i = 0$	$\sum d_i^2 = 20$	$\sum d_i^3 = 0$	$\sum d_i^4 = 164$

Central Moments:

$$\mu_r = \frac{1}{n} \sum (x_i - \bar{x})^r$$

$$\text{if } r=1, \quad \mu_1 = \frac{1}{n} \sum d_i = 0$$

$$\text{if } r=2, \quad \mu_2 = \frac{1}{n} \sum d_i^2 = \frac{20}{4} = 5$$

$$\text{if } r=3, \quad \mu_3 = \frac{1}{n} \sum d_i^3 = 0$$

$$\text{if } r=4, \quad \mu_4 = \frac{1}{n} \sum d_i^4 = \frac{164}{4} = 41$$

$$\text{Skewness} = \beta_1 = \frac{\mu_3^2}{\mu_2^3} = \frac{(0)^2}{(5)^3} = 0$$

$$\boxed{\beta_1 = 0} \rightarrow \text{Normal}$$

$$\text{Kurtosis} = \beta_2 = \frac{\mu_4}{\mu_2^2} = \frac{41}{25} = 1.64$$

$$\boxed{\beta_2 = 1.64} \rightarrow \text{Platykurtic}$$

x_i	$d_i = x_i - A$	d_i^2	d_i^3	d_i^4
2	$2-4 = -2$	4	-8	16
4	$4-4 = 0$	0	0	0
6	$6-4 = 2$	4	8	16
8	$8-4 = 4$	16	64	256
	$\sum d_i = 4$	$\sum d_i^2 = 24$	$\sum d_i^3 = 64$	$\sum d_i^4 = 288$

Non-central Moments:

$$\mu_r^1 = \frac{1}{n} \sum (x_i - A)^r$$

$$\text{If } r=1, \mu_1^1 = \frac{1}{4}(4) = 1$$

$$\text{If } r=2, \mu_2^1 = \frac{1}{4}(24) = 6$$

$$\text{If } r=3, \mu_3^1 = \frac{1}{4}(64) = 16$$

$$\text{If } r=4, \mu_4^1 = \frac{1}{4}(288) = 72$$

$$\mu_1 = 0$$

$$\mu_2 = \mu_2^1 - (\mu_1^1)^2 = 6 - 1 = 5$$

$$\begin{aligned} \mu_3 &= \mu_3^1 - 3\mu_2^1\mu_1^1 + 2(\mu_1^1)^3 \\ &= 16 - 3(6)(1) + 2(1)^3 \\ &= 0 \end{aligned}$$

$$\begin{aligned} \mu_4 &= \mu_4^1 - 4\mu_3^1\mu_1^1 + 6\mu_2^1(\mu_1^1)^2 - 3(\mu_1^1)^4 \\ &= 72 - 4(16)(1) + 6(6)(1)^2 - 3(1)^4 \\ &= 41 \end{aligned}$$

$$\text{Skewness} = \beta_1 = \frac{\mu_3^2}{\mu_2^3} = \frac{(0)^2}{(5)^3} = 0$$

$$\text{Kurtosis} = \beta_2 = \frac{\mu_4^2}{\mu_2^4} = \frac{41^2}{25^3} = 1.64$$

2. Calculate first 4 non-central moments for the following data! Also find coefficient of skewness &

X	f	$d_i = x_i - A$	d_i^2	d_i^3	d_i^4	$\sum f d_i$	Kurtosis		
							$\sum f d_i^2$	$\sum f d_i^3$	$\sum f d_i^4$
0	1	-4	16	-64	256	-4	16	-64	256
1	8	-3	9	-27	81	-24	72	-216	648
2	28	-2	4	-8	16	-56	112	-224	448
3	56	-1	1	-1	1	-56	56	-56	56
4	70	0	0	0	0	0	0	0	0
5	56	1	1	1	1	56	56	56	56
6	28	2	4	8	16	56	112	224	448
7	8	3	9	27	81	24	72	216	648

$$\mu_r^1 = \frac{1}{N} \sum f_i (x_i - A)^r ; x_i - A = d_i$$

$$\text{if } r=1, \mu_1^1 = \frac{1}{N} \sum f_i d_i = 0$$

$$\text{if } r=2, \mu_2^1 = \frac{1}{N} \sum f_i d_i^2 = \frac{512}{256} = 2$$

$$\text{if } r=3, \mu_3^1 = \frac{1}{N} \sum f_i d_i^3 = 0$$

$$\text{if } r=4, \mu_4^1 = \frac{1}{N} \sum f_i d_i^4 = \frac{2816}{256} = 11$$

$$\mu_1 = 0$$

$$\mu_2 = \mu_2^1 - (\mu_1^1)^2 = 2 - 0 = 2$$

$$\mu_3 = \mu_3^1 - 3\mu_2^1\mu_1^1 + 2(\mu_1^1)^3 = 0 - 3(4)(0) + 2(0)^3 = 0$$

$$\mu_4 = \mu_4^1 - 4\mu_3^1\mu_1^1 + 6\mu_2^1(\mu_1^1)^2 - 3(\mu_1^1)^4$$

$$= 11 - 4(0)(0) + 6(2)(0)^2 - 3(0)^4 = 11$$

$$\text{Coefficient of skewness} = \beta_1 = \frac{\mu_3^2}{\mu_2^3} = \frac{(0)^2}{(2)^3} = 0 \rightarrow \text{Normal}$$

$$\text{Kurtosis} = \beta_2 = \frac{\mu_4}{\mu_2^2} = \frac{11}{4} = 2.75 \rightarrow \text{platy kurtic}$$

3. Find skewness and kurtosis from the following data:

C-I	f	M _i	d _i = M _i - A	d _i ²	d _i ³	d _i ⁴	f _i d _i	f _i d _i ²
10-20	4	15	-20	400	-8000	160000	-80	1600
20-30	7	25	-10	100	-1000	10000	-70	700
30-40	16	35	0	0	0	0	0	0
40-50	20	45	10	100	1000	10000	200	2000
50-60	15	55	20	400	8000	160000	300	6000
60-70	8	65	30	900	27000	810000	240	7200
<u>N = 70</u>			<u>30</u>	<u>1900</u>	<u>27000</u>	<u>1150000</u>	<u>590</u>	<u>17500</u>

f _i d _i ³	f _i d _i ⁴
-32000	640000
-7000	10000
0	0
20000	200000
120000	2400000
216000	6480000

$$\mu_r' = \frac{1}{N} \sum f_i d_i r$$

$$\text{If } r=1, \mu_1' = \frac{1}{70} \times 590 = 8.42$$

$$\text{If } r=2, \mu_2' = \frac{1}{70} \times 17500 = 250$$

$$\text{If } r=3, \mu_3' = \frac{1}{70} \times 317000 = 4528.57$$

$$\text{If } r=4, \mu_4' = \frac{1}{70} \times 9190000 = 139857.14$$

$$\mu_1 = 0$$

$$\mu_2 = \mu_2' - (\mu_1')^2 = 250 - (8.42)^2 = 250 - 70.89 = 179.11$$

$$\begin{aligned}\mu_3 &= \mu_3' - 3\mu_2'\mu_1' + 2(\mu_1')^3 \\ &= 4528.57 - 3(250)(8.42) + 2(8.42)^3 \\ &= 4528.57 - 6315 + 1193.89 \\ &= -592.54\end{aligned}$$

$$\begin{aligned}\mu_4 &= \mu_4' - 4\mu_3'\mu_1' + 6\mu_2'(\mu_1')^2 - 3(\mu_1')^4 \\ &= 139857.14 - 4(4528.57)(8.42) + 6(250)(8.42)^2 \\ &\quad - 3(8.42)^4 \\ &= 139857.14 - 152522.23 + 106344.6 - 15078.89 \\ &= 78600.62\end{aligned}$$

$$\beta_1 = \frac{\mu_3'^2}{\mu_2'^3} = \frac{351103.65}{5745919.02} = 0.06 \rightarrow \text{positively skewed.}$$

$$\beta_2 = \frac{\mu_4'}{\mu_2'^2} = \frac{78600.62}{32080.39} = 2.45 \rightarrow \text{platy kurtosis}$$

Curve Fitting:

Suppose that a data is given in the two variables x and y , the problem of finding an analytical expression of the form $y=f(x)$ which fits the given data is called Curve fitting.

(Q)

Curve fitting means an exact relationship between two variables by algebraic equations. This relationship is equation of the curve. Curve fitting means to form an equation of the curve from the given data.

Method of Least Squares

The sum of the squares of the deviation between observed values and expected values should be minimum is called Residual Error (or Method of Least square).

$$E = \sum [y - f(x)]^2 \text{ is minimum}$$

- Some curve fitting equations are given below:

- Straight line: $y = a + bx$
- parabola: $y = a + bx + cx^2$
- power curve: $y = ax^b$
- exponential curve: $y = ab^x$ or $y = a \cdot e^{bx}$

Straight line equation:

$$\text{Let the straight line: } y = a + bx \quad \textcircled{1}$$

Taking summation on both sides of \textcircled{1}, we get

$$\sum y = na + b \sum x \quad \textcircled{2}$$

Taking ' x ' on both sides of \textcircled{1}, we get

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

Solve eq'n \textcircled{2} & \textcircled{3}, we get a & b

Problems:

- Fit a straight line for the following data:

x	y	x^2	xy
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

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

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

$$40 = 8a + 56b$$

$$364 = 56a + 524b$$

$$a = 0.5454,$$

$$b = 0.6363.$$

straight line is

2. Fit a straight line to the following data by using equation: $y = a_0 + a_1 x$

x	y	x^2	xy	$\Sigma y = n a_0 + a_1 \Sigma x$
0	-1	0	0	$\Sigma xy = a_0 \Sigma x + a_1 \Sigma x^2$
2	5	4	10	$36 = 4 a_0 + a_1 (14)$
5	12	25	60	$210 = a_0 (14) + a_1 (78)$
7	20	49	140	
$\Sigma x = 14$		$\Sigma y = 36$		$a_0 = -1.137$
$\Sigma x^2 = 78$		$\Sigma xy = 210$		$a_1 = 2.89$

$$y = a_0 + a_1 x$$

$$y = -1.137 + (2.89)x$$

Fitting of a second degree parabola:

Let the second degree parabola equation:

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

* Normal equations:

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

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

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

Problem:

1. Fit a second degree parabola for the following data:

x	y	x^2	x^3	x^4	Σxy	$\Sigma x^2 y$
0	1	0	0	0	0	0
1	1.8	1	1	1	1.8	1.8
2	2.6	4	8	16	2.6	5.2
3	5.2	9	27	81	7.5	22.5
4	10.8	16	64	256	25.2	100.8
$\Sigma x = 10$		$\Sigma y = 12.9$		$\Sigma x^2 = 30$		$\Sigma x^3 = 100$
$\Sigma x^4 = 354$		$\Sigma xy = 37.1$		$\Sigma x^2 y = 130.3$		

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

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

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

$$12.9 = 5a + b(10) + c(30)$$

$$a = 1.42$$

$$\Rightarrow 37.1 = a(10) + b(30) + c(100)$$

$$b = -1.07$$

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

$$c = 0.55$$

Second degree parabola eq'n is:

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

$$y = 1.42 - 1.07(x) + 0.55(x^2)$$

2. Fit a second degree parabola for the following data:

x	y	x^2	x^3	x^4	xy	x^2y
20	5.5	400	8000	160000	110	2200
40	9.1	1600	64000	2560000	364	14560
60	14.9	3600	216000	12960000	894	53640
80	22.8	6400	512000	40960000	1824	145920
100	33.3	10000	1000000	100000000	3330	333000
120	46	14400	1728000	207360000	5520	662400

$\Sigma x = 420$ $\Sigma y = 131.6$ $\Sigma x^2 = 36400$ $\Sigma x^3 = 3528000$ $\Sigma xy = 12042$ $\Sigma x^2y = 1211720$
 $\Sigma x^4 = 364000000$

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

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

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

$$131.6 = (6)a + b(420) + c(36400)$$

$$\Rightarrow 12042 = a(420) + b(36400) + c(3528000)$$

$$1211720 = a(36400) + b(3528000) + c(364000000)$$

$$a = 4.35$$

$$b = 2.410714 \times 10^{-3}$$

$$c = 2.8705 \times 10^{-3}$$

Second degree parabola eq'n is:

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

$$= (2.410714 \times 10^{-3})x + (2.8705 \times 10^{-3})x^2$$

Power Curve:

The power curve $y = ax^b$

Taking 'log' on both sides

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

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

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

Here, $\log y = Y$

$$\log a = A, a = \text{Antilog}(A)$$

$$\log x = X$$

$$Y = A + bX$$

Normal Equations:

$$\sum Y = nA + b \sum X \quad \dots \textcircled{1}$$

$$\sum XY = A \sum X + b \sum X^2 \quad \dots \textcircled{2}$$

Solve eq'n's $\textcircled{1}$ & $\textcircled{2}$, we get A & b

Problems:

1. Fit a power curve for the following data:

x	y	$X = \log x$	$Y = \log y$	x^2	XY
1	1200	0	3.0791	0	0
2	900	0.3010	2.9542	0.0906	0.8892
3	600	0.4771	2.7181	0.2276	1.3254
4	200	0.6020	2.3010	0.3624	1.3852
5	100	0.6989	2	0.4884	1.3978
6	50	0.7781	1.6989	0.6054	1.3219
$\sum x = 21$		$\sum y = 3050$	$\sum X = 2.8571$	$\sum Y = 14.81137$	$\sum X^2 = 1.7744$
					$\sum XY = 6.3195$

$$\sum Y = nA + b \sum X$$

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

$$\Rightarrow 14.81137 = 6A + 2.8571(b) \quad \dots \textcircled{1}$$

$$6.3195 = 2.8571A + 1.7744(b) \quad \dots \textcircled{2}$$

Solve $\textcircled{1}$ & $\textcircled{2}$, we get

$$A = 3.3123, b = -1.77189$$

$$a = \text{antilog}(A)$$

$$= 2052.57$$

2. Fit a power curve for the following data:

x	y	$X = \log x$	$Y = \log y$	x^2	XY
1	2.98	0	0.4742	0	0
2	4.26	0.3010	0.6294	0.0906	0.1894
3	5.21	0.4771	0.7168	0.2276	0.3419
4	6.1	0.6020	0.7853	0.3624	0.4727
5	6.8	0.6989	0.8325	0.4884	0.5818
6	7.5	0.7781	0.8750	0.6054	0.6808
$\Sigma X = 2.8571$		$\Sigma Y = 4.3132$		$\Sigma X^2 = 1.7744$	
$\Sigma XY = 2.2666$					

$$\Sigma Y = nA + b \Sigma X$$

$$\Sigma XY = A \Sigma X + b \Sigma X^2 \Rightarrow 2.2666 = A(2.8571) + b(1.7744)$$

$$A = 0.4741, b = 0.5139$$

$$a = \text{Antilog}(A)$$

$$= 2.9192$$

$$\text{power curve: } y = ax^b$$

$$y = (2.9192) x^{(0.5139)}$$

Exponential Curve:

$$\text{Let the exponential curve } y = ab^x$$

$$\log y = \log(ab^x)$$

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

$$Y = \log y$$

$$A = \log a$$

$$B = \log b$$

$$Y = A + Bx$$

Normal Equations:

$$\Sigma Y = nA + B\Sigma X \quad \textcircled{1}$$

$$\Sigma XY = A\Sigma X + B\Sigma X^2 \quad \textcircled{2}$$

Solve eq'n $\textcircled{1}$ & $\textcircled{2}$, we get A & B

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

1. Fit an exponential curve of the form: $y = ab^x$ to the following data:

x	y	x^2	$\log y$	xy			
1	1	1	0	0			
2	1.2	4	0.0791	0.1582			
3	1.8	9	0.2552	0.7656			
4	2.5	16	0.3979	1.5916			
5	3.6	25	0.5563	2.7815			
6	4.7	36	0.6720	4.032			
7	6.6	49	0.8195	5.7365			
8	9.1	64	0.9590	7.672			
$\sum x = 36$		$\sum x^2 = 204$		$\sum y = 3.7390$		$\sum xy = 22.7374$	

$$\sum y = nA + B\sum x$$

$$3.7390 = (8)A + B(36)$$

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

$$22.7374 = A(36) + B(204)$$

$$A = -0.1660, B = 0.1407$$

$$a = \text{antilog}(A), b = \text{antilog}(B)$$

$$a = 0.6823$$

$$b = 1.3826$$

Exponential curve: $y = ab^x$

$$y = (0.6823)(1.3826)^x$$

2. Fit an exponential curve of the form: $y = ab^x$ to the following data:

Year (x)	Production (y)	x^2	$\log y$	xy			
1901	3.9	3613801	0.5910	1123.491			
1911	5.3	3651921	0.7242	1383.9462			
1921	7.3	3690241	0.8633	1658.3993			
1931	9.6	3728761	0.9822	1896.6282			
1941	12.9	3767481	1.1105	2155.4805			
1951	17.1	3806401	1.2329	2405.3879			
1961	23.2	3845521	1.3654	2677.5494			
1971	30.5	3884841	1.4842	2925.3582			
$\sum x = 15488$		$\sum x^2 = 29988968$		$\sum y = 8.3537$		$\sum xy = 16226.2407$	

$$\Sigma Y = nA + BE^x$$

$$\Sigma XY = A\Sigma x + BE^x x^2$$

$$8.3537 = (B)A + B(15488)$$
$$16226.2407 = A(15488) + B(29988968)$$

$$A = -23.6063, B = 0.0127$$

$$a = \text{antilog}(A), b = \text{antilog}(B)$$

$$a = 2.4757 \times 10^{-24}, b = 1.0296$$

Exponential curve: $y = ab^x$

$$y = (2.4757 \times 10^{-24}) (1.0296)^x$$

Exponential Curve:

Let the exponential curve $y = ae^{bx}$

Taking log on both sides

$$\log y = \log(a \cdot e^{bx})$$

$$\log y = \log a + \log e^{bx}$$

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

$$[\log_{10} e = 0.4342]$$

$$\log y = Y$$

$$\log a = A$$

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

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

Normal equations:

$$\Sigma Y = nA + BE^x \quad \text{--- ①}$$

$$\Sigma XY = A\Sigma x + BE^x x^2 \quad \text{--- ②}$$

Solve eq'n ① & ②, we get A & B

$$a = \text{Antilog}(A) \quad \& \quad b = \frac{B}{\log_{10} e}$$

$$\Rightarrow b = \frac{B}{0.4342}$$

Problems:

- Fit an exponential curve of the form: $y = ae^{bx}$ to the following data:

x	y	x^2	$Y = \log y$	xy
1	1.6	1	0.2041	0.2041
2	4.5	4	0.6532	1.3064
3	13.8	9	1.1398	3.4194
4	40.2	16	1.6042	6.4168
5	12.5	25	1.0969	5.4845
6	30	36	1.4771	8.8626
$\sum x = 21$	$\sum x^2 = 91$	$\sum Y = 6.1753$		$\sum xy = 25.6938$

$$\sum Y = nA + B\sum x$$

$$\sum xy = A\sum x + B\sum x^2 \Rightarrow 25.6938 = A(21) + B(91)$$

$$A = 0.2131, B = 0.2331$$

$$a = \text{Antilog}(A), b = \frac{0.2331}{0.4342}$$

$$a = 1.6334$$

$$b = 0.5368$$

Exponential curve: $y = ae^{bx}$

$$y = (1.6334)e^{(0.5368)x}$$

2. Fit an exponential curve of the form: $y = ae^{bx}$
to the following data:

x	y	x^2	$Y = \log y$	xy
0	20	0	1.3010	0
1	30	1	1.4771	1.4771
2	52	4	1.7160	3.432
3	71	9	1.8864	5.6592
4	135	16	2.1303	8.5212
5	211	25	2.3242	11.621
6	326	36	2.5132	15.0792
7	550	49	2.7403	19.1821
$\sum x = 28$	$\sum x^2 = 140$	$\sum Y = 16.0885$		$\sum xy = 64.9118$

$$\Sigma Y = nA + B\Sigma x$$

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

$$16.0885 = (8)A + B(28)$$

$$64.9718 = A(28) + B(140)$$

$$A = 1.2892, B = 0.2062$$

$$a = \text{Antilog}(A), b = \frac{0.2062}{0.4342}$$

$$a = 19.4625, b = 0.4748$$

Exponential curve: $y = ae^{bx}$

$$y = (19.4625)e^{(0.4748)x}$$

Correlation:

In a bivariate distribution, the relationship between two variables is called correlation.

If the change in one variable affects the change in the other variable, then the variables are said to be correlated.

Ex: Price and demand
income and expenditure
height and weight of a group of people

Types of correlation:

There are 5 types of correlation. They are as follows:

1. Positive correlation:

In a bivariate distribution, one variable increase (or) decrease, the other variable also increase (or) decrease, then the relation is called positive correlation.

i.e., if $r > 0$

2. Negative correlation:

In a bivariate distribution, one

Negative correlation

i.e., if $r < 0$

3. Nonsense correlation/zero correlation:

The correlation which is no meaning is called Nonsense correlation

(or)

In a bivariate distribution, the two variables x and y are independent, then it is called Nonsense correlation or Zero correlation.

i.e., if $r=0$

4. Perfect positive correlation:

The two variables x and y are said to be perfect positively correlated, if they deviated same direction with constant ratio.

i.e., if $r=1$

Ex:	x	1	3	5	7
	y	2	4	6	8

5. Perfect negative correlation:

The two variables x and y are said to be perfect negatively correlated, if they deviated opposite direction with constant ratio.

i.e., if $r=-1$

Ex:	x	1	3	5	7
	y	8	6	4	2

Methods to measure of correlation:

There are different types of measures of correlation. Some of them are:

1. Karl Pearson's Coefficient of correlation (r)

Karl Pearson's coefficient of correlation:

The numerical measure of intensity or degree of linear relationship between the random variable x and y is called coefficient of correlation between the variables usually denoted by

$$r \text{ or } r(x,y) \text{ or } r_{xy}$$

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

Here,

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

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

\bar{x} = Mean in terms of 'x' & \bar{y} = Mean in terms of 'y'

$$\bar{x} = \frac{1}{n_1} \sum_{i=1}^{n_1} x_i$$

$$\bar{y} = \frac{1}{n_2} \sum_{j=1}^{n_2} y_j$$

Properties of correlation:

- The limits of the Karl Pearson's coefficient of correlation: ± 1
i.e., $-1 \leq r \leq 1$
- The coefficient of correlation is independent of change of origin and scale
- Two independent random variables are uncorrelated but converse is not true

Problems:

1. Find Karl Pearson's coefficient of correlation from the following data:

Height (x)	Weight (y)	$X = x - \bar{x}$	$Y = y - \bar{y}$	X^2	XY	Y^2
100	98	1	3	1	3	9
101	99	2	4	4	8	16
102	99	3	4	9	12	16
102	97	3	2	9	6	4
100	95	1	0	1	0	0
99	92	0	-3	0	0	9
97	95	-2	0	4	0	0
98	94	-1	-1	1	1	1
96	90	-3	-5	9	15	25
<u>95</u>	<u>91</u>	<u>-4</u>	<u>-4</u>	<u>16</u>	<u>16</u>	<u>16</u>
$\Sigma x = 990$	$\Sigma y = 950$	$\Sigma x = 0$		$\Sigma x^2 = 54$	$\Sigma xy = 61$	$\Sigma y^2 = 96$

$$\bar{x} = \frac{990}{10} = 99$$

$$\bar{y} = \frac{950}{10} = 95$$

$$r = \frac{\Sigma xy}{\sqrt{\Sigma x^2} \sqrt{\Sigma y^2}} = \frac{61}{\sqrt{54} \sqrt{96}} = 0.84$$

2. Calculate correlation coefficient from the following data:

x (height in inches)	y	$X = x - \bar{x}$	$Y = y - \bar{y}$	X^2	XY	Y^2
65	67	-3	-2	9	6	4
66	68	-2	-1	4	2	1
67	65	-1	-4	1	4	16
67	68	-1	-1	1	1	1
68	72	0	3	0	0	9
69	72	1	3	1	3	9
70	69	2	0	4	0	0
72	71	4	2	16	8	4
$\Sigma x = 544$	$\Sigma y = 552$			$\Sigma x^2 = 36$	$\Sigma xy = 24$	$\Sigma y^2 = 44$

$$\bar{x} = \frac{544}{8} = 68$$

$$\bar{y} = \frac{552}{8} = 69$$

Spearman's Rank coefficient of correlation:

This method is based on rank and is useful in dealing with qualitative characteristic such as mortality, intelligence, blindness, kindness, goodness etc..

The formula for Spearman Rank coefficient of correlation :

1. observations are not repeated / ranks are not repeated :

$$\rho = 1 - \left[\frac{6 \sum D^2}{N^3 - N} \right]; \text{ where } N = \text{pair of observation}$$

D = difference between two ranks
i.e., $D = R(x) - R(y)$

2. When Ranks are repeated :

$$\rho = 1 - \frac{6 \left(\sum D^2 + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) + \dots \right)}{N^3 - N}$$

Here N = pair of observations.

D = difference between two ranks

m = repeated observation value

Problems:

1. From the ranks obtained by 10 students in 2 subjects statistics and mathematics, to what extent the knowledge of the students in two subjects related:

(X) Statistics	(Y) Mathematics	$D = R(X) - R(Y)$	D^2
1	2	-1	1
2	4	-2	4
3	1	2	4
4	5	-1	1
5	3	2	4
6	9	-3	9
7	7	0	0
8	10	-2	4
9	6	3	9

Spearman's Rank coeff. of correlation $\rho = 1 - \frac{6 \sum D^2}{N^3 - N}$

$$= 1 - \frac{6 \times 40}{10^3 - 10} = 1 - \frac{240}{990} = 0.75$$

∴ Mathematics and statistics are positively correlated.

2. 10 competitors in a musical test were ranked by the 3 judges A, B and C in the following order

Ranks by A	B	C	$[P(A) - P(B)]$	D_1	D_1^2	D_2	D_2^2	D_3	D_3^2
1	3	6	-2	4	-3	9	-5	25	
6	5	4	+1	1	1	1	2	4	
5	8	9	-3	9	-1	1	-4	16	
10	4	8	6	36	-4	16	2	4	
3	7	1	-4	16	6	36	2	4	
2	10	2	-8	64	8	64	0	0	
4	2	3	2	4	-1	1	1	1	
9	1	10	8	64	-9	81	-1	1	
7	6	5	1	1	1	1	2	4	
8	9	7	-1	1	2	4	1	1	
				$\sum D_1^2 = 200$		$\sum D_2^2 = 214$		$\sum D_3^2 = 60$	

Using Rank coeff. of correlation, discuss which pair of judge approach a common liking in music

$$\rho_1 = 1 - \frac{6 \sum D_1^2}{N^3 - N} = 1 - \frac{6 \times 200}{990} = 1 - \frac{120}{99} = 1 - 1.21 = -0.21$$

$$\rho_2 = 1 - \frac{6 \sum D_2^2}{N^3 - N} = 1 - \frac{6 \times 214}{990} = 1 - \frac{1284}{990} = 1 - 1.29 = -0.29$$

$$\rho_3 = 1 - \frac{6 \sum D_3^2}{N^3 - N} = 1 - \frac{6 \times 60}{990} = 1 - \frac{36}{99} = 1 - 0.36 = 0.64$$

3. Find Rank coefficient of correlation to the following data:

X	Y	$D = S(X) - S(Y)$	D^2	$\rho = 1 - \frac{6\sum D^2}{N^3 - N}$
85	93	1	1	
60	75	1	1	$= 1 - \frac{6 \times 4}{125 - 5}$
73	65	-1	1	$= 1 - 0.2$
40	50	0	0	$\rho = 0.8$
90	80	-1	1	
			<u>$\sum D^2 = 4$</u>	

4. Obtain the rank coefficient of correlation for the following data:

X	Y	$S(X)$	$S(Y)$	$S(X) - S(Y)$	D	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	
					<u>$\sum D^2 = 72$</u>	

In X-series, 75 is repeated 2-times $\therefore m=2$

" , 64 " 3 " $\therefore m=3$

In Y-series, 68 " 2 " $\therefore m=2$

$$\rho = 1 - \frac{6 \left(\sum D^2 + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) \right)}{N^3 - N}$$

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

$$6 (72 + 0.5 + 2 + 0.5) \quad 6 (75)$$

5. Obtain the rank coefficient of correlation for the following data:

X	Y	P(X)	S(Y)	D	ΣD^2
48	13	2	6.5	-4.5	20.25
33	13	4	6.5	-2.5	6.25
40	24	3	1	2	4
9	6	10	9	1	1
16	15	8	5	3	9
25	20	5	2	3	9
24	9	6	8	-2	4
16	16	8	4	4	16
51	19	1	3	-2	4
16	4	8	10	-2	4
<hr/>					
$\frac{7+8+9}{3} = 8$	$\frac{6+7}{2} = 6.5$			$\Sigma D^2 = 77.5$	

In X-series, 16 repeated 3 times $\Rightarrow m=3$

In Y-series, 13 " 2 times $\Rightarrow m=2$

$$r = 1 - \frac{6(\Sigma D^2 + \frac{1}{12}(m^3-m) + \frac{1}{12}(m^3-m))}{N^3-N}$$

$$= 1 - \frac{6(77.5 + \frac{1}{12}(3^3-3) + \frac{1}{12}(2^3-2))}{10^3-10}$$

$$= 1 - \frac{6(77.5 + 2 + 0.5)}{990}$$

$$= 1 - 0.080$$

$$= 0.92$$

Regression Analysis:

Regression Analysis in general sense the estimation of prediction of the unknown value of one variable from the known value of other variable. It is very important statistical tool which is extensively used in almost all sciences - natural, social and physical science.

Definition:

Regression Analysis is a mathematical measure of the average relationship between two or more variables in terms of original units of data.

In Regression Analysis, there are two types of variables. The variable whose value is influenced / is to be predicted is called "dependent variable" and the variable which influences the values is used for prediction is called "Independent variable".

Lines of Regression:

are Lines of Regression (of) the line which gives the best estimate of the value of one variable for any given value of the other variable. In case of two variables 'X' and 'Y' we shall have two lines of Regression. They are

1. Regression Line X on Y
2. Regression Line Y on X

$$1 \text{ Regression Line } X \text{ on } Y = X - \bar{X} = Y \frac{\sigma_x}{\sigma_y} (Y - \bar{Y})$$

(or)

$$X - \bar{X} = b_{XY} (Y - \bar{Y})$$

(or)

Here $x = \bar{x} - \frac{y - \bar{y}}{r}$ mean in terms of x

$y = \bar{y} + r \frac{x - \bar{x}}{r}$ mean in terms of y

$r \frac{\sigma_x}{\sigma_y}$ on b_{xy} on $\frac{\sum xy}{\sum y^2}$ = Regression coefficient of x on y

2. Regression line y on $x = Y - \bar{Y} = r \cdot \frac{\sigma_y}{\sigma_x} (x - \bar{x})$

$r \cdot \frac{\sigma_y}{\sigma_x}$ on b_{yx} on $\frac{\sum xy}{\sum x^2}$ = Regression coefficient of y on x

3. Coefficient of correlation: $r = \sqrt{b_{xy} \cdot b_{yx}}$ on $\sqrt{\frac{\sum xy}{\sum y^2} \times \frac{\sum xy}{\sum x^2}}$

4. Probable error = P.E = $0.6745 \times \frac{1-r^2}{\sqrt{N}}$

Properties of Regression Coefficient:

1. The correlation coefficient is geometric mean of two regression coefficients.

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

2. The arithmetic mean of two regression coefficients is greater than or equal to correlation coefficient

$$\frac{b_{xy} + b_{yx}}{2} \geq r$$

3. If the coefficient of correlation cannot exceed 1 in case of regression, one of the regression coefficient is greater than 1, then the other must be less than 1.

4. Both the regression coefficient will have the same sign either positive or negative.

5. If one regression coefficient is '+' then the other should be '-' and vice versa

Uses of Regression Analysis:

- 1) The Regression Analysis technique is very useful in predicting the probable value of an unknown variable in response to some known related variable.
- 2) The Regression Analysis is useful in establishing the nature of relationship between 2 variables.
- 3) Regression Analysis is extensively used for measurement and estimating the relationship among variables.
- 4) Regression Analysis provides regression coefficient which are generally used in calculation of correlation coefficient.

Problems:

1 Calculate the Regression coefficient of X on Y and Y on X and also find probable error from the following data:

price (x)	Demand (y)	$x = x - \bar{x}$	x^2	$y = y - \bar{y}$	y^2	xy	
10	40	-3	9	-1	1	3	
12	38	-1	1	-3	9	3	
13	43	0	0	2	4	0	
12	45	-1	1	4	16	-4	
16	37	3	9	-4	16	-12	
15	43	2	4	2	4	4	
$\sum x^2 = 24$			$\sum y^2 = 50$			$\sum xy = -6$	

Estimate the likely demand when the price is £ 20

$$\bar{x} = \frac{10+12+13+12+16+15}{6} ; \bar{y} = \frac{40+38+43+45+37+43}{6}$$

$$\bar{x} = 13$$

$$\bar{y} = 41$$

Regression line X on Y:

$$x = a + b y \quad (x = \bar{x} + \bar{x}^r y)$$

Regression line Y on X:

$$y = a + b x \quad (y = \bar{y} + \bar{y}^r x)$$