

INTRODUCTION OF STATISTICS UNIT-2

Statistics:

- Statistic is the science of collection, organization, presentation, analysis and interpretation of numerical data.
- Father of Modern statistic is "Prof. Ronald Alfred Fisher."
- Father of Indian statistic is "Prof. Prasanta Chandra Mahalanobis."

Data:

It means numerical information about a particular inquiry.

Measures of central tendency / Average:

It is a single value which represents the entire distribution. In other words it gives the idea about the central part of the distribution.

Characteristics of Average:

- It should be rigidly defined.
- It should be based on all the observation.
- It should be capable of further mathematical treatment. → It is easy to calculate and understand.
- It should not be affected by extreme items.

Various Measures of Central Tendency:

There are different types of average. The most commonly used average are

1. Arithmetic Mean (A.M)

2. Median

3. Mode

4. Geometric Mean (G.M)

5. Harmonic Mean (H.M)

Variable:-

The quantitative phenomenon under study like marks in a test, height, weight, numbers of student in a class etc. is known as variable.

→ There are two types of variables.

1. Discrete Variable

2. Continuous Variable

Discrete Variable:-

Those variables which can not take all possible values within a given specified range is known as discrete variable.

Ex:- Marks in a test; number of student in a class

Continuous Variable:-

Those Variables which can take all possible values within a given specified range is known as continuous variable.

Ex:- Age, height, weight etc.

Frequency:-

The number of times an item occurs in a series.

Ex:- 18, 19, 20, 18, 19, 30

Age

Frequency

18

2

19

2

20

1

30

1

(M.A) mean distribution

(M.H) median distribution

(M.M) mode distribution

(M.F) mid range distribution

Frequency Distribution:

- The primary data are generally in raw form and are produced in an unorganized manner and it leads to be organized and processed.
- Frequency distribution is one such technique to process the raw data.
- Frequency distribution is a least table or graph that displays the frequency of various outcomes in a sample.
- Each entry in the table contains the frequency of the occurrences of values within a particular group or interval.
- There are different types of frequency distribution.
 1. Discrete / ungrouped frequency distribution
 2. Grouped frequency distribution
 3. Continuous frequency distribution

1. Discrete / ungrouped frequency Distribution:

In discrete frequency distribution the values of the variable are determined individually.

Ex:— Age frequency

18 1
19 2
20 3
21 4
22 5
23 6
24 7
25 8
26 9
27 10
28 11
29 12
30 13

2. Grouped frequency Distribution:

In grouped frequency distribution the data are classified in to different class intervals with gaps and their respective frequency are assigned as the class intervals.

Ex:-

<u>Age</u>	<u>No. of Student</u>
1-9	6
10-19	5
20-29	3
30-39	4

3. Continuous Frequency Distribution:-

In Continuous frequency Distribution, the data are classified into different class intervals without gaps and their respective frequency are assigned as the class intervals.

Ex:-

Age

No. of Student

10-20

10

20-30

5

30-40

3

40-50

4

Class Interval:-

- Class interval is a numerical width of a class in a frequency distribution each class is specified by two extreme values is called as class limits.
- The smaller one is known as Lower limit & larger one is known as upper limit.
- There are two types of class interval.

1. Inclusive

2. Exclusive

1. Inclusive:-

The classes of the type like 1-9, 10-19, 20-29 etc. in which both lower and upper limits are included

in the class are known as Inclusive class interval.

2. Exclusive:—

The classes of the type like 10-20, 20-30, 30-40 etc. in which the lower limit is included and the upper limit is excluded from the respective classes and immediately included in the next class is known as exclusive.

Mid value:—

The mid value of any class is obtained on dividing the sum of the upper and lower limits.

$$\text{Mid value} = \frac{\text{Lower Limit} + \text{Upper Limit}}{2}$$

Ex:—

Age	Mid value
10-20	15
20-30	25
30-40	35

Class Boundaries:—

→ In Inclusive type classification or grouped frequency distribution there are gap between the upper limit of any class and lower limit of the succeeding class. So, there is need to convert the data into a continuous frequency Distribution.

→ The upper and lower limit of the new exclusive type classes are called as Class Boundaries.

→ If d is the gap between the upper limit of the any class & lower limit of the succeeding class, the class boundaries of any class is given by

Upper class boundary = Upper class limit + $\frac{d}{2}$

Lower class boundary = Lower class limit - $\frac{d}{2}$

($\frac{d}{2}$ is the correction factor)

Ex:— Age New class interval

$$1-9 \quad 1-0.5 - 9+0.5 = 0.5 - 9.5$$

$$10-19 \quad 10-0.5 - 19+0.5 = 9.5 - 19.5$$

$$20-29 \quad 20-0.5 - 29+0.5 = 19.5 - 29.5$$

$$30-39 \quad 30-0.5 - 39+0.5 = 29.5 - 39.5$$

Cumulative frequency Distribution:—

Cumulative frequency Distribution is obtained on successively adding the frequencies the values of the variables.

Ex:—

Age

Frequency

cumulative frequency (c.f.)

0-10

5

5

0P-0P

10-20

3

8

0P-0P

20-30

2

10

0P-0P

30-40

1

11

0P-0P

40-50

2

13

0P-0P

Arithmetic Mean (A.M.):—

→ It is defined by sum of the observations divided by the total no. of observations.

→ Generally A.M is denoted as (\bar{x})

so, Mathematically = we can write, if a series consists of n observation, say x_1, x_2, \dots, x_n ,

→ The Arithmetic Mean (\bar{x}) is given by

$$\bar{x} = \frac{\text{Sum of the observation}}{\text{Total no. of observation}} = \frac{x_1 + x_2 + \dots + x_n}{n}$$

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

→ In case of Discrete Frequency Distribution: If a series consist of n observations say x_1, x_2, \dots, x_n & the corresponding frequency is f_1, f_2, \dots, f_n .

The Arithmetic Mean is given by

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

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

$$\boxed{\bar{x} = \frac{\sum_{i=1}^n f_i x_i}{N}}$$

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

→ In Case of Continuous Frequency Distribution x is taken as the mid value of corresponding classes.

Example - 1 :— (Discrete Series)

Calculate the Arithmetic Mean the wages of 100 workers of a firm from the following frequency distribution.

2.00	2.04	2.08 - 2.12	Ed	32 - 1P
2.08	2.12	2.16 - 2.20	7P	33 - 12
2.12	2.16	2.20 - 2.24	3 FG	34 - 13
2.16	2.20	2.24 - 2.28	2 FG	35 - 1F
2.20	2.24	2.28 - 2.32	1 FG	36 - 1

<u>Wages (in RS.) (X)</u>	<u>No. of workers frequency (F)</u>	<u>$\sum FX$</u>
80	5	400
82	8	656
85	17	1445
90	30	2700
95	18	1710
98	11	1078
100	8	800
110	3	330
<u>$\sum F = 100$</u>		<u>$\sum FX = 9119$</u>

$$\bar{X} = \frac{\sum FX}{\sum F} = \frac{9119}{100} = \text{RS. } 91.19$$

Example-2: — (Continuous Series)

The marks of 200 school student is given. In the following frequency distribution calculate the average mark.

<u>Marks</u>	<u>No. of Students (F)</u>	<u>New Marks</u>	<u>Mid value (X)</u>	<u>FX</u>
11-20	8	10.5 - 20.5	15.5	124
21-30	12	20.5 - 30.5	25.5	306
31-40	35	30.5 - 40.5	35.5	1242.5
41-50	63	40.5 - 50.5	45.5	2866.5
51-60	45	50.5 - 60.5	55.5	2497.5
61-70	27	60.5 - 70.5	65.5	1768.5
71-80	10	70.5 - 80.5	75.5	755
<u>$\sum F = 200$</u>				<u>$\sum FX = 9560$</u>

$$\bar{X} = \frac{\sum f_x}{\sum f} = \frac{9560}{200} = 47.8 \text{ Marks}$$

Example-3: —

Calculate the Arithmetic mean of Marks of 10 students are 50, 60, 52, 70, 57, 52, 53, 63, 62, 70.

Ans: — $\bar{X} = \frac{\text{sum of observation}}{\text{Total no. of observation}}$

$$\begin{aligned} &= \frac{50+60+52+70+57+52+53+63+62+70}{10} \\ &= \frac{589}{10} = 58.9 \end{aligned}$$

Change of origin and scale: —

→ When the values of the variable x in the data are large or the value involves decimal figures, the computation of arithmetic mean becomes time consuming and tedious (difficult).

→ So, the computation can be made simple by setting the origin to a convenient point and reducing the magnitude of the figures by changing the Shortcut method or assumed Mean Method or Step deviation method.

→ It is given by

$$\bar{X} = A + \bar{h}d$$

$$\text{where, } \bar{A} = \frac{\sum fd}{\sum f}$$

$$\text{where, } d = \frac{x-A}{h}$$

(Individual Series)

A is any arbitrary value
 h is the magnitude of the Class Interval

Note:-1 (Continuous series)

If we take. $d = x - A$

so, A.M is given by $\bar{x} = A + \bar{d}$

then, $\bar{d} = \frac{\sum d}{n}$, ($n = \text{no. of observation}$).

Example :-

Find the average of the salary 210 persons of a firm given in the following frequency distribution.

Monthly
salary in Rs.

No. of
Person (f)

Mid
value (x)

$d = \frac{x - A}{h}$

fd

3000 - 3999

15

3499.5

-3

-45

4000 - 4999

13

4499.5

-2

-26

5000 - 5999

46

5499.5

-1

-46

6000 - 6999

67

6499.5

0

0

7000 - 7999

32

7499.5

1

32

8000 - 8999

23

8499.5

2

46

9000 - 9999

14

9499.5

3

42

$\sum fd = 33$

$A = 6499.5$

$h = 1000$

$$\bar{x} = A + h \cdot \frac{\sum fd}{\sum f}$$

$$= 6499.5 + 1000 \cdot \frac{3}{210}$$

$$= 6513.78$$

Discrete Series:

$$d = x - A$$

$$\bar{x} = A + \frac{\sum fd}{\sum f}$$

$$N = \sum f$$

$$d = x - A$$

$$fd + A = x$$

$$\frac{fd}{N} + A = \bar{x}$$

$$\frac{fd}{N} = b$$

Properties of arithmetic mean:-

① The sum of the deviations of set of values from their arithmetic mean is set to zero.

So, mathematically we can write as

$$\sum f_i(x_i - \bar{x}) = 0$$

$$\text{OR } \sum f_i(x_i - \bar{x}) = 0$$

Proof:-

$$\begin{aligned} \text{L.H.S.} &= \sum f_i(x_i - \bar{x}) \\ &= \sum f_i x_i - \bar{x} \sum f_i \\ &= \bar{x} \sum f_i - \bar{x} \sum f_i \\ &= N\bar{x} - N\bar{x} \\ &= 0 \quad \text{R.H.S.} \end{aligned}$$

$$\boxed{\begin{aligned} \bar{x} &= \frac{\sum f_i x_i}{\sum f_i} \\ \Rightarrow \sum f_i \bar{x} &= \sum f_i x_i \\ (\bar{N}\bar{x}) & \end{aligned}}$$

② Let \bar{x}_1 be the arithmetic mean of a set of data having small n_1 observations & \bar{x}_2 be the arithmetic mean of another set of data having n_2 observations the arithmetic mean of the combined data of $(n_1 + n_2)$ observations is given by

$$\boxed{\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}}$$

Example:-

① In a class of 2 sections A & B the average of the marks in statistics of the students of section A is 62 & of section B is 60 if the no. of student of section A & B is 80 & 70 find the average marks of all 150 students taken together.

Ans:- $\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$

$$\begin{aligned} &= \frac{62 \times 80 + 60 \times 70}{150} = \frac{4960 + 4200}{150} = \frac{9160}{150} = 61.066 \end{aligned}$$

Q) The average weekly wages of all the workers in a factory is Rs - 600 & the average weekly wages of male & female workers are Rs - 630 & Rs - 530. Find the percentages of male & female workers in a factory.

Ans:— Given $\bar{X} = 600$

$$\bar{x}_1 = 630$$

$$\bar{x}_2 = 530$$

$$\bar{X} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$$

$$\Rightarrow 600 = \frac{n_1 \times 630 + n_2 \times 530}{n_1 + n_2}$$

$$\Rightarrow 600n_1 + 600n_2 = 630n_1 + 530n_2$$

$$\Rightarrow 600n_1 - 630n_1 = 530n_2 - 600n_2$$

$$\Rightarrow -30n_1 = -70n_2$$

$$\Rightarrow \frac{n_1}{n_2} = \frac{70}{30} = \frac{7}{3}$$

$$\% \text{ of male workers} = \frac{7}{7+3} \times 100 = 70\%$$

$$\% \text{ of female workers} = \frac{3}{7+3} \times 100 = 30\%$$

Proof:

We know that if \bar{X} is the mean of n observations,

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

$$\Rightarrow \sum x = n\bar{X} \quad (\text{i.e. the sum of } n \text{ observations} = n \times \text{A.M})$$

If \bar{x}_1 is the mean of n_1 observations of the 1st group.

The sum of n_1 observations in the 1st group = $n_1 \bar{x}_1$

Similarly, \bar{x}_2 is the mean of n_2 observations of the 2nd group.

The sum of n_1+n_2 observations of the combine group = $n_1 \bar{x}_1 + n_2 \bar{x}_2$

Hence, the mean \bar{x} of the combine group n_1+n_2 observation is given by

$$\bar{x} = \frac{\text{sum of } (n_1+n_2) \text{ observation}}{n_1+n_2}$$

$$\boxed{\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1+n_2}}$$

(proved)

③ The sum of the squares of deviation of the given set of observations is minimum, when taken from Arithmetic mean.

Mathematically, we can write for a giving frequency distribution the sum,

$$(b) S = \sum f(x-A)^2$$

which represents the sum of the squares of deviation of given observation from any arbitrary value 'A' is minimum when we $A = \bar{x}$

PROOF (By using the principle of maxima and minima in differential calculus).

For S to be minimum,

$$\frac{dS}{dA} = 0 \quad \& \quad \frac{d^2S}{dA^2} > 0$$

But if $\frac{d^2S}{dA^2} < 0$ then it is said to be a point of inflection.

Differentiating ' S ' w.r.t A , we get

$$\begin{aligned}\frac{ds}{da} &= \frac{d}{da} [\Sigma f(x-a)^2] \\&= 2\Sigma f(x-a) \cdot \frac{d}{da}(x-a) \\&= 2\Sigma f(x-a)(-1) \\&= -2\Sigma f(x-a)\end{aligned}$$

$$\frac{ds}{da} = 0 \Rightarrow \Sigma f(x-a) = 0$$

$$\Rightarrow \Sigma fx - \Sigma fa = 0$$

$$\Rightarrow \Sigma fx - A\Sigma f = 0$$

$$\Rightarrow \Sigma fx - AN = 0$$

$$\Rightarrow \Sigma fx = NA$$

$$\Rightarrow \frac{\Sigma fx}{N} = A$$

$$\Rightarrow \boxed{\bar{x} = A}$$

Differentiating eqn ① w.r.t, A , we get

$$\frac{d^2S}{da^2} = 2\Sigma f$$

$$= 2N > 0 \quad (\text{It is always positive})$$

Weighted Arithmetic Mean: $\bar{x} = \frac{\sum fx}{\sum f}$ (Proved)

→ For computing Arithmetic mean are based on the assumption that all the items in the distribution are of equal importance. But in real life it can not happen.

→ In such cases Proper weightage is should be given to various items. The weights attach to each items being proportional to the important of the item in the distribution.

→ Let $w_1, w_2, w_3, \dots, w_n$ be the weights attached to variables say $x_1, x_2, x_3, \dots, x_n$.

The weighted arithmetic mean usually denoted by \bar{x}_w is given by

$$\bar{x}_w = \frac{w_1 x_1 + w_2 x_2 + \dots + w_n x_n}{w_1 + w_2 + \dots + w_n}$$

$$= \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}$$

Example:-

- ① A candidate obtained the following percentages of marks in an examination.

English: 63, Hindi: 75, Math: 63, Phy: 59, Chem: 55
Find the weighted arithmetic mean if weights 1, 2, 1, 3 & 3 respectively are allocated in the subject.

<u>Sub</u>	<u>x (% of marks)</u>	<u>weights (w)</u>	<u>$w \cdot x$</u>
Eng	60	1	60
Hin	75	2	150
Math	63	1	63
Phy	59	3	177
Chem	55	3	165

$$\bar{x}_w = \frac{\sum w_i x_i}{\sum w_i} = \frac{615}{10} = 61.5$$

- ② The nos. 3.2, 5.8, 7.9, 4.5 & have frequency $x, x+2, x-3, x+6$ if the arithmetic mean is 4.876.50, find the value of x .

<u>No. (x)</u>	<u>f</u>	<u>fx</u>
3.2	x	$3.2x$
5.8	$x+2$	$5.8x + 11.6$
7.9	$x-3$	$7.9x - 23.7$
4.5	$x+6$	$4.5x + 27$
	$\sum f = 4x + 5$	

$$\bar{x} = \frac{\sum fx}{\sum f} = \frac{3.2x + 5.8x + 11.6 + 7.9x - 23.7 + 4.5x + 27}{4x + 5}$$

$$\Rightarrow 4.876 = \frac{21.4x + 14.9}{4x + 5}$$

$$\Rightarrow 19.504x + 24.38 = 21.4x + 14.9$$

$$\Rightarrow 19.504x - 21.4x = 14.9 - 24.38$$

$$\Rightarrow 1.896x = -9.48$$

$$\Rightarrow x = \frac{-9.48}{1.896} = 5$$

Merits of Arithmetic Mean:

- (i) It should be rigidly defined.
- (ii) It should be easy to understand and simple to calculate.
- (iii) It should be based on all the observations.
- (iv) It should be suitable for further mathematical treatment.
- (v) It should not be affected much by extreme observations.

Demerits of Arithmetic Mean:-

- (i) The strongest drawback of arithmetic mean is that it is very much affected by extreme observations. Two or three very large values of the variable may unduly affect the value of the arithmetic mean.
- (ii) It can not be determined by inspection nor can it be located graphically.
- (iii) Arithmetic mean can not be used if we are dealing with qualitative characteristics which cannot be measured quantitatively such as intelligence, honesty etc.
- (iv) Arithmetic mean cannot be obtained if a single observation is missing or lost or is illegible unless we drop it out and compute the arithmetic means of the remaining values.
- (v) In extremely asymmetrical distribution, usually arithmetic mean is not representative of the distribution and hence is not a suitable measure of location.
- (vi) Arithmetic mean may lead to wrong conclusion if the details of the data from which it is obtained are not available. In this connection it is worthwhile to quote the words of H. Sechrist:

"If an average is taken as a substitute for the details, then the arithmetic mean, in spite of the simplicity and ease of calculation, has little to recommend when series are non-homogeneous."

To understand a heterogeneous series
is to rob it of its meaning

$$a = \bar{x}$$

Ex. 0.01, 0.0101, 0.010101

Median:-

Median is the middle most observation which divide the series into two equal parts. It is the value such that the no. of observations above it and no. of observations below it. The median is only positional average that is its value depends on the position occupied by a value in the frequency distribution.

Calculation of Median:- (Individual series)

- IF the no. of observation is odd then the median is $(n+1/2)^{\text{th}}$ item.
- IF the no. of observation is even then the median is obtain by the A.M of two middle observations after they are arranged in either ascending or descending order of magnitude.

Ex:-

Find the median of 5 observations 35, 12, 40, 8,

60.

Descending Order of 5 observations

8, 12, 35, 40, 60.

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

$$= \left(\frac{5+1}{2}\right)^{\text{th}}$$

$$= \left(\frac{6}{2}\right)^{\text{th}}$$

$$= 3^{\text{rd}} \text{ item}$$

Find the median of 6 observations 60, 40, 12, 8, 35, 50

$n=6$

Descending order of 6 observations

8, 12, 35, 40, 50, 60

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

= $\frac{(\frac{6}{2})^{\text{th}} \text{ item} + (\frac{6+1}{2})^{\text{th}} \text{ item}}{2}$

= $\frac{3^{\text{rd}} \text{ item} + 4^{\text{th}} \text{ item}}{2}$

$$= \frac{35 + 40}{2}$$

$$= \frac{75}{2} = 37.5$$

∴ Median is 37.5.

* IF the no. of observations is odd then median is $(\frac{n+1}{2})^{\text{th}}$ item. After the observations having arranged in ascending or descending order of magnitude.

* In case of even no. of observations median is the arithmetic mean of $(\frac{n}{2})^{\text{th}}$ observation & $(\frac{n+1}{2})^{\text{th}}$ observation.

Discrete series: → If the series consist of 'n' observations x_1, x_2, \dots, x_n , the corresponding frequency are f_1, f_2, \dots, f_n $N = \sum f_i$ the median is the size of the capital $(\frac{N+1}{2})^{\text{th}}$ item.

→ In case if you use cumulative frequency distribution the steps for calculating median are
(1) Prepare less than cumulative frequency distribution.
(2) find $\frac{N}{2}$
(3) See the cumulative frequency just greater than $\frac{N}{2}$.

(iv) The corresponding value of the variable

gives median.

Ex:- calculate median.

8 coins were tossed together and the no. of heads(x) resulting was noted the operation was repeated 256 times and the frequency distribution of no. of head is given below.

<u>No. of heads(x)</u>	<u>Frequency (F)</u>	<u>Less than C.F</u>
0	1	1
1	9	10
2	26	36
3	59	95
4	72	167
5	52	219
6	29	248
7	7	255
8	1	256

$$N = \sum F = 256$$

$$\frac{N}{2} = \frac{256}{2} = 128$$

The cumulative frequency is just greater than 128 is 167. The value of x corresponding to 167 is 4. Hence, the median no. of x is 4.

Continuous Series :-

(i) Prepare 'less than' cumulative frequency distribution.

(ii) Find $\frac{N}{2}$.

(iii) See the cumulative frequency just greater than $\frac{N}{2}$.

(iv) The corresponding class contains the median value it is called the median class. The median is given by

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

where,

l is the lower limit of the median class,

f is the frequency of the median class,

h is the magnitude or width of the median class,

$N = \sum f$, is the total frequency,

c is the cumulative frequency of the class preceding the median class.

Ex:-

Find the missing frequency from the following distribution of daily sales of shops, given that the median sale of shops is ₹2,400.

Sales in hundred ₹.	No. of shops (f)	cumulative frequency
0-10	5	5
10-20	25	30
20-30	a	30+a
30-40	18	48+a
40-50	7	55+a

Let the missing frequency be 'a'.

Given, the median sale is 2400.

20-30 is the median class.

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

$$20 + \frac{10}{10} (50 - 30) = 24$$

where, $l = 20$, $h = 10$, $f = a$, $\frac{N}{2} = \frac{55+a}{2}$, $C = 30$.

$$\text{Median} = 20 + \frac{10}{a} \left(\frac{55+a}{2} - 30 \right)$$

$$\Rightarrow 24 = 20 + \frac{10}{a} \left(\frac{55+a-60}{2} \right)$$

$$\Rightarrow 24 = 20 + \frac{10}{a} \left(\frac{a-5}{2} \right)$$

$$\Rightarrow 24 = 20 + \frac{10a-50}{2a}$$

$$\Rightarrow 24 = \frac{40a+10a-50}{2a}$$

$$\Rightarrow 24 \times 2a = 50a - 50$$

$$\Rightarrow 48a = 50a - 50$$

$$\Rightarrow 48a - 50a = -50$$

$$\Rightarrow -2a = -50$$

$$\Rightarrow a = \frac{-50}{-2} = 25$$

\therefore Hence, the missing frequency is 25.

Ex:-

In the frequency distribution of 100 families given below the number of families corresponding to expenditure groups 20-40 & 60-80 are missing from the table. However, the median is known to be 50. find the missing frequencies.

Expenditure (in Rupees)	No. of families (f)	cumulative frequency (C.F) (less than)
0-20	14	14
20-30	f_1	$14 + f_1$
30-40	27	$41 + f_1$
40-60	f_2	$41 + f_1 + f_2$
60-80		$41 + f_1 + f_2$
80-100	$15 - (f_1 + f_2)$	$56 + f_1 + f_2$
$N = 100 = 56 + f_1 + f_2$		

Given; $N=100$

$$N = 56 + f_1 + f_2$$

$$\Rightarrow 100 = 56 + f_1 + f_2$$

$$\Rightarrow f_1 + f_2 = 100 - 56$$

$$\Rightarrow f_1 + f_2 = 44 \quad \text{--- (1)}$$

Since median is given to be 50, which lies in the class 40-60, therefore, 40-60 is the median class.

where, $l=40$, $h=20$, $f=27$, $c=14+f_1$

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

$$\Rightarrow 50 = 40 + \frac{20}{27} [50 - (14 + f_1)]$$

$$\Rightarrow 50 - 40 = \frac{20}{27} (36 - f_1)$$

$$\Rightarrow 10 = \frac{20}{27} (36 - f_1)$$

$$\Rightarrow 270 = 720 - 20f_1$$

$$\Rightarrow 270 - 720 = -20f_1$$

$$\Rightarrow -450 = -20f_1$$

$$\Rightarrow f_1 = \frac{-450}{-20} = \frac{45}{2} = 22.5 \approx 23$$

$$f_2 = 44 - f_1 = 44 - 23 = 21$$

Merits: ——————
It is rigidly defined.

(ii) Median is easy to understand and easy to calculate for a non-mathematical person.

(iii) Since median is a positional average, it is not affected at all by extreme observations and such is very useful in the case of skewed distributions, J-shaped or inverted J-shaped distribution such as the distribution of wages, incomes & wealth. So in case

- of extreme observations, median is a better average to use than the arithmetic mean since the latter gives a distorted picture of the distribution.
- (i) Median can be computed while dealing with a distribution with open end classes.
- (ii) Median can sometimes be located by simple inspection and can also be computed graphically.

Demerits:

- (i) In case of even number of observations for an ungrouped data, median cannot be determined exactly. We merely estimate it as the arithmetic mean of the two middle terms. In fact any value lying between the two middle observations can serve the purpose of median.
- (ii) Median, being a positional average, is not based on each and every item of the distribution depends on all the observations only to the extent whether they are smaller than or greater than it; the exact magnitude of the observations being immaterial.
- (iii) Median is not suitable for further mathematical treatment i.e. given the sizes and the median values of different groups, we cannot compute the median of the combined group.
- (iv) Median is relatively less stable than mean, particularly for small samples since it is affected more by fluctuations of sampling as compared with arithmetic mean.

Mode :-

Mode is the value of a series which is predominant in it. In other words mode is the value which occurs most frequently in a set of observations. In case of discrete series mode is the value of variable corresponding to the maximum frequency.

<u>Age</u>	<u>frequency</u>
15	20
16	30
17	35 → Maximum frequency
18	15

mode is 17.

→ In case of continuous series mode is defined by

$$\text{Mode} = l + \frac{h(f_1 - f_0)}{2f_1 - f_0 - f_2}$$

The class corresponding to the maximum frequency is called as modal class.

where 'l' is the lower limit of the class

'f₁' is the frequency of the modal class

'f₀' is the frequency of class preceding the modal class

'f₂' is the frequency of the class succeeding the modal class

'h' is the magnitude of modal class.

calculate mode.

<u>(1)</u>	<u>X</u>	<u>F</u>
	1	3
	2	1
	3	18
	4	25
	5	40
	6	30

Here maximum frequency is 40
the value corresponding to maximum frequency is 5.
Hence, mode is 5.

7	22
8	10
9	6

Q2

<u>Class</u>	<u>F</u>
0 - 10	5
10 - 20	12
20 - 30	25
30 - 40	10
40 - 50	9
50 - 60	8

Here the maximum frequency is 25, the class corresponding to maximum frequency is 20-30. The modal class is 20-30.

$$\text{Where, } l = 20$$

$$h = 10$$

$$f_1 = 25$$

$$f_0 = 12$$

$$f_2 = 10$$

$$\text{Mode} = l + \frac{h(f_1 - f_0)}{2f_1 - f_0 - f_2}$$

$$= 20 + \frac{10(25 - 12)}{2 \times 25 - 12 - 10}$$

$$= 20 + \frac{130}{28}$$

$$\therefore \text{Mode} = \frac{560 + 130}{28}$$

$$= \frac{690}{28}$$

$$\therefore \text{Mode} = 24$$

\therefore Mode is 24.

Merits:-

- (i) Mode is easy to calculate & understand. In some cases it can be located merely by inspection. It can also be estimated graphically from a histogram.
- (ii) Mode is not at all affected by extreme observations and as such is preferred to arithmetic mean while dealing with extreme observations.
- (iii) It can be conveniently obtained in the case of open end classes which do not pose any problems here.

Demerits:-

- (i) Mode is not rigidly defined.
- (ii) Since mode is the value of x corresponding to the maximum frequency, it is not based on all the obs' of the series. Even in the case of the continuous frequency distribution mode depends on the frequencies of modal class and the classes preceding and succeeding it.
- (iii) Mode is not suitable for further mathematical treatment.
- (iv) As compared with mean, mode is affected to a greater extent by the fluctuations of sampling.

Geometric Mean:-

If the series consists of ' n ' observations say x_1, x_2, \dots, x_n then geometric mean is defined as n th root of their product.

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

Taking logarithms both sides, we get

$$\log G.M = \log [(x_1, x_2, \dots, x_n)^{1/n}]$$

$$\rightarrow \log G.M = \frac{1}{n} [\log x_1 + \log x_2 + \dots + \log x_n]$$

$$\rightarrow \log G.M = \frac{1}{n} \sum_{i=1}^n \log x_i$$

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

* In case of frequency distribution Geometric mean is defined by

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

Ex:-

Find the Geometric mean of 24, 8, 12, 16 & 24.

<u>X</u>	<u>log xⁱ</u>
2	log 2 = 0.301
4	log 4 = 0.6020
8	log 8 = 0.9030
12	log 12 = 1.079
16	log 16 = 1.204
24	log 24 = 1.380
<hr/>	
	$\sum \log x_i = 5.469$

$$G.M = \text{Antilog} \left[\frac{1}{n} \sum \log x \right]$$

$$= \text{Antilog} \left[\frac{1}{5} \times 5.469 \right]$$

$$= \text{Antilog} (0.9115)$$

- ② Find the Geometric mean for the following frequency distribution

<u>Marks</u>	<u>No. of Students(f)</u>	<u>Mid value (x)</u>	<u>log x</u>	<u>f.log x</u>
0-10	5	5	log 5 = 0.699	3.495
10-20	7	15	log 15 = 1.176	8.232
20-30	15	25	log 25 = 1.398	20.97
30-40	25	35	log 35 = 1.544	38.6
40-50	8	45	log 45 = 1.653	13.294

$$N = \Sigma f = 60$$

$$\Sigma f \log x_i = 84.521$$

$$G.M = \text{Antilog} \left[\frac{1}{N} \sum f \log x_i \right]$$

$$\Rightarrow \text{Antilog} \left[\frac{1}{60} \times 84.521 \right]$$

$$= \text{Antilog} (1.4087)$$

$$= 25.628$$

Merits:

- (i) Geometric mean is rigidly defined.
- (ii) It is based on all the observations.
- (iii) It is suitable for further mathematical treatment.
If G_1 & G_2 are the geometric means of two groups of sizes n_1 & n_2 respectively, then the geometric mean G of the combined group of size $n_1 + n_2$ is given by
- (iv) As compared with mean, G.M is affected to a lesser extent by extreme observations.
- (v) It is not affected much by fluctuations of sampling.

Demerits:

- (i) Because of its abstract mathematical character, geometric mean is not easy to understand to calculate for a non-mathematical person.
- (ii) If any one of the observation is zero, geometric mean becomes zero and if any one of the observation is negative, geometric mean becomes imaginary regardless of the magnitude of the other items.

Harmonic mean:-

Harmonic mean is the reciprocal of arithmetic mean and reciprocal of given value.

→ We can write, if x_1, x_2, \dots, x_n is a set of n observation the Harmonic mean is given by

$$\boxed{H.M = \frac{1}{\frac{1}{n} \left[\sum_{i=1}^n \frac{1}{x_i} \right]} = \frac{n}{\sum_{i=1}^n \left(\frac{1}{x_i} \right)}}$$

→ In case of frequency distribution, the Harmonic mean is given by

$$\boxed{H.M = \frac{1}{\frac{1}{N} \left[\sum_{i=1}^n \frac{f_i}{x_i} \right]} = \frac{N}{\sum_{i=1}^n \left(\frac{f_i}{x_i} \right)}} \quad (\text{where } N = \sum f_i)$$

Ex:-

The following table if the weights of 31 persons in a sample inquiry. calculate the mean weight using Harmonic mean.

<u>weights(x)</u>	<u>No. of persons(f)</u>	<u>(f/x)</u>
130	3	0.0230
135	4	0.0296
140	6	0.0428
145	6	0.0414
146	3	0.0205
148	5	0.0338
149	2	0.0134
150	1	0.0066
157	1	0.0063
	$\sum f = N = 31$	$\sum (f/x) = 0.2177$

$$H.M = \frac{N}{\frac{1}{E(F/x)}} = \frac{32}{0.2577} = 124.36$$

② A cyclist pedals from his house to his college at a speed of 10 km. p.hrc. and back from the college to his house at 15 km. p.hrc. Find the average speed.

→ Let the distance from the house to college be x kms.
In going from house to college, the distance is covered in $\frac{x}{10}$ hours.

While coming from college to house, the distance is covered in $\frac{x}{15}$ hours.

Thus a total distance of $2x$ kms is covered in $\left(\frac{x}{10} + \frac{x}{15}\right)$ hours.

Hence, average speed = $\frac{\text{Total distance travelled}}{\text{Total time taken}}$

$$\text{avg} = \frac{\text{dist}}{\text{time}} = \frac{2x}{\left(\frac{x}{10} + \frac{x}{15}\right)} = \frac{2}{\left(\frac{1}{10} + \frac{1}{15}\right)} = 12 \text{ km p.hrc.}$$

Ex:-

Given below is the frequency distribution of marks obtained by 90 students. Compute the arithmetic mean, median and mode.

$$\frac{Z.C + Z.P.E}{2E} =$$

<u>Marks</u>	<u>Marks</u>	<u>n</u>	<u>f</u>	<u>fx</u>	<u>C.F</u>
15-19	14.5-19.5	17	6	102	6
20-24	19.5-24.5	22	14	308	20
25-29	24.5-29.5	27	12	324	32
30-34	29.5-34.5	32	10	320	42
35-39	34.5-39.5	37	10	370	52
40-44	39.5-44.5	42	9	378	61
45-49	44.5-49.5	47	9	423	70
50-54	49.5-54.5	52	10	520	80
55-59	54.5-59.5	57	5	285	85
60-64	59.5-64.5	62	4	248	89
65-69	64.5-69.5	67	1	67	90

Mean

$$\bar{x} = \frac{\sum f x}{\sum f} = \frac{3345}{90} = 37.166$$

$$\text{Median} = l + \frac{h}{f} (N/2 - C)$$

$$\text{where, } \frac{N}{2} = \frac{90}{2} = 45, f = 10, C = 42$$

$$h = 5, l = 34.5$$

C.F is just greater than 52 and the corresponding median class is 34.5-39.5.

$$\text{median} = 34.5 + \frac{5}{10} (45 - 42)$$

$$= 34.5 + \frac{1}{2} (3)$$

$$= 34.5 + 1.5$$

$$= 36$$

$$\begin{aligned}
 \text{Mode} &= 3\text{median} - 2\text{mean} \\
 &= 3(36) - 2(37.16) \\
 &= 108 - 74.32 \\
 &= 33.68
 \end{aligned}$$

Relationship between mean, median & mode:

→ In case of symmetrical distribution

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

$$M = M_d = M_o$$

$$\text{Mean} - \text{Median} = \frac{1}{3} (\text{Mean} - \text{Mode})$$

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

Relationship between AM, GM & HM

Prove that $AM \geq GM \geq HM$

Proof:

Let a & b be two real positive numbers.

i.e. $a > 0, b > 0$.

$$AM = \frac{a+b}{2}$$

$$GM = \sqrt{ab}$$

$$HM = \frac{1}{\frac{1}{a} + \frac{1}{b}} = \frac{2ab}{a+b}$$

$$AM - GM = \frac{a+b}{2} - \sqrt{ab}$$

$$= \frac{a^2 + 2ab + b^2 - 2ab}{2} = \frac{a^2 - 2ab + b^2}{2} = \frac{(a-b)^2}{2}$$

$$> 0 \quad \text{as } (a-b)^2 \geq 0$$

i.e. the square of real quantity is always non-negative.

$$AM - GM \geq 0$$

$$AM \geq GM \quad \text{--- } \textcircled{1}$$

$$GM - HM = \sqrt{ab} - \frac{2ab}{a+b}$$

$$= \sqrt{ab} \left(1 - \frac{2\sqrt{ab}}{a+b} \right)$$

$$= \sqrt{ab} \left(\frac{a+b-2\sqrt{ab}}{a+b} \right)$$

$$= \sqrt{ab} \left(\frac{(a-\sqrt{b})^2}{a+b} \right) \geq 0$$

i.e. the square of a real quantity is always non-negative.

$$GM - HM \geq 0$$

$$GM \geq HM \quad \text{--- } \textcircled{11}$$

Combining eq' $\textcircled{1}$ & $\textcircled{11}$

$$AM \geq GM \geq HM, (\text{Proved})$$

Partition value:—

The value which divide the series into number of equal parts are known as Partition value.

→ There are 3 types of partition value.

1. Quartiles

2. Deciles

3. Percentiles

Quartiles:—

→ The three points which divide the series into four equal parts are known as quartiles.

→ There are 3 points Q_1, Q_2, Q_3 such that $Q_1 \leq Q_2 \leq Q_3$, known as the quartiles.

- Q_1 is the lower or 1st quartile is a value which has 25% of the items of the distribution below it & consequently 75% of the item are greater than it.
- Q_2 is the 2nd quartile & coincides with the median. It has an equal no. of observation above it & below it.
- Q_3 is the 3rd quartile. It has 75% of the observation below it and 25% of the observation above it.

Calculation of quartiles:

1. Individual series:

In case of ungrouped (Individual) data having n observations arranged in the ascending or descending order i^{th} quartile is given by

$$Q_i = \text{value of } \frac{i(n+1)}{4}^{\text{th}} \text{ item (where } i=1,2,3)$$

2. Discrete series:

(i) find $\frac{iN}{4}$, where $N = \sum f_i$

(ii) C is the cumulative frequency just greater than $\frac{iN}{4}$.

(iii) The corresponding value of x keeps the value of Q_i .

3. Continuous series:

In case of continuous or grouped frequency distribution the i^{th} quartile is given by

$$Q_i = l + \frac{h}{f} \left(\frac{iN}{4} - c \right) \quad (\text{where } i=1,2,3)$$

Ex:-

Marks secured by 11 students in statistics are as follows 34, 27, 64, 89, 75, 82, 54, 49, 68, 84, 45 then calculate Q_1 , Q_2 & Q_3 .

Ans:- Arranging in descending order =

27, 34, 45, 49, 54, 64, 68, 75, 82, 84, 89

$$Q_1 = \text{value of } \frac{i(11+1)}{4} \text{ th item}$$

$$Q_1 = \text{value of } \frac{1(11+1)}{4} \text{ th item}$$

$$= \text{value of } \frac{1 \times 12^3}{4} \text{ th item}$$

= value of 3rd. item = 45

$$Q_2 = \text{value of } \frac{2(11+1)}{4} \text{ th item}$$

$$= \text{value of } \frac{2 \times 12^3}{4} \text{ th item}$$

$$= \text{value of } 6^{\text{th}} \text{ item} = 64$$

$$Q_3 = \text{value of } \frac{3(11+1)}{4} \text{ th item}$$

$$= \text{value of } \frac{3 \times 12^3}{4} \text{ th item}$$

$$= \text{value of } 9^{\text{th}} \text{ item} = 82$$

Deciles:-

→ The 9 values which divides the series in to 10 equal parts. So, there are 9 deciles such as

D_1, D_2, \dots, D_9 such that $D_1 \leq D_2 \leq \dots \leq D_9$.

→ D_5 coincide with median

Calculation of Deciles:-

1. Individual Series:-

In case of ungrouped data the deciles is

given by

$$D_i = i^{\text{th}} \text{ decile} = \text{value of } \frac{i(n+1)}{10} \text{-th item}$$

(where $i = 1, 2, 3, 4, \dots, 9$)

Discrete series:—

(i) Find $\frac{in}{10}$, where $N = \Sigma f$

(ii) C is the cumulative frequency just greater than $\frac{in}{10}$.

(iii) The corresponding value of x keeps the value of D_i .

Continuous Series:—

In case of continuous series i^{th} deciles is given by

$$D_i = l + \frac{h}{f} \left(\frac{in}{10} - C \right) \quad (i = 1, 2, 3, \dots, 9)$$

Percentiles:—

→ The 99 values which divides the series in to 100 equal parts. So, there are 99 percentiles are P_1, P_2, \dots, P_{99} . such that $P_1 \leq P_2 \leq \dots \leq P_{99}$

→ P_{50} coincide with median.

Calculation of Percentiles:—

1. Individual Series:—

In case of ungrouped data the percentiles are given by

$$P_i = i^{\text{th}} \text{ percentile} = \text{value of } \frac{i(n+1)}{100} \text{-th item}$$

(where $i = 1, 2, 3, \dots, 99$)

Discrete series:

- (i) Find $\frac{in}{100}$, where $N = EF$
- (ii) c is the cumulative frequency just greater than $\frac{in}{100}$.
- (iii) The corresponding value of x keeps the value of P_i .

Continuous series:

In case of continuous series the i th percentile is given by

$$P_i = l + \frac{h}{F} \left(\frac{in}{100} - c \right)$$

Ex:-

Annual sales No. of stores

0-100

4

100-200

28

200-300

35

300-400

25

400-500

7

500-600

1

Compute the Quartile, 6th deciles and 20th Percentiles.

Ans:-

Annual sales

0-100

No. of stores

C.F

4

100-200

28

4

200-300

35

9

300-400

25

92

400-500

7

99

500-600

1

100

$$N = 100 \quad (\text{E.B.I.E. = 1 month})$$

$$(i) Q_1 = \frac{N}{4} = \frac{100}{4} = 25$$

C.F just greater than $\frac{N}{4}$ is 32.
The 1st quartile lies in the class interval 100-200.

$$Q_1 = l + \frac{h}{f} \left(\frac{in}{4} - c \right)$$

$$= 100 + \frac{100}{28} (25 - 4)$$

$$= 100 + \frac{100}{28} \times 21$$

$$= \frac{2800 + 2100}{28} = \frac{4900}{28} = 175$$

$$(ii) Q_2 = \frac{100}{2} = 50$$

C.F just greater than $\frac{N}{2}$ is 67.

The 2nd quartile lies in the class interval 200-300.

$$Q_2 = l + \frac{h}{f} \left(\frac{in}{2} - c \right)$$

$$= 200 + \frac{100}{35} (50 - 32)$$

$$= 200 + \frac{100}{35} (18)$$

$$= \frac{7000 + 1800}{35} = 251.43$$

$$(iii) Q_3 = \frac{3(N+1)}{4} = \frac{3(100+1)}{4} = \frac{30001}{4} = 7500.25 = 7500.75$$

C.F is just greater than $\frac{3(N+1)}{4}$ is 92.

The 3rd quartile lies in the class interval 300-400.

$$Q_3 = l + \frac{h}{f} \left(\frac{in}{4} - c \right)$$

$$= 300 + \frac{100}{25} (77.5 - 67)$$

$$= 300 + \frac{100}{25} \times 8.85$$

$$= \frac{7500 + 885}{25} = 332$$

(iv) $D_6 = ?$

$$\frac{6N}{10} = \frac{6 \times 100}{10} = 60$$

The cumulative frequency is just greater than 60. The corresponding interval class is 200-300.

$$D_6 = l + \frac{h}{f} \left(\frac{6N}{10} - c \right)$$

$$l = 200, h = 100, f = 35, \frac{6N}{10} = 60, c = 32$$

$$D_6 = 200 + \frac{100}{35} (60 - 32)$$

$$= 200 + \frac{100}{35} (28)$$

$$= \underline{\underline{200+2800}}$$

$$= \underline{\underline{2800}}$$

$$= \underline{\underline{280.42}}$$

(v) $P_{20} = ?$

$$\frac{20N}{100} = \frac{20 \times 100}{100} = 20$$

The cumulative frequency just greater than 20 is 32. The corresponding interval class is 100-200.

$$P_{20} = l + \frac{h}{f} \left(\frac{20N}{100} - c \right)$$

$$l = 100, h = 100, f = 28, \frac{20N}{100} = 20, c = 4$$

$$P_{20} = 100 + \frac{100}{28} (-20 - 4)$$

$$= 100 + \frac{100}{28} \times 16$$

$$= \underline{\underline{2800+1600}} \over 28$$

$$= 157.14$$

MEASURES OF DISPERSION

Measures of Dispersion:

→ The measures of dispersion is a measure of extent to which individual items vary from central value.

→ Literally, dispersion means scatteredness.

→ If the dispersion is small it means that the given data values are closer to the central value.

Objectives or significance of the Measures of Dispersion:

* The main objectives of studying dispersion may be summarised as :—

- (1) To find out the reliability of an average.
- (2) To control the variation of the data from the central value.
- (3) To compare two or more sets of data regarding their variability.
- (4) To obtain other statistical measures for further analysis of data.

Characteristics for an ideal Measure of Dispersion:

- (i) It should be rigidly defined.
- (ii) It should be easy to calculate & easy to understand.
- (iii) It should be based on all the observations.
- (iv) It should be amenable to further mathematical treatment.
- (v) It should be affected as little as possible by fluctuations of sampling.
- (vi) It should not be affected much by extreme observations.

$$\text{below minimum} - \text{above maximum} = \text{spread}$$

Types of Dispersion:

(1) Absolute measures of dispersion

(2) Relative measures of dispersion

(1) Absolute measures of dispersion:

- The measures of dispersion which are expressed in terms of the original units of a series are known as Absolute measures of dispersion.
- Such measures are not suitable for comparing the variability of the two distributions which are expressed in different units of measurement.

(2) Relative measures of dispersion:

- The Relative measures of dispersion are obtained as ratios or percentage and it is independent of units of measurement for comparing of the variability of the two distributions, we can use Relative measures of dispersion instead of Absolute measures of dispersion.

Various Measures of dispersion:

- The Various measures of dispersion are:-

(i) Range

(ii) Quartile deviation or semi-Interquartile range

(iii) Mean deviation

(iv) Standard deviation

(v) Variance

Range:-

- Range is defined as the difference between the maximum value and minimum value of the distribution.

$$\text{Range} = \text{Maximum value} - \text{Minimum value}$$

Relative Measures of Range:

→ If we want to compare the variability of the distribution given in different units of measurements we can use Relative measures. So, this relative measure is called as Co-efficient of range.

Let maximum value = A

Minimum value = B

$$\text{Co-efficient of range} = \frac{\text{Maximum value} - \text{Minimum value}}{\text{Maximum value} + \text{Minimum value}}$$

$$\rightarrow \text{Co-efficient of range} = \frac{A-B}{A+B}$$

Example:

Calculate the Range and the co-efficient of range of A.S monthly earning for a year.

<u>Month</u>	<u>Monthly earning (in hundred)</u>
--------------	-------------------------------------

1 13900

2 15000

3 15100

4 15100

5 15700

6 15800

7 16000

8 16100

9 16200

10 16200

11 17300

12 17500

Range = Maximum earning - Minimum earning

$$= 17500 - 13900 = 3600$$

* Co-efficient of Range

$$= \frac{\text{Maximum earning} - \text{Minimum earning}}{\text{Maximum earning} + \text{Minimum earning}}$$

$$= \frac{3600}{31400} = 0.115$$

Quartile Deviation or semi Inter-quartile Range:-

→ It is given by
$$Q.D = \frac{Q_3 - Q_1}{2}$$

where, Q_3 = upper quartile or third quartile
 Q_1 = lower quartile

→ For comparative studies of variability of two distributions, we need a relative measure which is known as Co-efficient of Quartile Deviation.

$$\text{Co-efficient of } Q.D = \frac{\frac{Q_3 - Q_1}{2}}{\frac{Q_3 + Q_1}{2}} = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

Merits of Quartile Deviation:-

→ Quartile deviation is quite easy to understand and calculate.

→ It has a number of obvious advantages over range as a measure of dispersion.

Ex:-

(a) As against range which was based on two observations only, Q.D. makes use of 50% of the data & such is obviously a better measure than Range.

(b) Since Q.D ignores 25% of the data from the beginning of the distribution and another 25% of the data from the top end, it is not affected at all by extreme observations.

(c) Q.D can be computed from the frequency distribution with open end classes. In fact, Q.D is the only measure of dispersion, which can be obtained while dealing with a distribution open end classes.

Demerits of Quartile deviation:

- Q.D is not based on all the observations since it ignores 25% of the data at the lower end and 25% of the data at the upper end of the distribution. Hence, it can not be regarded as a reliable measure of variability.
- Q.D is affected considerably by fluctuations of sampling.
- Q.D is not suitable for further mathematical treatment.

Mean Deviation: — (Average deviation)

Individual series:

Mean deviation is given by

$$M.D = \frac{1}{n} \sum_{i=1}^n |x_i - A|$$

where, n = no. of observation

A = mean / mode / median

- Absolute measures of mean deviation:
- In case of frequency distribution

$$M.D = \frac{1}{N} \sum_{i=1}^N f_i |x_i - A|$$

Mean deviation about mean:

$$M.D = \frac{1}{N} \sum F |x_i - \bar{x}|$$

Mean deviation about median:

$$M.D = \frac{1}{N} \sum F |x_i - M_d|$$

Mean deviation about mode:

$$M.D = \frac{1}{N} \sum F |x_i - M_o|$$

Relative Measures of Mean deviation:-

- Relative Measures of dispersion of mean deviation is called the Co-efficient of mean deviation.
- It is given by

$$\text{Co-efficient of M.D.} = \frac{\text{Mean Deviation}}{\text{Average about which it is calculated}}$$

$$\text{Co-efficient of M.D. about mean} = \frac{\text{M.D.}}{\text{Mean}}$$

$$\text{Co-efficient of M.D. about median} = \frac{\text{M.D.}}{\text{median}}$$

NOTE:-

Mean deviation about mean is always greater than mean deviation about median.

OR

Mean deviation is least when taken about median.

Merits of Mean deviation:-

- Mean deviation is rigidly defined.
- It is easy to understand and calculate.
- Mean deviation is based on all observations.
- It is thus definitely a better measure of dispersion than the range and quartile deviation.
- Mean deviation provides an accurate and true measure of dispersion.
- As compared with standard deviation, it is less affected by extreme observation.

Demerits of Mean deviation:

- The strongest objection against mean deviation is that while computing its value, we take the absolute value of the deviations about an average and ignore the signs of the deviations.
- It is rarely used in sociological studies.
- It can not be computed for distributions with open end classes.
- Mean deviation tends to increase with the size of the sample though not proportionately and not so rapidly as range.

Ex:-

MEASURE OF DISPERSION

calculate the M.D from the median for the following data.

Marks	Frequency	C.F	Mid value (x)	$ x - Md $	$f x - Md $
0-10	5	5	5	41.43	207.15
10-20	8	13	15	31.43	251.44
20-30	12	25	25	21.43	150
30-40	12	37	35	11.43	137.16
40-50	28	65	45	1.43	40.04
50-60	20	85	55	8.57	171.14
60-70	10	95	65	18.57	185.70
70-80	10	100	75	28.57	285.70
					$N = \sum f = 100$

$$\frac{N}{2} = \frac{100}{2} = 50$$

$$\left[f(x - x) \sum_{i=1}^n \frac{f_i}{N} \right] = 0$$

$$\sum f |x - Md| = 1428.6$$

The cumulative frequency is greater than equal to 50 is 260. The corresponding class (40-50) is the median class.

$$l=40, h=10, f=28, c=32, N/2=50$$

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

$$= 40 + \frac{10}{28} (50 - 32)$$

$$= 40 + \frac{10}{28} \times 18$$

$$= 40 + 6.43$$

$$= 46.43$$

$$\text{M.D above median} = \frac{1}{N} \sum f |x - \text{Md}|$$

$$= \frac{1}{100} \times 1428.6$$

$$= \frac{1428.6}{100} = 14.286 = 14.29$$

Standard deviation:

→ Standard deviation is defined as the positive square root of the arithmetic mean of the square of the deviation of the given deviations from their arithmetic mean.

→ It is usually denoted by ' σ '.

Individual series:

In case of individual series the standard deviation is given by

$$\boxed{\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$\bar{x} = \frac{1}{n} \sum x_i, n = \text{no. of observation}$$

(Arithmetic mean of the given values)

frequency distribution:— appears with tabulations

In case of frequency distribution the standard deviation is given by

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^n f_i (x_i - \bar{x})^2}$$

where, $N = \sum f_i$, $\bar{x} = \frac{\sum f_i x_i}{\sum f_i}$

$$\text{Co-efficient of Standard Deviation} = \frac{\text{S.D}}{\text{mean}} = \frac{\sigma}{\bar{x}}$$

Variance:

→ The square of the standard deviation is known as variance.

→ It is denoted by σ^2 .

Individual Series:

In Case of Individual series Variance is given

by

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

frequency distribution:

In Case of frequency distribution variance is given by

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

Co-efficient of variation:

→ hundred times the co-efficient of dispersion based on standard deviation is called the co-efficient of variance.

$$C.V = \frac{\sigma}{\bar{x}} \times 100$$

→ For comparing the variability of two series we

calculate the co-efficient of variation for each series. The series having greater C.V is said to be more variability than the others series. And the series having lesser C.V is said to be more consistent or homogeneous than the others.

Combined Standard deviation:

Let $\sigma_1, \sigma_2, \dots, \sigma_k$ are the S.D $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_k$ are the A.M. and n_1, n_2, \dots, n_k are the sizes of k group then the S.D σ of the combined group up size $N = n_1 + n_2 + n_3 + \dots + n_k$ is given by

$$\sigma = \sqrt{\frac{n_1(\sigma_1^2 + d_1^2) + n_2(\sigma_2^2 + d_2^2) + \dots + n_k(\sigma_k^2 + d_k^2)}{n_1 + n_2 + \dots + n_k}}$$

where $d_1 = \bar{x}_1 - \bar{x}$ and $\bar{x} = \frac{n_1\bar{x}_1 + n_2\bar{x}_2}{n_1 + n_2}$
 $d_2 = \bar{x}_2 - \bar{x}$

Moments:

→ Moment gives the idea about the Central part of the distribution and their variability. It also gives the idea about the skewness and kurtosis of the distribution.
 → There are 4 types of moments.

(i) Raw moment

(ii) Central moment

(iii) Factorial moment

(iv) Absolute moment

Raw moment:

The n^{th} raw moment of a variable x about any point $x = A$ is usually denoted by m_n is given by

$$M_n' = \frac{1}{N} \sum f_i (x_i - A)^n \quad \text{--- ①}$$

→ Raw moment is denoted by (M_n')

putting $n=1$

$$M_1' = \frac{1}{N} \sum f_i (x_i - A) \quad \text{--- ②}$$

$$= \frac{1}{N} \sum f_i x_i - \frac{1}{N} \sum f_i A \quad \text{--- ③}$$

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

$$\text{putting } A=0, \quad M_1' = \bar{x} \quad \text{--- ④}$$

i.e. the 1st order of raw moment is equivalent to mean.

Taking $n=2, 3, 4$ in eq ①

we get,

$$M_2' \text{ (2nd order raw moment about } A) = \frac{1}{N} \sum f_i (x_i - A)^2$$

$$M_3' \text{ (3rd order raw moment about } A) = \frac{1}{N} \sum f_i (x_i - A)^3$$

$$M_4' \text{ (4th order raw moment about } A) = \frac{1}{N} \sum f_i (x_i - A)^4$$

Central moment:

The n th moment of a variable x about the mean \bar{x} usually denoted as M_n .

It is given by

$$M_n = \frac{1}{N} \sum f_i (x_i - \bar{x})^n \quad \text{--- ⑤}$$

$$\text{putting } n=0, \quad M_0 = \frac{1}{N} \sum f_i (x_i - \bar{x})^0 = 1$$

$$(x) \text{ mean} = n=1, \quad M_1 = \frac{1}{N} \sum f_i (x_i - \bar{x})^1 = 0$$

$$n=2, \quad M_2 = \frac{1}{N} \sum f_i (x_i - \bar{x})^2 = \sigma^2 \text{ (variance)}$$

The 2nd order central moment is equivalent to variance.

$$M_2 - M_1^2 = \sigma^2$$

Relation between raw moment & central moment:

$$M'_2 = M_2 + (M'_1)^2$$

$$M'_3 = M_3 + 3M_2M'_1 + (M'_1)^3$$

$$M'_4 = M_4 + 4M_3M'_1 + 6M_2(M'_1)^2 + (M'_1)^4$$

Relation between central moment & raw moment:

$$M_2 = M'_2 - (M'_1)^2$$

$$M_3 = M'_3 - 3M'_2M'_1 + (M'_1)^3$$

$$M_4 = M'_4 + 4M'_3M'_1 + 6M'_2(M'_1)^2 - 3(M'_1)^4$$

Factorial moment:

Factorial moment of order r about the origin of the distribution α_i/f_i ($i=1, 2, \dots, n$) is obtained as

$$M_{(r)}' = \frac{1}{N} \sum_{i=1}^n f_i x_i^{(r)}$$

$$\text{Where, } x^{(r)} = \alpha(\alpha-1)(\alpha-2) \dots (\alpha-r+1)$$

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

Thus, the factorial moment of order r about any point $\alpha=a$ is given by,

$$M_{(r)} = \frac{1}{N} \sum f_i (x_i - a)^{(r)}$$

$$\text{Where } (x-a)^{(r)} = (x-a)(x-a-1)(x-a-2) \dots (x-a-r+1)$$

Put. $r=1, 2, 3, 4$ in eq(i), we get

$$M'_1 = \frac{1}{N} \sum f_i x_i = M_1 \quad (\text{about origin}) = \text{Mean}(X)$$

$$(2) M'_2 = \frac{1}{N} \sum f_i x_i^2$$

$$= \frac{1}{N} \sum f_i x_i (x_i - 1) = \frac{1}{N} \sum f_i x_i^2 - \frac{1}{N} \sum f_i x_i$$

$$M_{(2)} = M_2 - M_1^2$$

$$M_3' = \frac{1}{N} \sum f_i x_i^{(3)}$$

$$= \frac{1}{N} \sum f_i x_i (x_i - 1)(x_i - 2)$$

$$\text{Endo} = \frac{1}{N} \sum f_i [(x_i^2 - x_i)(x_i - 2)]$$

$$= \frac{1}{N} \sum f_i [x_i^3 - 2x_i^2 - x_i^2 + 2x_i]$$

$$= \frac{1}{N} \sum f_i x_i^3 - 3 \frac{1}{N} \sum f_i x_i^2 + 2 \frac{1}{N} \sum f_i x_i$$

$$M_3' = M_3 - 3M_2 + 2M_1$$

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

$$= \frac{1}{N} \sum f_i [(x_i(x_i - 1)(x_i - 2)(x_i - 3))]$$

$$= \frac{1}{N} \sum f_i x_i (x_i^3 - 6x_i^2 + 11x_i - 6)$$

$$= \frac{1}{N} \sum f_i x_i^4 - 6 \frac{1}{N} \sum f_i x_i^3 + 11 \frac{1}{N} \sum f_i x_i^2 - 6 \frac{1}{N} \sum f_i x_i$$

$$M_4' = M_4 - 6M_3 + 11M_2 - 6M_1$$

Conversely, we will get,

$$M_1' = M_{(1)}'$$

$$M_2' = M_{(2)}' + M_{(4)}'$$

$$M_3' = M_{(3)}' + 3M_{(2)}' - 2M_{(4)}'$$

$$M_4' = M_{(4)}' + 3(M_{(2)}' + M_{(3)}) - 2M_{(4)}'$$

$$M_3' = M_{(3)}' + 3M_{(3)}' + M_{(1)}'$$

$$M_4' = M_{(4)}' + 6M_{(3)}' - 11M_{(2)}' + 6M_{(1)}'$$

$$= M_4 + 6(M_{(3)}' + 3M_{(2)}' + M_{(1)}') - 11(M_{(2)}' + M_{(1)}) - 2M_{(4)}$$

$$M_4' = M_{(4)} + 6M_{(3)} + 7M_{(2)} + M_{(1)}$$

Absolute moments:

For the frequency distribution x_i/f_i ($i=1, 2, 3, \dots, n$) the r th absolute moment of the variable x about the origin is given by,

$$\frac{1}{N} \sum_{i=1}^n f_i |x_i|^r, \quad N \neq f_i$$

where, $|x_i|^r$ represents the absolute / modulus value of x_i^r .

The r th absolute moment of the variable about the mean \bar{x} is given by,

$$\frac{1}{N} \sum_{i=1}^n f_i |x_i - \bar{x}|^r$$

NOTE:

The moments are used to describe the various characteristics of a frequency distribution like Central tendency, variation, skewness and kurtosis.

Sheppard's correction: → This was introduced

→ While calculating moments it is assumed that all the values of a variable in a class interval are concentrated at the center of that interval which is mid value.

→ So, this assumption is an approximation to facilitate calculation and it introduced some error, which is known as grouping error.

$$W.L.C = (M_{(4)} + 2M_{(3)} + 3M_{(2)} + M_{(1)}) - \left(\frac{(M_{(4)} + 2M_{(3)} + 3M_{(2)} + M_{(1)})}{N} \right)^2$$

Conditions for applying Sheppard's correction:

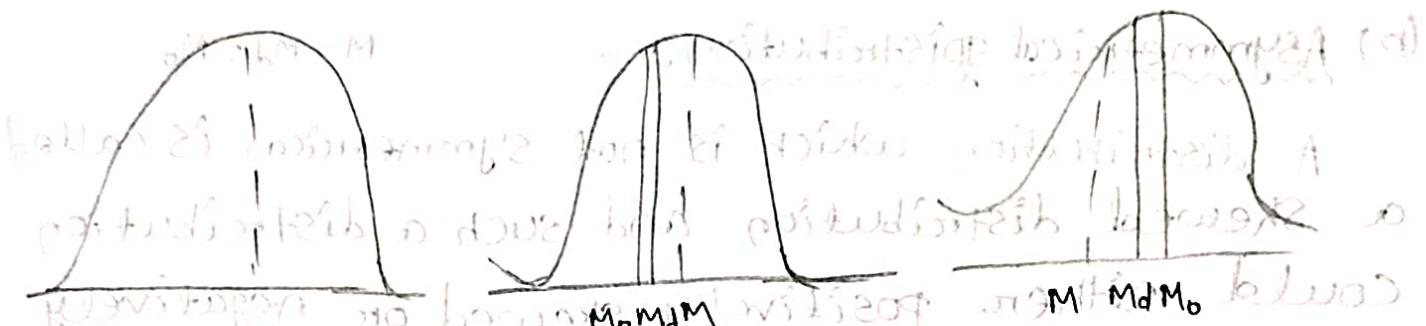
- (i) Total frequency should be large.
 - (ii) The frequency distribution should be continuous.
 - (iii) The distribution must have finite range.
- The effect due to grouping at the mid point of the interval can be corrected by the following formulae. So, this formula is known as Sheppard's correction.

$$M_2 (\text{Corrected}) = M_2 - \frac{h^2}{12}$$

$$M_3 (\text{Corrected}) = M_3$$

$$M_4 (\text{Corrected}) = M_4 - \frac{h^2}{2} \times M_2 + \frac{7}{240} \times h^4$$

where, h is the width of the class interval.



Symmetrical

Distribution

moderately skewed

Distribution

very skewed

Distribution

(M > M_d > M_o)

(M < M_d < M_o)

Skewness:

→ When a series is not symmetrical it said to be asymmetrical or skewed.

→ Literally, skewness means lack of symmetry; we study skewness to have an idea about the shape of the curve which can draw with the help of the given data.

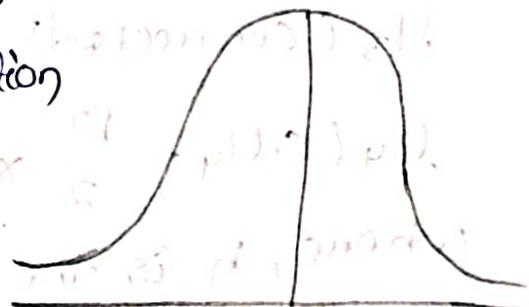
→ A distribution is said to be skewed if

- (i) Mean, median and mode fall at different points i.e. $\text{Mean} \neq \text{Median} \neq \text{Mode}$
- (ii) Quartiles are not equidistant from median i.e. $Q_3 - M_d \neq M_d - Q_1$
- (iii) The curve drawn with the help of the given data is not symmetrical but stretched more to one side than to the other.

(a) Symmetrical Distribution:

In symmetrical distribution

the values of mean, median and mode coincide.



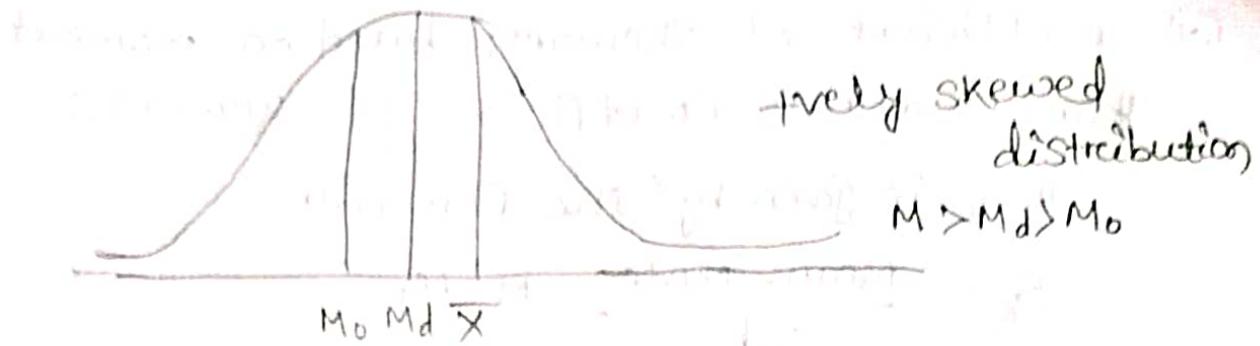
(b) Asymmetrical Distribution:

$$M = M_d = M_o$$

A distribution which is not symmetrical is called a skewed distribution and such a distribution could either positively skewed or negatively skewed distribution.

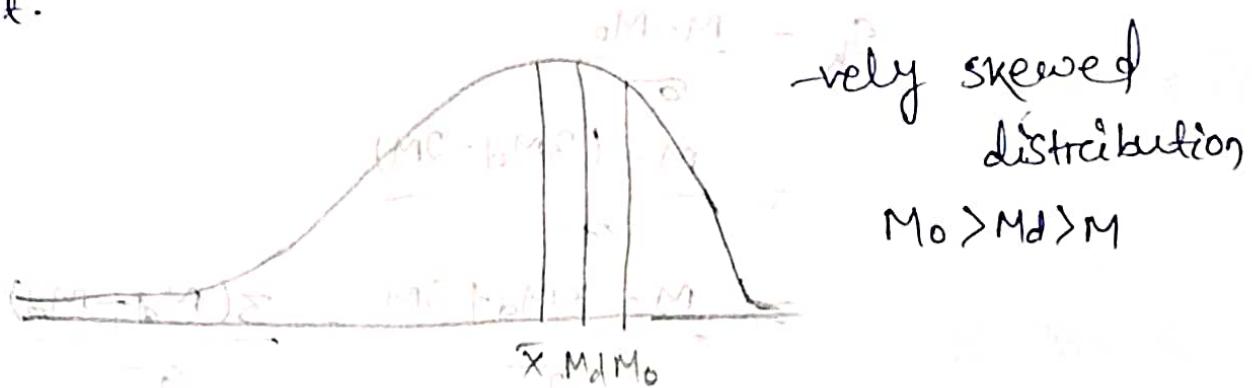
(c) Positively Skewed Distribution:

In the positively skewed distribution, the values of the mean is maximum, mode is least and median lies in between the two. A frequency distribution for which the curve has a longer tail towards the right is said to be truly skewed distribution.



(d) Negatively skewed distribution:—

In a negatively skewed distribution the value of mode is maximum, mean is least and median is lies in between the two. In negatively skewed distribution if the longer tail towards the left.



Measures of Skewness:—

Various measures of Skewness are—

$$(i) S_K = M - Md \quad (ii) S_K = M - M_o \quad (iii) (Q_3 - M_d) - (M_d - Q_1)$$

→ There are the absolute measures of skewness. As in dispersion for comparing two series. So, do not calculate these absolute measures but we calculate the relative measures called the Co-efficient of Skewness.

→ There are three important co-efficient of skewness are—

- (i) Prof Karl Pearson's Co-efficient of skewness
- (ii) Prof Bowley's Co-efficient of skewness

(iii) Co-efficient of Skewness based on moment

Karl Pearson's Co-efficient of Skewness:—

This is given by the formula

$$S_K = \frac{\text{Mean} - \text{Mode}}{\text{S.D.}} = \frac{M - M_0}{\sigma}$$

where, σ is the S.D. of the distribution

→ If mode is ill-defined, then using the empirical relation $M_0 = 3M_d - 2M$, for a moderately asymmetrical distribution we get,

$$\text{MODE} = 3M_d - 2M$$

$$S_K = \frac{M - M_0}{\sigma}$$

$$= \frac{M - (3M_d - 2M)}{\sigma}$$

$$= \frac{M - 3M_d + 2M}{\sigma} = \frac{3(M_d - M)}{\sigma}$$

We observed data, $S_K = 0$, if $M = M_0 = M_d$

Hence for a Symmetrical distribution, mean, median, and mode coincide.

Skewness is positive if $M > M_0$ or $M > M_d$

Skewness is negative if $M < M_0$ or $M < M_d$

The limits of Karl Pearson's Co-efficient of Skewness are $-3 \leq S_K \leq 3$.

Bowley's Co-efficient of Skewness:—

$$S_K = \frac{(Q_3 - M_d) - (M_d - Q_1)}{(Q_3 - M_d) + (M_d - Q_1)}$$

$$= \frac{Q_3 + Q_1 - 2M_d}{Q_3 - Q_1}$$

2) If $S_K = 0$, $Q_3 - M_d = M_d - Q_1$
 This implies that for a symmetrical distribution ($S_K = 0$), median is equidistant from the upper and lower quartiles.

Skewness is positive if

$$Q_3 - M_d > M_d - Q_1 \Rightarrow Q_3 + Q_1 > 2M_d$$

Skewness is negative if

$$Q_3 - M_d < M_d - Q_1 \Rightarrow Q_3 + Q_1 < 2M_d$$

The limits of Bowley's co-efficient of skewness is $-1 \leq S_K \leq 1$.

3) Based upon moments, co-efficient of skewness is

$$\text{skewness } S_K = \frac{\beta_1 (\beta_2 + 3)}{2(5\beta_2 - 6\beta_1 - 9)}$$

where, $\beta_1 = \frac{M_3^2}{M_2^3}$ is the co-efficient of skewness.

$\beta_2 = \frac{M_4}{M_2^2}$ is the co-efficient of skewness.
Pearson's β and γ co-efficients:

Karl Pearson's defined the following four co-efficient, based upon the first moment about mean:

$$\beta_1 = \frac{M_3^2}{M_2^3}, \quad \gamma_1 = +\sqrt{\beta_1}$$

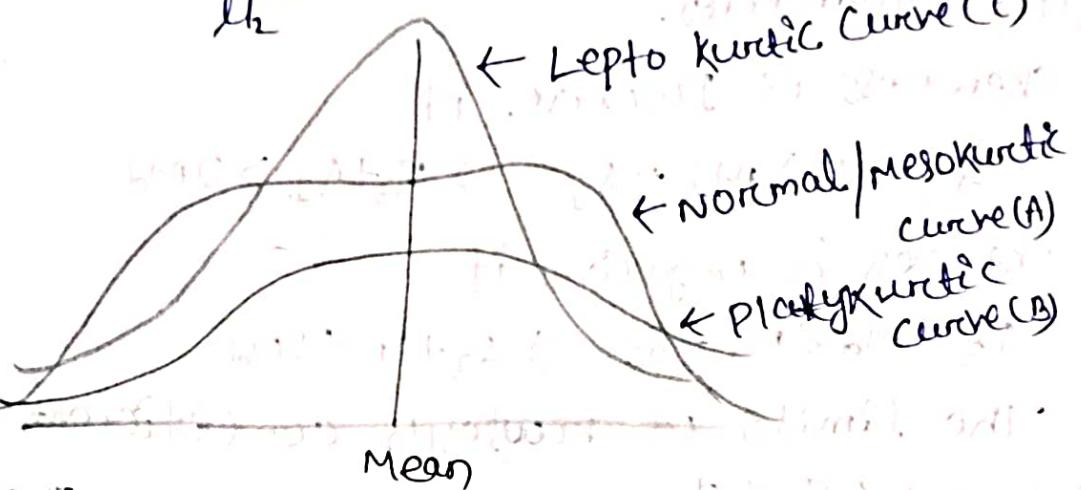
$$\beta_2 = M_4 / M_2^2, \quad \gamma_2 = \beta_2 - 3$$

Kurtosis:

Kurtosis gives us to have an idea about the "flatness or peakedness" of the frequency curve.

It is measured by the co-efficient β_2 or its deviation γ_2 is given by

$$\beta_2 = \frac{M_4}{M_2^2}, \quad \gamma_2 = \beta_2 - 3$$



→ Curve of the type 'A' which is neither flat nor packed is called the normal curve or mesokurtic curve and for such a curve $\beta_2 = 3$ i.e. $\gamma_2 = 0$.

→ Curve of the type 'B' which is flatter than the normal curve is known as platykurtic curve and for such a curve, $\beta_2 < 3$ i.e. $\gamma_2 < 0$.

→ Curve of the type 'C' which is more packed than the normal curve is called leptokurtic curve and for such a curve $\beta_2 > 3$ i.e. $\gamma_2 > 0$.

$$M_4 = M_2 + \frac{\sigma^4}{\mu^2} = 3\sigma^4 + \mu^4$$

$$\therefore \beta_2 = \frac{M_4}{M_2^2} = \frac{3\sigma^4 + \mu^4}{3\sigma^4} = 1 + \left(\frac{\mu}{\sigma}\right)^4$$

$\therefore \text{Leptokurtic}$

with lower β_2 as result of low spread of data
and vice versa with β_2 higher to spread of data.

UNIT-3

Correlation:

In a bivariate distribution we have studied two variables. If we change in one variable affects the change in other variables. so, we say that these two variables are said to be Correlated.

→ Correlation is the linear relationship between two variables.

Types of Correlation:

There are 3 types of Correlation.

1. Positive Correlation

2. Negative Correlation

3. Zero Correlation

Positive Correlation:

If the two variables deviate in same direction i.e. increase or decrease of one variables result in increase or decrease of the other variables. The correlation is said to be Positive.

Ex:- Income and Expenditure

Heights and weights of a group of person

Negative Correlation:

If the two variables deviate in opposite direction i.e. increase or decrease in one variables result in increase or decrease of another variables. The correlation is said to be Negative.

Ex:- Price and demand of a commodity

Sales of woolen garments and temperature

Volume and pressure of a perfect gas

Zero Correlation:

If there is no relationship between two variables is known as zero Correlation.

Ex:- The rate of marriage and agricultural production.

- Methods of studying correlation:—
- There are 3 methods for studying correlation.
1. Karl Pearson's Co-efficient of Correlation
 2. Scatter Diagram
 3. Spearman's Rank Correlation

Karl Pearson's Co-efficient of Correlation:—

- Correlation Co-efficient between two variables X and Y usually denoted by r_{xy} or r .
- It is a numerical measure of linear relationship between two variables.
 - The Karl Pearson's developed a formula called Co-efficient of Correlation.
 - It is defined as

$$r_c = \frac{\frac{1}{n} \sum (x - \bar{x})(y - \bar{y})}{\sqrt{\frac{1}{n} \sum (x - \bar{x})^2} \sqrt{\frac{1}{n} \sum (y - \bar{y})^2}}$$

Or

$$r_c = \frac{\frac{1}{n} \sum xy - \bar{x}\bar{y}}{\sqrt{\frac{1}{n} \sum x^2 - (\bar{x})^2} \sqrt{\frac{1}{n} \sum y^2 - (\bar{y})^2}}$$

- The limits of Correlation Co-efficient is -1 to $+1$. i.e. $-1 \leq r_c \leq 1$.
- It means the Correlation Co-efficient can not exceed unity.

Interpretation of r_c :—

- If $r_c = 0$; No correlation.
- If $r_c > 0$, the correlation is positive.
- If $r_c < 0$, the correlation is Negative.
- If $r_c = +1$, the correlation is perfect positive.
- If $r_c = -1$, the correlation is perfect negative.

Ex:-

Calculating the correlation co-efficient for the following heights of fathers (X) and their sons (Y).

X	Y	XY	X ²	Y ²
65	67	4355	4225	4489
66	68	4488	4356	4624
67	65	4355	4489	4225
67	68	4556	4489	4624
68	72	4896	4684	5184
69	72	4968	4761	5184
70	69	4830	4900	4761
72	71	5112	5184	5041

$$\sum X = 554 \quad \sum Y = 552 \quad \sum XY = 37560 \quad \sum X^2 = 37028 \quad \sum Y^2 = 38132$$

$$\bar{X} = \frac{\sum X}{n} = \frac{554}{8} = 68$$

$$\bar{Y} = \frac{\sum Y}{n} = \frac{552}{8} = 69$$

$$r = \frac{1}{n} \sum XY - \bar{X}\bar{Y}$$

$$\text{mean} \sqrt{\frac{1}{n} \sum X^2 - (\bar{X})^2} \sqrt{\frac{1}{n} \sum Y^2 - (\bar{Y})^2}$$

$$\sqrt{\frac{1}{8} \times 37560 - 4692}$$

$$\sqrt{\frac{1}{8} \times 37028 - 4624} \sqrt{\frac{1}{8} \times 38132 - 4761}$$

$$= \frac{4965 - 4692}{\sqrt{4628.5 - 4624} \sqrt{4765.5 - 4761}}$$

$$= \frac{3}{\sqrt{4.5} \sqrt{5.5}} = \frac{3}{2.12 \times 2.34} = 0.60$$

⇒ The correlation coefficient is 0.60.

Properties of Correlation Co-efficient:

- The correlation co-efficient lies in between -1 to +1,
i.e. $-1 \leq r \leq 1$.
 - The correlation co-efficient is independent of the change of origin and scale.
 - Two variables are ~~are~~ uncorrelated.
- IF $r=0$, then the variables are independent
- $r(x+bx, cy+d) = \frac{ac}{pc} r(x, y)$

Scatter Diagram method:

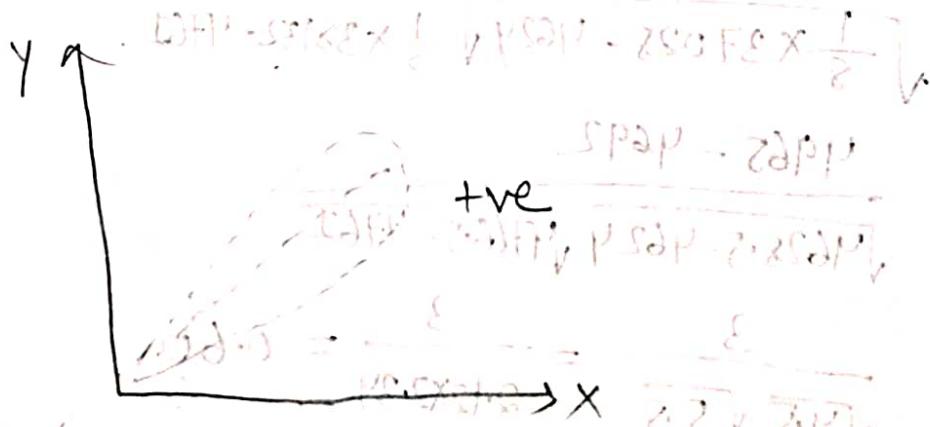
Scatter Diagram is one of the simplest way of diagrammatic representation of bivariate distribution.

Suppose we are given n pairs of values

$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ of two variables 'x' & 'y'. The values of x and y are plotted as dots on the x-axis and y-axis in xy Plane. The diagram of dots are obtain. So, it is known as Scatter Diagram.

Positive Correlation:

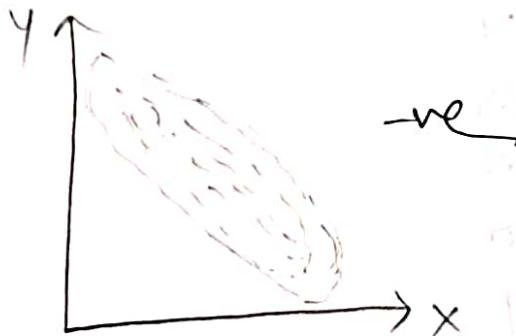
If there is an upward trend raising from lower left hand corner and going upward to the upper right hand corner the correlation is positive.



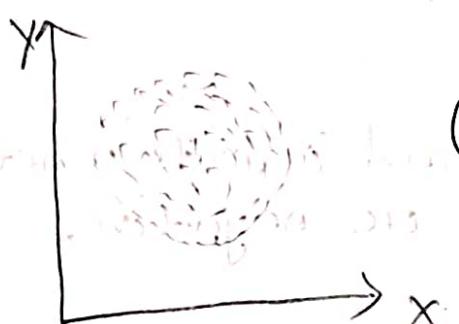
$r > 0$ \Leftrightarrow $r > 0$ \Rightarrow Positive Correlation

Negative Correlation :-

If there is a downward trend from the upper left hand corner to the lower right corner then it is negative correlation.

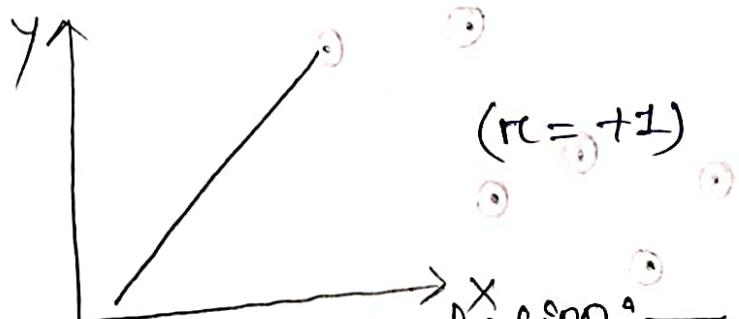


No Correlation :-



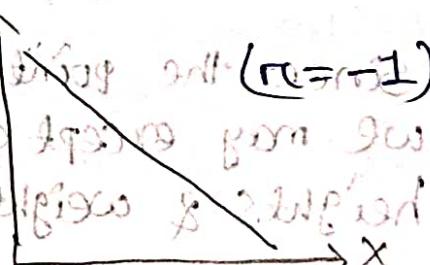
Perfect Positive Correlation :-

If all the points lie on a straight line starting from the left bottom and going upwards the right top the correlation is perfect positive.



Perfect negative correlation :-

If all the points lie on a straight line starting from left up and coming down to right bottom the correlation is perfect negative.

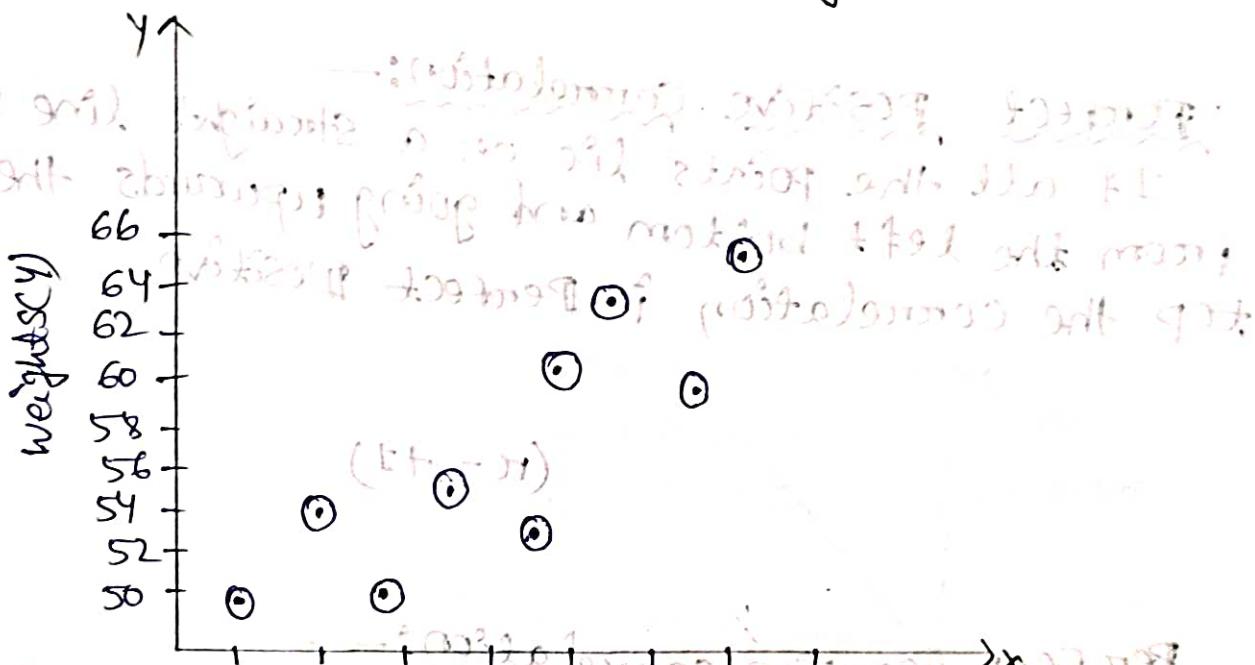


Ex:-

The following are the heights and weights of 10 students of a Bcom class.

<u>heights(x)</u>	<u>weights(y)</u>
62	50
72	65
68	63
58	50
65	54
70	60
66	61
63	55
60	54
72	65

Draw a scatter diagram and indicate whether the correlation is positive or negative.



(Since) the points are dense i.e close to each other so, we may expect a high degree of correlation b/w the heights & weights further the points are relevant upward.

end, starting from the left bottom going up towards the right top. The correlation is positive. So, hence we may expect a high degree of +ve correlation between the series of heights and weights in the class of BCom Student.

Rank Correlation:

- It is used to measure the degree of association between two sets of qualitative observations which can be rank ~~not~~ graded through scores.
- The Spearman's Rank Correlation Co-efficient usually denoted by (r).
- It is given by
$$r = 1 - \frac{6 \sum d_i^2}{n(n^2-1)}$$

(Where, $d_i = \text{Rank of } X - \text{Rank of } Y$
 $n = \text{no. of observation}$)

Computation of Correlation Co-efficient:

- (i) When actual ranks are given
- (ii) When ranks are not given

Case-1: — When Actual Ranks are Given

The following steps are:—

- Compute d , the difference of ranks.
- Compute d^2
- Obtain the sum $\sum d^2$
- Use the formula to get the value of r .

Ex: —

The Ranks of the 15 students in two subjects A and B are given below. The two numbers within the brackets denoting the ranks of the same student in A and B respectively.

$(1,10), (2,7), (3,2), (4,6), (5,4), (6,8), (7,3), (8,1), (9,11), (10,15)$,
 $(11,9), (12,5), (13,14), (14,12), (15,13)$

Ans: Rank of X Rank of Y $\text{d} = X - Y$ $\sum d^2$

1	10	-9	81
2	7	-5	25
3	2	10	100
4	6	-2	4
5	4	1	1
6	8	-2	4
7	3	4	16
8	1	7	49
9	11	-2	4
10	15	5	25
11	9	2	4
12	5	7	49
13	14	-1	1
14	12	-2	4
15	13	2	4

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

$$= 1 - \frac{6 \times 272}{15(15^2-1)}$$

$$= 1 - \frac{6 \times 272}{15 \times 225}$$

$$\sum d^2 = 0$$

$$\sum d^2 = 272$$

$17 \times 17 = \frac{18}{35} = \frac{18}{35}$

case-2 When Ranks are not Given

- When the actual data but not ranks relating to two variables are given. In such a case we shall convert data into ranks.
- The highest (smallest) observation is given the rank 1. Then the highest (next lower) observation is given in rank 2 and so on.
- In case of repeated ranks the rank correlation co-efficient is given by

$$r = 1 - \frac{6}{n(n^2-1)} \left[\sum d^2 + \frac{m(m^2-1)}{12} \right]$$

(where n = no. of times an item is repeated)

Ex: —

A psychologist wanted to compare two methods of teaching, he selected a random sample of 20 students, he grouped them in to 10 pairs. So that the students in a pair have approximately equal scores on an intelligence scale in each pair. One student was taught by method A and method B and examine after the course. The marks often by them are tabulated below. Then find the correlation co-efficient.

Ans:

Pairs	A	B	Rank in A (x)	Rank in B (y)	d = x - y	d^2
1	24	37	6	1	5	25
2	29	35	3	2	1	1
3	19	16	8.5	9.5	-1	1
4	14	26	10	9	1	1
5	30	23	1.5	5	-3.5	12.25
6	19	27	8.5	3	5.5	30.25

7	27	19	8	30	20	15	7	5.5	30.25
9	20	16	9.5	25	6.25				
10	28	11	4	22	11	6	5	25	
11	22	11	6	5	25				

$$\sum d^2 = 225$$

We see that in series A, the item 19 is repeated in two times. The correction factor is given by $\frac{m(m^2-1)}{12}$

Where, $m=2$

$$C.F = \frac{2(2^2-1)}{12} = \frac{1}{3}$$

$$\text{Similarly } 30, m=2, C.F = \frac{1}{3}$$

In Series B, the item 36 is repeated in

two times. The correction factor is given

by $\frac{m(m^2-1)}{12}$

for 16, where $m=2$

$C.F = \frac{2(2^2-1)}{12} = \frac{1}{3}$

So, the rank correlation is given by,

$$r_s = 1 - \frac{6}{n(n^2-1)} \left[\sum d^2 - \frac{m(m^2-1)}{12} \right]$$

$$= 1 - \frac{6}{11(11^2-1)} \left(225 + \frac{1}{2} + \frac{1}{2} + \frac{1}{3} \right)$$

$$r_s = \frac{6 \times 226.5}{11 \times 120} = 0.8$$

	A	B	C	D	E	F	G	H	I
26	2	1	11	120	FE	PG	E		
5	5	1	1	226.5	E	28	PG	E	
5	5	1	220	2.8	E	21	PI	E	
25.51	2.8	2.8	26.5	220	= $\frac{26.5}{220} = 0.0295$	AS	PC	P	
25.05	2.2	2	2.1	2.8		ES	OS	Z	

Probable Error: —

- If r is the correlation co-efficient in a sample of n pairs of observations. The standard error is usually denoted as $S.E(r)$.
- It is given by

$$S.E(r) = \frac{1-r^2}{\sqrt{n}}$$

- The probable error of correlation co-efficient is given by

$$P.E(r) = \frac{1-r^2}{\sqrt{n}} \times 0.6745$$

Multiple Correlation:

- The multiple correlation is the study of combined effect of two or more variables on a single variable.

Ex: Yield of crop (say y_1) depends upon quality of seed (x_2), fertility of soil (x_3) and so on.

- The multiple correlation co-efficient is denoted as $R_{1.23}$.

- It is given by

$$R_{1.23} = \frac{r_{12}^2 + r_{13}^2 - 2r_{12}r_{13}r_{23}}{\sqrt{1 - r_{23}^2}}$$

(r_{12} is the total correlation co-efficient between two variables x_1 & x_2)

Where, $R_{1.23}$ is the multiple correlation co-efficient. We study the multiple impacts of 2nd & 3rd independent variable on 1st dependent variable.

$$R_{2.13} = \frac{r_{12}^2 + r_{23}^2 - 2r_{12}r_{23}r_{13}}{\sqrt{1 - r_{13}^2}}$$

$$R_{3.12} = \frac{r_{13}^2 + r_{23}^2 - 2r_{13}r_{23}r_{12}}{\sqrt{1 - r_{12}^2}}$$

- The limits of multiple correlation coefficient is 0 to 1 i.e. $0 \leq R_{1.23} \leq 1$.
- $R_{1.23}$ is always greater than equal to total Correlation co-efficient, i.e $R_{1.23} \geq r_{12} r_{13} r_{23}$.

Partial Correlation:

The Correlation between two variables x_1 and x_2 after eliminating the linear effect of other variable x_3 is known as Partial Correlation.

→ It is denoted by $r_{12.3}$.

→ It is given by,

$$r_{12.3} = \frac{r_{12} - r_{13} r_{23}}{\sqrt{(1-r_{13}^2)(1-r_{23}^2)}}$$

Where $r_{12.3}$ is the Partial Correlation between x_1 and x_2 . we study the partial impact of 2nd & 3rd variable keeping 1st independent variable constant.

$$r_{13.2} = \frac{r_{13} - r_{12} r_{32}}{\sqrt{(1-r_{12}^2)(1-r_{32}^2)}}$$

$$r_{23.1} = \frac{r_{23} - r_{21} r_{31}}{\sqrt{(1-r_{21}^2)(1-r_{31}^2)}}$$

Ex:-

From the data relating to the yield of dry bark (x_1), height (x_2) and girth (x_3). for 18 cinchona plants. The following Correlation Co-efficient were obtained.

$$r_{12} = 0.77, r_{13} = 0.72, r_{23} = 0.52$$

Find the partial Correlation and multiple Correlation co-efficient.

Ans:—

$$\begin{aligned}\text{Partial Correlation } (\rho_{12.3}) &= \frac{\rho_{12} - \rho_{13}\rho_{23}}{\sqrt{(1-\rho_{13}^2)(1-\rho_{23}^2)}} \\ &= \frac{0.77 - 0.72 \times 0.52}{\sqrt{(1-(0.72)^2)(1-(0.52)^2)}} \\ &= \frac{0.77 - 0.37}{\sqrt{(1-0.51)(1-0.27)}} \\ &= \frac{0.4}{\sqrt{0.49 \times 0.73}} \\ &= \frac{0.4}{0.59} = 0.672(3)\end{aligned}$$

Multiple correlation co-efficient ($R_{1.23}$) =

$$\begin{aligned}&\sqrt{\frac{\rho_{12}^2 + \rho_{13}^2 - 2\rho_{12}\rho_{13}\rho_{23}}{1-\rho_{23}^2}} \\ &= \sqrt{\frac{(0.77)^2 + (0.72)^2 - 2 \times 0.77 \times 0.72 \times 0.52}{1-(0.52)^2}} \\ &= \sqrt{\frac{0.59 + 0.51 - 0.57}{1-0.27}}\end{aligned}$$

$$= \sqrt{\frac{1.1 - 0.57}{0.73}}$$

$$\begin{aligned}&= \sqrt{\frac{0.53}{0.73}} \\ &= \sqrt{0.726}\end{aligned}$$

Simple regression:—
The regression analysis confined to the study of only two variable at a time is known as simple regression.

Regression Analysis:

- Regression analysis is a mathematical measure of the average relationship between two or more variables in terms of original units of data.
- Regression analysis is the technique for estimating the relationship among variables.
- In Regression analysis there are two types of variable.
 - (i) Dependent Variable
 - (ii) Independent Variable

Dependent Variable:

The Variable whose value is to be predicted or implemented is called as dependent variable.

Independent Variable:

The variable whose value is used for prediction is called as independent variable.

Multiple regression:

The regression analysis study more than two variables at a time is known as multiple regression.

Ex:- The cost of a product depends on the production and advertising expenditure.

Linear regression:

If the regression curve is a straight line we say that there is linear regression between the variables under study.

Lines of regression:

In regression there are always two lines of regression

→ Lines of regression of Y on X and X on Y

s.t it is given by

$$Y - \bar{Y} = n \cdot \frac{\sigma_y}{\sigma_x} (X - \bar{X}) \quad | \quad Y = aX + b$$

→ Lines of regression of X on y

it is given by

$$X - \bar{X} = n \cdot \frac{\sigma_x}{\sigma_y} (Y - \bar{Y}) \quad | \quad X = aY + b$$

Note:-

The two lines of regression passes through the point \bar{X}, \bar{Y} .

Why there are two lines of regression?

→ There are always two lines of regression one is regression line of y on x another is regression line of x on y .

→ The line of regression of y on x is used to predict the value of y for any given values of x , where x is independent variable and y is dependent variable.

→ Similarly the line of regression of x on y is used to predict the value of x for any given values of y , where x is dependent variable y is independent variable.

Ex:-

Obtain the equations of two lines of regression for the following data also obtain the estimate of n for $y = 70$,

ΣPAP

$\Sigma PAP = 0.0072 X$

$PX - PXP =$

ΣPPL

$\Sigma PPL = 0.0058 X PXP =$

$PAP =$

<u>X</u>	<u>y</u>	<u>xy</u>	<u>x^2</u>	<u>y^2</u>
65	67	4355	4225	4489
66	68	4488	4356	4624
67	65	4355	4489	4225
67	68	4556	4489	4624
68	72	4896	4624	5184
69	72	4968	4761	5184
70	69	4836	4900	4761
72	71	5112	5184	5041

$$\Sigma x = 544 \quad \Sigma y = 552 \quad \Sigma xy = 37,560 \quad \Sigma x^2 = 37,028 \quad \Sigma y^2 = 38,132$$

$$\bar{y} = \frac{\Sigma y}{n}, \text{ to find out } \bar{x} = \frac{\Sigma x}{n}$$

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

$$= \frac{5844}{8} = 68$$

$$\sigma_y = \sqrt{\frac{1}{n} \sum y^2 - (\bar{y})^2}$$

$$= \sqrt{\frac{1}{8} \times 38132 - (69)^2}$$

$$\text{standard deviation of } y = \sqrt{4766.5 - 4761} = \sqrt{5.87} = 2.345$$

$$\sigma_x = \sqrt{\frac{1}{n} \sum x^2 - (\bar{x})^2}$$

$$= \sqrt{\frac{1}{8} \times 37028 - (68)^2}$$

$$= \sqrt{4628.5} = 68$$

$$= \sqrt{4.5} = 2.121$$

$$r = \frac{\frac{1}{n} \sum xy - \bar{x} \bar{y}}{\sigma_y \sigma_x} = \frac{\frac{1}{8} \times 37560 - 4692}{2.345 \times 2.121} = \frac{4695 - 4692}{4.973}$$

$$= \frac{3}{4.973}$$

$$= 0.605$$

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

$$\Rightarrow Y - 69 = 0.6x \cdot \frac{2.35}{2.12} (x - 68)$$

$$\Rightarrow Y = 0.66(x - 68) + 69 \\ = 0.66x - 44.88 + 69$$

$$= 0.66x + 24.12$$

Eqn of line of regression x on y

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

$$\Rightarrow x - 68 = 0.6x \cdot \frac{2.12}{2.35} (Y - 69)$$

$$\Rightarrow x = 0.54(Y - 69) + 68 \\ = 0.54Y - 37.26 + 68 \\ = 0.54Y + 30.74$$

To estimate x for given y , we use the line of regression of x on y .

If $y = 70$, estimated value of x is given by

$$x = 0.54 \times 70 + 30.74 \\ = 68.54$$

Regression Co-efficient:

The regression co-efficient of y on x represent the increment in the value of dependent variable y corresponding to a unit change in the value of independent variable x .

$$\left\{ \begin{array}{l} b_{yx} = \text{Regression Co-efficient } y \text{ on } x \\ = r \cdot \frac{\sigma_y}{\sigma_x} \end{array} \right.$$

The regression co-efficient of X on Y represent the increment in the value of dependent variable X corresponding to a unit change in the value of independent variable Y .

$$\boxed{b_{xy} = \text{Regression Co-efficient } X \text{ on } Y}$$

$$= r \cdot \frac{\sigma_x}{\sigma_y}$$

Properties of Regression Variable:

- ① Correlation Coefficient (r) is the geometric mean between the regression co-efficients (b_{yx}, b_{xy}).
- i.e. $r = \pm \sqrt{b_{yx} * b_{xy}}$
- ② If one of the regression co-efficient is greater than unity other must be less than unity.
- $b_{yx} > 1, b_{xy} < 1$
- ③ The regression co-efficient are independent of Change of origin but not scale.
- ④ The Modulus value of A.M of the regression co-efficient is not less than the modulus value of correlation co-efficient.

$$\text{i.e. } \left| \frac{b_{yx} + b_{xy}}{2} \right| \geq |r|$$

- NOTE:- → If $r=0$, the two variables are uncorrelated.
 → The lines of regression perpendicular to each other.
 → If $r=\pm 1$, the two variables are coincide. The lines of regression parallel to each other.

Ex:-

In a partially destroyed laboratory, record of an analysis of correlation data, the following results are obtained. The variance of $X=9$. Regression equations $8X - 10Y + 66 = 0$, $40X - 18Y = 214$. find (a) mean value of X & Y ,
 (b) the Correlation Co-efficient bet' X & Y
 (c) the standard deviation of Y .

Ans:-

We know that both the line of regression passes through (\bar{X}, \bar{Y}) , we have $8\bar{X} - 10\bar{Y} + 66 = 0$ & $40\bar{X} - 18\bar{Y} - 214 = 0$

$$8\bar{X} - 10\bar{Y} + 66 = 0 \quad \text{--- (1)}$$

$$\underline{40\bar{X} - 18\bar{Y} - 214 = 0} \quad \text{--- (2)}$$

$$\text{eqn(1)} \times 5 = 40\bar{X} - 50\bar{Y} + 330 = 0$$

$$\underline{\underline{40\bar{X} - 18\bar{Y} - 214 = 0}}$$

$$-32\bar{Y} + 116 = 0$$

$$\Rightarrow -32\bar{Y} = -116$$

$$\Rightarrow \bar{Y} = \frac{-116}{-32} = 17$$

Putting the value of \bar{Y} in eqn (1)

$$8\bar{X} - 10\bar{Y} + 66 = 0$$

$$\Rightarrow 8\bar{X} - 10 \times 17 + 66 = 0$$

$$\Rightarrow 8\bar{X} - 170 + 66 = 0$$

$$\Rightarrow 8\bar{X} = 104$$

$$\Rightarrow \bar{X} = \frac{104}{8} = 13$$

Solving above eqn, we get $\bar{X} = 13$, $\bar{Y} = 17$

Let $8x - 10y + 66 = 0$ & $40x - 18y - 214 = 0$ be the line of regn of y on x and x on y .

So, These eqn can be put in the form

$$Y = \frac{8}{10}x + \frac{66}{10}$$

$$X = \frac{18}{40}y + \frac{214}{40}$$

So, b_{yx} will be regression co-efficient of y on x .

$$b_{yx} = \text{regression co-efficient of } y \text{ on } x = \frac{8}{10}$$

$$b_{xy} = \text{regression co-efficient of } x \text{ on } y = \frac{18}{40}$$

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

$$= \pm \sqrt{\frac{8}{10} \times \frac{18}{40}} = \pm \sqrt{\frac{36}{50}} = \pm 0.6$$

$$= \pm \sqrt{\frac{9}{25}} = \pm \frac{3}{5} = \pm 0.6$$

Since both the regression co-efficients are positive, we take $r = \pm 0.6$

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

$$\Rightarrow \frac{4}{5} = \frac{3}{5} \times \frac{\sigma_y}{\sigma_x}$$

$$\Rightarrow \frac{4}{5} \times \frac{3}{5} \times 3 = \sigma_y$$

$$\Rightarrow \sigma_y = 4$$

Co-efficient of Determination:

→ The square of the correlation co-efficient is called as Co-efficient of determination.

Co-efficient of determination = $r^2 = \frac{\text{Explained Variance}}{\text{Total Variance}}$

Ex:-

If the value of $r=0.9$ then $r^2=0.81$. It means that 80% of the variation in the dependent variable has been explained by the independent variable. It gives the percentage variation in the dependent variable that is accounted by the independent variable.

$$y = ax + b$$

↓ ↓
Dependent variable Independent variable

Difference bet' correlation and regression;

Correlation	Regression
① Correlation Co-efficient can be used to measure the relationship between two variables.	① Regression analysis is to study the nature of relationship between two variables.
② In correlation analysis we cannot say that one variable is the cause and the other is the effect.	② In regression analysis it is possible to study cause and effect relationship.
③ In correlation analysis both r_{xy} and r_{yx} are symmetric i.e. ($r_{xy} = r_{yx}$)	③ In regression analysis b_{yx} and b_{xy} are not symmetric, i.e. ($b_{yx} \neq b_{xy}$)
④ Correlation co-efficient is independent of change of origin & scale.	④ Regression Co-efficient are independent of change of origin but not scale.

UNIT - 4

Principle of least square:—

The Principle of least square states that by getting the sum of square of the errors of a minimum value, the most probable value of a system of unknown quantity can be obtain.

→ The principle of least square states that the based value of a_1, a_2, \dots, a_m in the relation

$$Y_i = f(x_i, a_1, a_2, \dots, a_m) + E_i^2$$

where $i = 1, 2, 3, \dots, n$

are those values of a_1, a_2, \dots, a_m values which minimize the sum of square of the errors.

Hence, to find the values of a_1, a_2, \dots, a_m such that

$$\text{minimized } E = \sum_{i=1}^n E_i^2$$

$$= \sum_{i=1}^n [Y_i - f(x_i, a_1, a_2, \dots, a_m)]^2$$

so, this can be done by solving the normal eq

$$\frac{\partial E}{\partial a_1} = 0, \frac{\partial E}{\partial a_2} = 0$$

$$Y = ax + b + e$$

$$E = \sum_{i=1}^n e$$

$$E = \sum (Y_i - ax_i - b)^2$$

The Partial differentiation of E with respect to a and b and equating to zero

$$\frac{dE}{da} = -2\sum x(Y - ax - b) = 0 \text{ and } \frac{dE}{db} = -2\sum (Y - ax - b) = 0$$

$$\Rightarrow \sum xy - a\sum x^2 - b\sum x = 0 \quad \Rightarrow \sum Y - a\sum x - nb = 0$$

$$\Rightarrow \sum xy = a\sum x^2 + b\sum x \rightarrow (i) \quad \Rightarrow \sum Y = a\sum x + nb \rightarrow (ii)$$

Solving eqn (i) and eqn (ii) we get \hat{a} and \hat{b} , the best fitting line is $y = \hat{a}x + \hat{b}$.

Fitting of straight line eqn:

The straight line eqn is given by $y = ax + b$ using the principle of least square, we have to determine the constant a and b . Then

$$E = \sum (Y - ax - b)^2$$
 is minimum

The partial differentiation of E with respect to a and b and equating to zero.

$$\frac{dE}{da} = -2\sum x(Y - ax - b) = 0 \quad \frac{dE}{db} = -2\sum (Y - ax - b) = 0$$

$$\Rightarrow \sum xy - a\sum x^2 - b\sum x = 0 \quad \Rightarrow \sum Y - a\sum x - nb = 0$$

$$\Rightarrow \sum xy = a\sum x^2 + b\sum x \rightarrow (i) \quad \Rightarrow \sum Y = a\sum x + nb \rightarrow (ii)$$

Solving eqn (i) and eqn (ii) we get \hat{a} and \hat{b} , the best fitting straight line is $y = \hat{a}x + \hat{b}$.

Ex: — Fit the straight line of the form $y = ax + b$.

	<u>x</u>	<u>y</u>	<u>x^2</u>	<u>xy</u>
1	8	8	64	64
3	12	36	144	108
5	15	75	225	187.5
7	17	119	289	233
8	18	144	324	259.2
10	20	200	400	300
$\sum y = 34$		$\sum y = 90$	$\sum xy = 582$	$\sum x^2 = 948$

$$\text{The normal eq's are } 582 = 248a + 34b \quad (i)$$

$$90 = 34a + 6b \quad (ii)$$

Solving eq's (i) & (ii)

$$\text{eq}^{\prime\prime} (i) \times 3 = 1746 = 744a + 102b$$

$$\text{eq}^{\prime\prime} (ii) \times 17 = 1530 = 578a + 102b$$

$$\text{eq}^{\prime\prime} (i) - (ii) = 216 = 166a$$

$$= \frac{216}{166} = a$$

$$= 1.301 = a$$

Pulling value of a in eq (ii), we get

$$90 = 34a + 6b$$

$$\Rightarrow 90 = 34 \times 1.301 + 6b$$

$$\Rightarrow 90 - 44.2 = 6b$$

$$\Rightarrow \frac{45.8}{6} = b$$

$$\Rightarrow 7.63 = b$$

Solving eq (i) & (ii) we get a & b ; $y = 1.301x + 7.63$

Fitting of 2nd degree polynomial eq:

The 2nd degree polynomial eq is given by,

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

Using the principle of least square, we have to determine the constant a , b and c .

$$E = \sum (y - a - bx - cx^2)^2 \text{ is minimum}$$

The partial difference of E with respect to a , b & c and equating to zero.

$$\frac{dE}{da} = 2 \sum (y - a - bx - cx^2) = 0 \text{ is minimum}$$

$$\Rightarrow \sum y - na - b \sum x - c \sum x^2 = 0 \quad (i)$$

$$\frac{dE}{db} = -2 \sum x (y - a - bx - cx^2) = 0$$

$$\Rightarrow \sum xy - a \sum x - b \sum x^2 - c \sum x^3 = 0 \quad (ii)$$

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

$$\frac{dE}{dc} = -2\sum x^2 (y - a - bx - cx^2)$$

$$\Rightarrow \sum xy - a\sum x^2 - b\sum x^3 - c\sum x^4$$

$$\Rightarrow \sum xy = a\sum x^2 + b\sum x^3 + c\sum x^4 \quad \text{--- (ii)}$$

Solving eqn (i), (ii) and (iii) we get \hat{a}, \hat{b} & \hat{c} the best

fitting line eqn is $y = \hat{a} + \hat{b}x + \hat{c}x^2$.

Fitting of exponential curve—

$$(i) \quad y = ab^x \quad \text{--- (i)}$$

Taking logarithms both the side,

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

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

$$\text{Let } u = \log y, A = \log a \text{ & } B = \log b$$

$$u = A + Bx$$

The normal eqn for estimating A and B are

$$\sum u = nA + B\sum x \quad \text{--- (ii)}$$

$$\sum ux = A\sum x + B\sum x^2 \quad \text{--- (iii)}$$

Solving eqn (ii) and (iii), we get

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

$$b = \text{Antilog}(B)$$

With these value of a and b in eqn (i), we get the curve of best fit.

$$(ii) \quad \text{ab}^x = ae^{bx}$$

Taking logarithms both the side,

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

$$\log y = \log a + bx \log e$$

Let $U = \log y$, $A = \log a$ & $B = \log b$

$$U = A + BX$$

The normal eqn for estimating A and B are

$$\sum U = nA + B\sum X \quad (\text{ii})$$

$$\sum UX = A\sum X + B\sum X^2 \quad (\text{iii})$$

Solving eqn (ii) and (iii), we get a and b

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

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

With these value of a and b in eqn (i) we get

the curve of best fit.

Ex: Fit an exponential curve of the form

$y = ab^x$ to the following data.

<u>x</u>	<u>y</u>	<u>$U = \log x$</u>	<u>$\sum u$</u>	<u>$\sum u^2$</u>
1	1.0	0.000	0	0
2	1.2	0.0792	0.1584	0.0000
3	1.8	0.253	0.7659	0.0000
4	2.5	0.3979	1.5916	0.0000
5	3.6	0.5536	2.7815	0.0000
6	4.7	0.6721	4.0326	0.0000
7	6.6	0.8193	5.7365	0.0000
8	9.1	0.9590	7.6720	0.0000
$\Sigma x = 36$		$\Sigma U = 3.7393$	$\Sigma u = 22.7385$	$\Sigma u^2 = 204$

The normal eqn for estimating A and B are

$$\sum U = nA + B\sum x \quad (\text{i})$$

$$3.7393 = 8A + 36B \quad (\text{i})$$

$$\sum UX = A\sum x + B\sum x^2 \quad (\text{ii})$$

$$\Rightarrow 22.7385 = 36A + 204B \quad (\text{ii})$$

Solving eq (i) and (ii), we get.

$$\text{eqn (i)} \times 9 = 33.65 = 72A + 324B$$

$$(\text{ii}) \times 2 = 45.47 = 72A + 408B$$

$$\Rightarrow -11.82 = -84B$$

$$\Rightarrow \frac{-11.82}{-84} = B$$

$$\Rightarrow 0.1407 = B \Rightarrow b = \text{Antilog}(0.1407)$$

$$= 1.380$$

Putting the value of B in eqn (i)

$$\text{we get, } 3.7393 = 8A + 36 \times 0.1407$$

$$\Rightarrow 3.7393 = 8A + 5.0652$$

$$\Rightarrow 3.7393 = 5.0652 = 8A$$

$$\Rightarrow \frac{-1.3259}{8} = A$$

$$\Rightarrow 0.1657 = A \Rightarrow a = \text{Antilog}(0.1657) = 0.6827$$

The exponential curve eqn is

$$Y = (0.6827 \cdot 1.380)^X$$

Theory of attributes:

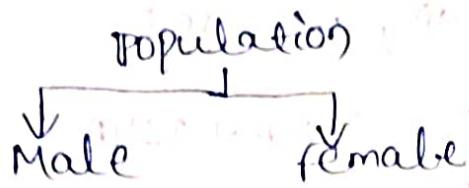
→ Theory of attributes deals with qualitative characteristics calculated by using quantitative measurement. e.g. Honesty, beauty, habit of smoker, intelligence.

→ In the study of attributes, the objects are classified according to the absence or presence of the attribute in them.

Dichotomous classification:

If the universe is divided into two sub class and known as Dichotomous classification.

e.g -

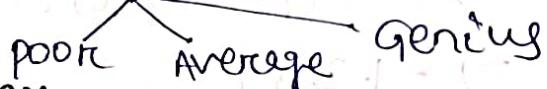


Manifold classification:-

→ If the universe is divided in to more than two classes are known as manifold classification.

e.g - The attribute intelligent is divided in to 4 types.

1. Genius
2. Average intelligent
3. very intelligent
4. Below average intelligent



Types of attributes:-

There are two types of attributes.

- (1) Positive attribute. (presence of attribute)
- (2) Negative attribute. (Absence of attribute)

Symbols and Notations:-

→ The presence of attributes are represented by Capital letters i.e. A, B, C, D, ..., Z.

→ The absence of attributes are represented by greek letters i.e. α , β , γ etc.

Number of attributes:-

1 - Attributes - 3 classes - N, A, d

2 - Attributes - 9 classes - N, A, B, α , β , γ , AB, α B, α A

3 - Attributes - 27 classes - N, A, B, C, α , β , γ , δ , α B, α C, α A, β B, β C, β A, γ B, γ C, γ A, δ B, δ C, δ A, α B γ , α C γ , α A γ , β B γ , β C γ , β A γ , γ B δ , γ C δ , γ A δ , δ B γ , δ C γ , δ A γ

4 - Attributes - 81 classes - N, A, B, C, D, α , β , γ , δ , ϵ , α B, α C, α D, α A, β B, β C, β D, β A, γ B, γ C, γ D, γ A, δ B, δ C, δ D, δ A, ϵ B, ϵ C, ϵ D, ϵ A, α B γ , α C γ , α D γ , α A γ , β B γ , β C γ , β D γ , β A γ , γ B δ , γ C δ , γ D δ , γ A δ , δ B γ , δ C γ , δ D γ , δ A γ , α B δ , α C δ , α D δ , α A δ , β B δ , β C δ , β D δ , β A δ , γ B ϵ , γ C ϵ , γ D ϵ , γ A ϵ , δ B ϵ , δ C ϵ , δ D ϵ , δ A ϵ , α B γ δ , α C γ δ , α D γ δ , α A γ δ , β B γ δ , β C γ δ , β D γ δ , β A γ δ , γ B δ ϵ , γ C δ ϵ , γ D δ ϵ , γ A δ ϵ , δ B γ ϵ , δ C γ ϵ , δ D γ ϵ , δ A γ ϵ

5 - Attributes - 243 classes - N, A, B, C, D, E, α , β , γ , δ , ϵ , ζ , α B, α C, α D, α E, α A, β B, β C, β D, β E, β A, γ B, γ C, γ D, γ E, γ A, δ B, δ C, δ D, δ E, δ A, ϵ B, ϵ C, ϵ D, ϵ E, ϵ A, ζ B, ζ C, ζ D, ζ E, ζ A, α B γ , α C γ , α D γ , α E γ , α A γ , β B γ , β C γ , β D γ , β E γ , β A γ , γ B δ , γ C δ , γ D δ , γ E δ , γ A δ , δ B γ , δ C γ , δ D γ , δ E γ , δ A γ , ϵ B δ , ϵ C δ , ϵ D δ , ϵ E δ , ϵ A δ , ζ B δ , ζ C δ , ζ D δ , ζ E δ , ζ A δ , α B γ δ , α C γ δ , α D γ δ , α E γ δ , α A γ δ , β B γ δ , β C γ δ , β D γ δ , β E γ δ , β A γ δ , γ B δ ϵ , γ C δ ϵ , γ D δ ϵ , γ E δ ϵ , γ A δ ϵ , δ B γ ϵ , δ C γ ϵ , δ D γ ϵ , δ E γ ϵ , δ A γ ϵ , ϵ B δ ζ , ϵ C δ ζ , ϵ D δ ζ , ϵ E δ ζ , ϵ A δ ζ , ζ B δ ϵ , ζ C δ ϵ , ζ D δ ϵ , ζ E δ ϵ , ζ A δ ϵ

Class frequency:—

Class Frequency means no. of frequency assign in each others.

e.g. (A), (AB), (BC), (AC) \rightarrow +ve frequencies

(D), (B), (D), (DB) \rightarrow -ve frequencies

(ABD), (BDC), (ABC) \rightarrow Contingency frequencies

Some algebraic Expressions:—

$$(A) = (AB) + (AD)$$

$$(B) = (AB) + (DB)$$

$$(D) = (DB) + (DD)$$

$$(P) = (AP) + (DP)$$

formula for any frequency classes:—

No. of classes of n th order frequency = $C_n \times 2^n$ classes

n = no. of attributes

$$\sum C_n = 1, 0! = 1, 1! = 1$$

n = order of the class

// find the no. of order 2 classes if the no. of attributes is 2.

No. of 2 order class = $C_2 \times 2^2$

$$= 2 \times 4$$

Contingency Table:—

Contingency matrix shows the relation between two variables or attributes.

\rightarrow It is first used by Karl Pearson.

Attributes	A	D	Total	(AB), (A) (B) \rightarrow +ve frequency
B	(AB)	(DB)	(B)	(DB), (D), (B) \rightarrow -ve frequency
P	(AP)	(DP)	(P)	(DP), (A), (P) \rightarrow Contingency
Total	(A)	(D)	N	Frequency

Conditions for consistency of data:—

Rule:— (i) Frequency of every class is greater than or equals to zero. i.e $A \geq 0$,
(ii) Class frequency is less than or equals to N (population). i.e $A \leq N$.

Ex:— Determine the consistency in the given data.

$$(AB) = 70, (\alpha B) = 40, \alpha P = 60, (B) = 100$$

Attributes	A	α	Total
B	NB 60	αB 40	B 100
P	NP 70	αP 60	P 130
Total	A 130	α 100	$N=230$

$$AB = B - \alpha B = 100 - 40 = 60, AP = \alpha P + AB = 70 + 60 = 130$$

$$B = 100, P = 130, A = 130, \alpha = 100, N = 230.$$

Non frequency is negative and all class frequency is less than N . so the data is inconsistent independent of attributes.

→ If the attributes are said to be independent because there does not exist any relationship between them.

ex:— beauty and intelligence, gender and success

→ If the attributes ~~are~~ A & B are said to be independent then the proportion of AB in the population is equal to product of the proportion of A & B in the population. $\frac{(AB)}{N} = \frac{(A)}{N} \cdot \frac{(B)}{N}$

Methodology to check association and independence attributes:

① Proportion Method:

$$(i) \frac{(AB)}{(B)} = \frac{(AB)}{(P)} \rightarrow \text{Two attributes are independent}$$

$$(ii) \frac{(AB)}{(B)} > \frac{(AB)}{(P)} \rightarrow \text{Positive association}$$

$$(iii) \frac{(AB)}{(B)} < \frac{(AB)}{(P)} \rightarrow \text{Negative Association}$$

Ex: — $(AB) = 100$, $(B) = 10$, $(A) = 150$, $(P) = 15$.
Find out how A and B are associated.

Ans: —

$$\frac{(AB)}{(B)} = \frac{100}{10} = 10$$

$$\frac{(AB)}{(P)} = \frac{100}{15} = 10$$

$$\frac{(AB)}{(B)} = \frac{(AB)}{(P)}$$

\Rightarrow Two attributes are independent.

② Comparison Method:

$$(i) (AB) = \frac{(A)(B)}{N} \rightarrow \text{Two attributes are independent}$$

$$(ii) (AB) > \frac{(A)(B)}{N} \rightarrow \text{Positive association}$$

$$(iii) (AB) < \frac{(A)(B)}{N} \rightarrow \text{Negative association}$$

Ex: — find out the association in the given data.

$$N = 106, (A) = 70, (B) = 36, (AB) = 20$$

Ans: — $(AB) = 20$

$$\frac{(A)(B)}{N} = \frac{70 \times 36}{106} = 23.77$$

$$(AB) < \frac{(A)(B)}{N} \rightarrow \text{Negative association}$$

③ Yule's co-efficient of association Method:

$$Q = \frac{(AB)(\alpha\beta) - (A\beta)(\alpha B)}{(AB)(\alpha\beta) + (A\beta)(\alpha B)}$$

→ Yule's Co-efficient of association is lies between -1 to +1.

Representation of values of Q:

- If Q lies bet' 0 to 1 → positive association
 - If Q lies bet' -1 to 0 → negative association
 - If $Q = 0$ → independent association
 - If $Q = +1$ → perfectly positive association
 - If $Q = -1$ → perfectly negative association
- Q// $(AB) = 60$, $(A\beta) = 20$, $(A\beta) = 80$, $(\alpha B) = 30$. Find the association.

$$Q = \frac{(AB)(\alpha\beta) - (A\beta)(\alpha B)}{(AB)(\alpha\beta) + (A\beta)(\alpha B)} = \frac{60 \times 20 - 80 \times 30}{60 \times 20 + 80 \times 30}$$

$$= \frac{1200 - 2400}{1200 + 2400} = \frac{-1200}{3600}$$

⇒ Negatively association.

Co-efficient of colligation:

It is another measure of association suggested by Prof. Yule's. It is denoted by γ .

$$1 - \frac{(AB)(\alpha\beta)}{\sqrt{(AB)(\alpha\beta)}} \leftarrow (i) \rightarrow (ii)$$

$$\text{or } 1 + \frac{(AB)(\alpha\beta)}{\sqrt{(AB)(\alpha\beta)}} \leftarrow (iii) \rightarrow (iv)$$

Relation bet' Q & γ :

$$Q = \frac{2\gamma}{1 + \gamma^2}$$

limits of γ is -1 to +1.

UNIT-1

Defⁿ of Statistics:

Statistics is a science of collection, organization, presentation, analysis and interpretation of numerical data.

→ Father of Modern Statistics is Prof. Sir R.A. Fisher.

→ Father of Indian Statistics is Prof. Prasanta Chandra Mahalanobis.

Scope of Statistics:

→ Statistics are numerical statement of facts capable of analysis and interpretation as well as study of method used in collection, organization, presentation, analysis and interpretation of numerical data.

→ Statistics are used in different sectors such as statistics in Economics, in education, in planning, in commerce, in business etc.

Limitation of Statistics:

→ Statistics does not study for qualitative phenomenon.

→ Statistics does not study for individuals.

→ Statistics laws are not exact.

→ Statistics is liable to misuse.

Function of Statistics:

(1) Collection of data:

→ The 1st step of an investigation is a collection of data.

→ Careful collection is needed because further analysis is based on this.

→ Data must be reliable. If the data is not reliable, the collected data are faulty. Therefore the good investigator must take special care in collection of data.

(2) Organization of Data:-

→ Organization of data refers to the arrangement of failures. In such a form of comparisons of the mass similar data may be facilitated and further analysis may be possible.

(3) Presentation of Data:-

→ The collection of data are generally in raw form and need to be classified and tabulated before they can analysed.

→ Therefore the collected data are to be presented in tabular or diagrammatic or graphical form.

(4) Analysis of Data:-

→ After the presentation of data the next step is to analyse the presented data.

→ Analysis includes some condensation, summarization and conclusion through measures of central tendency, dispersion, skewness etc.

(5) Interpretation of Data:-

→ Interpretation of data takes the result of analysis and makes inference to a particular research studies and draw conclusions about these relations.

What do you mean by Data:

- Data is a collection of information or facts on a certain enquiry is called as Data.
- There are two types of Data.
 - (i) Primary Data
 - (ii) Secondary Data

(i) Primary Data:

- Primary Data refers to the first hand data gathered by researcher himself.
- The primary data are usually in a raw or bulky form (wide, heavy, large)

(ii) Secondary Data:

- Secondary Data is a data that has already been gathered or collected from others.
- Secondary data is accessible in the form of data collected from different sources such as government, publication, census, books, journals, articles, websites etc.

Methods of collecting primary Data:

- (1) Direct personal investigation
- (2) Indirect oral interview
- (3) Information received through local agencies
- (4) Mailed questionnaire Method
- (5) Schedules sent through enumerators

Direct Personal Investigation:

Direct Personal Investigation or personal interview is a method of collecting primary data through

which the investigator contacts the informant directly to collect data by conducting on the spot enquiry.

(ii) Indirect oral interview:

Indirect oral investigation is a method of collecting primary data through which the investigator approaches third parties who are in possession of required information about the subject of enquiry.

(iii) Information received through local agencies:

It is a method of collecting primary data under which the investigator appoints local persons or correspondents at different places. These correspondents collect information in their own way and provide it to the investigator.

(iv) Method questionnaire Method:

The questionnaires are the fundamental instrument for gathering information in research. fundamentally, it is a bunch of standardized questions, frequently called items, which follow a decent plan to gather individual information around at least one explicit theme.

(v) Schedule sent through enumerators:

In this method, list of questions or schedules are sent to the informants through the enumerators.

Classification:—

The arrangement of data in to different classes, which are to be determined depending upon the nature, objectives and scope of the enquiry.

Ex:— The number of students registered in UNIC college in the academic year 2022-23 may be classified on the basis of any of the following criteria.

- | | |
|----------------|------------------------------|
| (i) Sex | (iv) State which they belong |
| (ii) Age | (v) Height |
| (iii) Religion | (vi) Weight |

→ The data are classified primarily depending upon the objective and the purpose of the enquiry. The data can be classified in to following 4 basis.

- (i) Geographical Classification
- (ii) Chronological Classification
- (iii) Qualitative Classification
- (iv) Quantitative Classification

(i) Geographical classification:—

In this classification the basis of classification is the geographical or locational differences between the various items in the data like:—states, cities, regions, areas etc.

Ex:— The yield of agricultural output for hector for different countries in given period; Density of population in different cities of India

(ii) Chronological classification:—

Chronological Classification is one of the data are classified on the basis of differences time.

Ex:- The population of any country for different years.

(iii) Qualitative classification:— When the data are classified according to some qualitative phenomenon which are not capable of quantitative measurement like;— honesty, beauty, intelligence etc.

(iv) Quantitative classification:— When the data are classified on the basis of Quantitative Phenomenon which are capable of Qualitative measurement like;— height, weight, income, expenditure etc.

Variable:— The qualitative phenomenon under study like marks in a test, height, weight, number of student in a class etc. is known as Variable.

→ There are two types of Variable:

- ① Discrete variable
- ② Continuous variable

Discrete Variable:—

Those variables which can not take all possible values within a given specified range is known as Discrete variable.

Ex:- Marks in a test, number of student in a class

Continuous Variable:—

Those variables which can take all possible values within a given specified range is known as Continuous Variable.

Ex:- Height, weight, Age etc.

Class Interval:

Class interval in a numerical width in a class and frequency distribution each class is specified with two extreme value. This is called as class interval.

- The upper limit smaller one is Lower limit and the bigger one is upper limit.
- There are two type of class interval

(i) Inclusive

(ii) Exclusive

Inclusive Class Interval:

The class of type like : 1-9, 10-19, 20-29 etc. in which both upper limit and lower limit includes in a class is known as inclusive class interval.

Exclusive Class Interval:

The class of type like - 1-10, 10-20, 20-30 etc. in which lower limit is include and the upper limit is excluded from the respective class and immediately included in the next class. It is known as exclusive class interval.

Mid Value:

Mid value of any class can be obtained by dividing the sum of lower limit and upper limit.

$$\text{Midvalue} = \frac{\text{Lower limit} + \text{Upper limit}}{2}$$

Ex:- Age Midvalue

10-20	15	15 is a good breaking point. A set of data 20-30 contains 25 which is not older than 30.
30-40	35	35 is not 40

Class boundaries:

- In exclusive type classification or ungrouped frequency distribution there are the gap between the upper limit of any class and lower limit of the succeeding class. So, there is need to convert the data in to a continuous frequency distribution.
 - The upper and lower limit of new exclusive type class are called as class boundaries.
 - If d is the gap between upper limit of any class and lower limit of succeeding class, so, the class boundaries of any class is given by
 - * Upper class boundaries = upper limit + $\frac{d}{2}$
 - * Lower class boundaries = Lower limit - $\frac{d}{2}$
- Open END Class:
- If the lower limit of first class and upper limit of last class are not specified and this classes one of the limits is missing is called as open end classes.
- Ex:- Less than 30, age above 60.

Tabulation of Data:

After classification the next step in the purpose of summarisation of data is to put the classification data in rows and columns having special characteristics on a piece of paper. Such a representation of data is known as Tabulation.

Salient features of a Table:

- A table should have a table number, table title and subtitle for rows and columns of the body of the table.

- A good table should be simple to understand.
- The numerical figures represented in the table should be compare the other related figures without include waiting for further calculation.
- A good table should be proper size and shape with proper spacing for rows and columns and should have attractive get up.
- The data incorporated in the table should be order alphabetically, geographically or depending upon the necessity of facilitate ready for comparison.

Types of Table:

(Classification of tables according to purpose—

(1) General purpose table:

It also known as information table and reference table. It providing a symmetric manner usually in chronological order. All information connected with enquiry. This table are also known as ~~General purpose table~~ Master table.

(2) Summary Table:

This table are constructed and derived from the master table with certain purpose and are useful for analytical and comparing studies, involving the study relationships among variables. calculating of analytical statistics like ratio, percentage, etc. are in co-operated to this tables.

Classification of tables according to nature of data

There are two types.

(1) Primary table:

The primary tables represent data, in original form without arithmetical rounding up.

(2) Derived table:

This table condensed figures and entries.

E.g. Ratio, percentages etc. from the primary data.

Classification of table according to characteristics:

(1) Simple table:

A simple table is the one way classification dealing with only one characteristics.

(2) Complex table:

A complex table is based on more than one characteristics.

Graphical Representation of data:

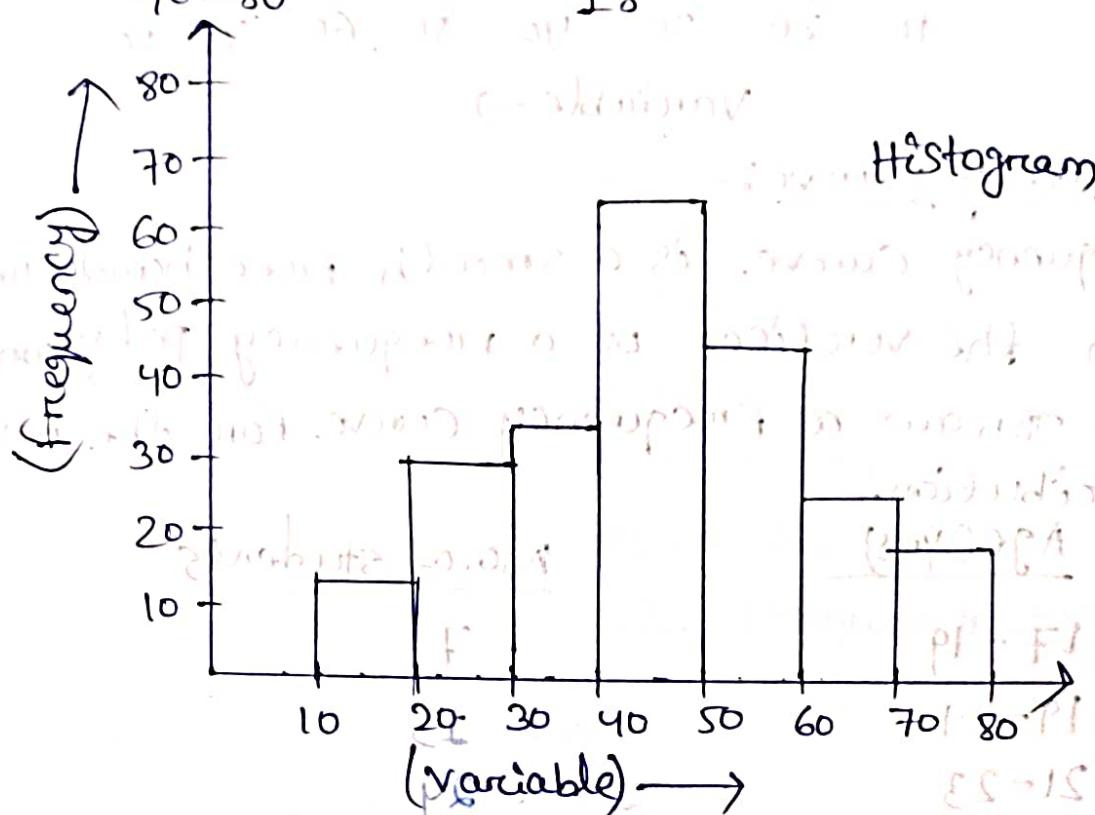
(1) Histogram:

Histogram is a graphical representation of continuous frequency distribution. It consists of series of adjacent vertical rectangle on the sections of the horizontal axis (x-axis) with bases sectors equal to width of the corresponding class interval and heights are taken the areas of a rectangles are equal to the frequencies of the corresponding class.

The variable values are taken along the x-axis and the frequency along the y-axis.

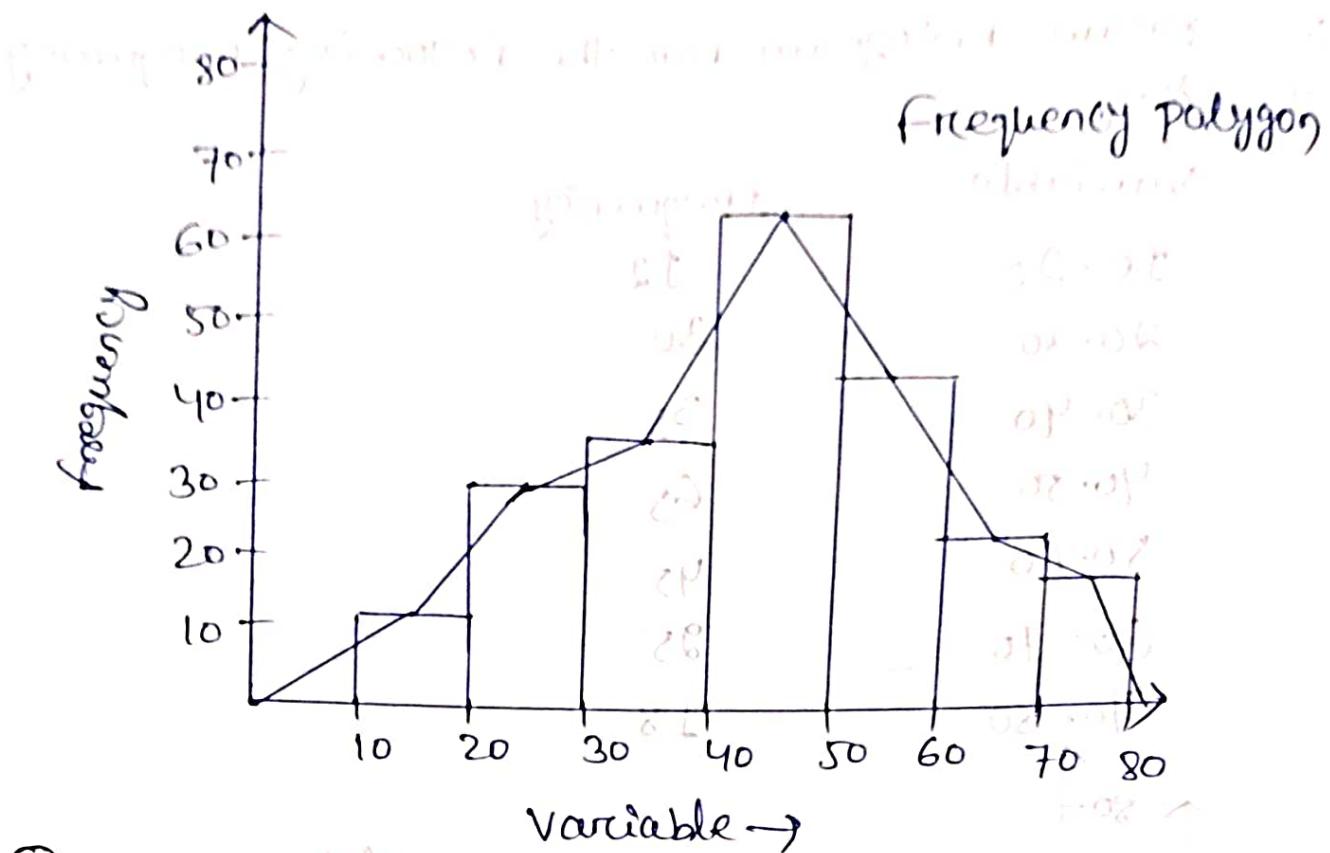
Ex:- Draw histogram for the following frequency distribution.

<u>variable</u>	<u>frequency</u>
10-20	12
20-30	30
30-40	35
40-50	65
50-60	45
60-70	25
70-80	18



(2) Frequency polygon:

Frequency polygon is obtain on plotting the frequencies of the vertical axis (Y-axis) against the corresponding values of the variable on the horizontal axis (X-axis) and joining the midpoints obtained by the straight line.



③ Frequency Curve:

frequency curve is a smooth free hand curve through the vertices of a frequency polygon.

Ex: — Draw a frequency curve for the following distribution.

Age (Yrs)

No. of students

17-19

7

19-21

13

21-23

11

23-25

20

25-27

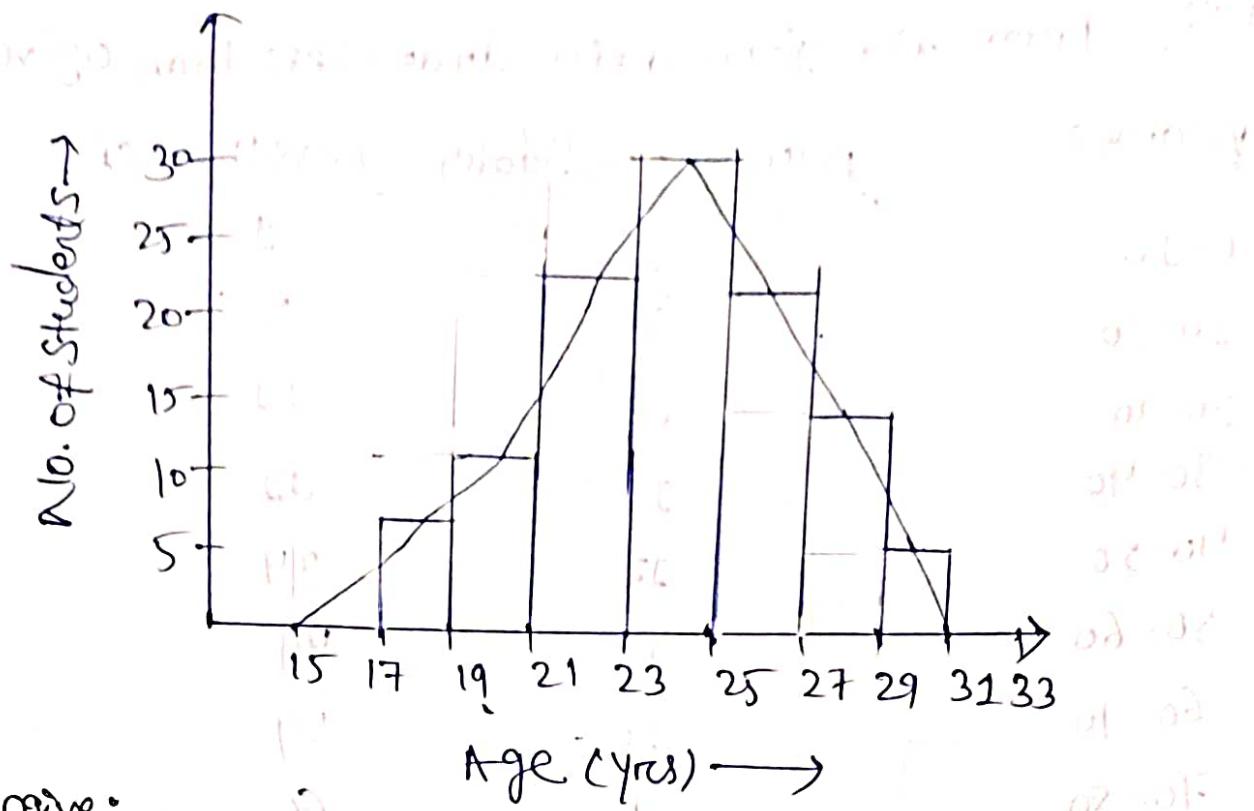
22

27-29

15

29-31

6



→ The graphical representation of the partition value is known as ogive. In other words, ogive is a graphical representation of cumulative frequency distribution of continuous variable. It consist plotting the cumulative frequency (along the y-axis) against the class boundaries of the x-axis's. So, there are two types of cumulative frequency distribution.

- (i) Less than ogive
- (ii) More than ogive

Less than ogive:

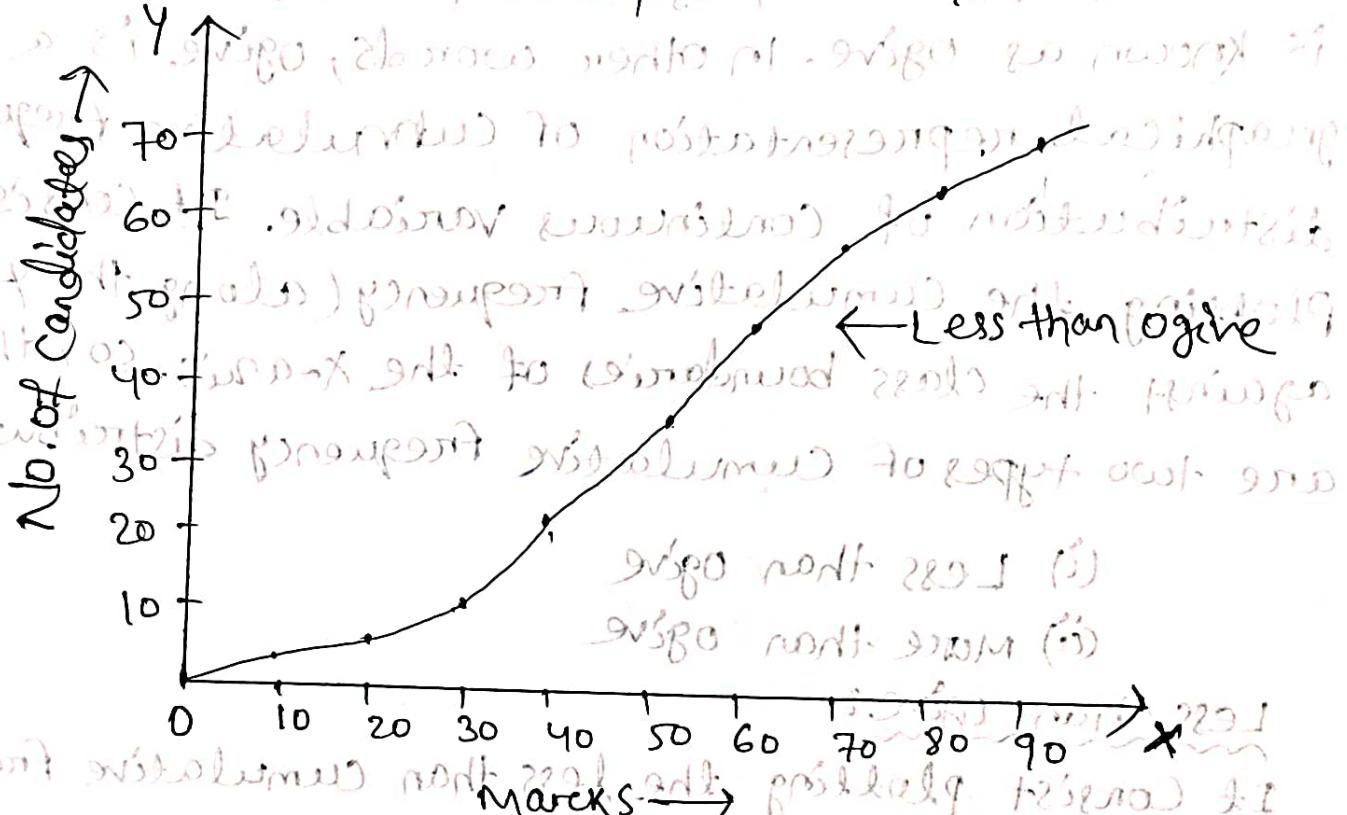
It consist plotting the less than cumulative frequency against the upper class boundaries of the respective classes. The points are obtained and are joined by smooth hand curve.

→ Less than ogive is an increasing curve, sloping upward from left to right and the shape is 'S' type.

Ex:- From the given data draw less than ogive.

Marks

<u>Marks</u>	<u>No. of candidates</u>	<u>Less than cf</u>
0-10	2	2
10-20	3	5
20-30	6	11
30-40	11	22
40-50	12	34
50-60	15	49
60-70	10	59
70-80	7	66

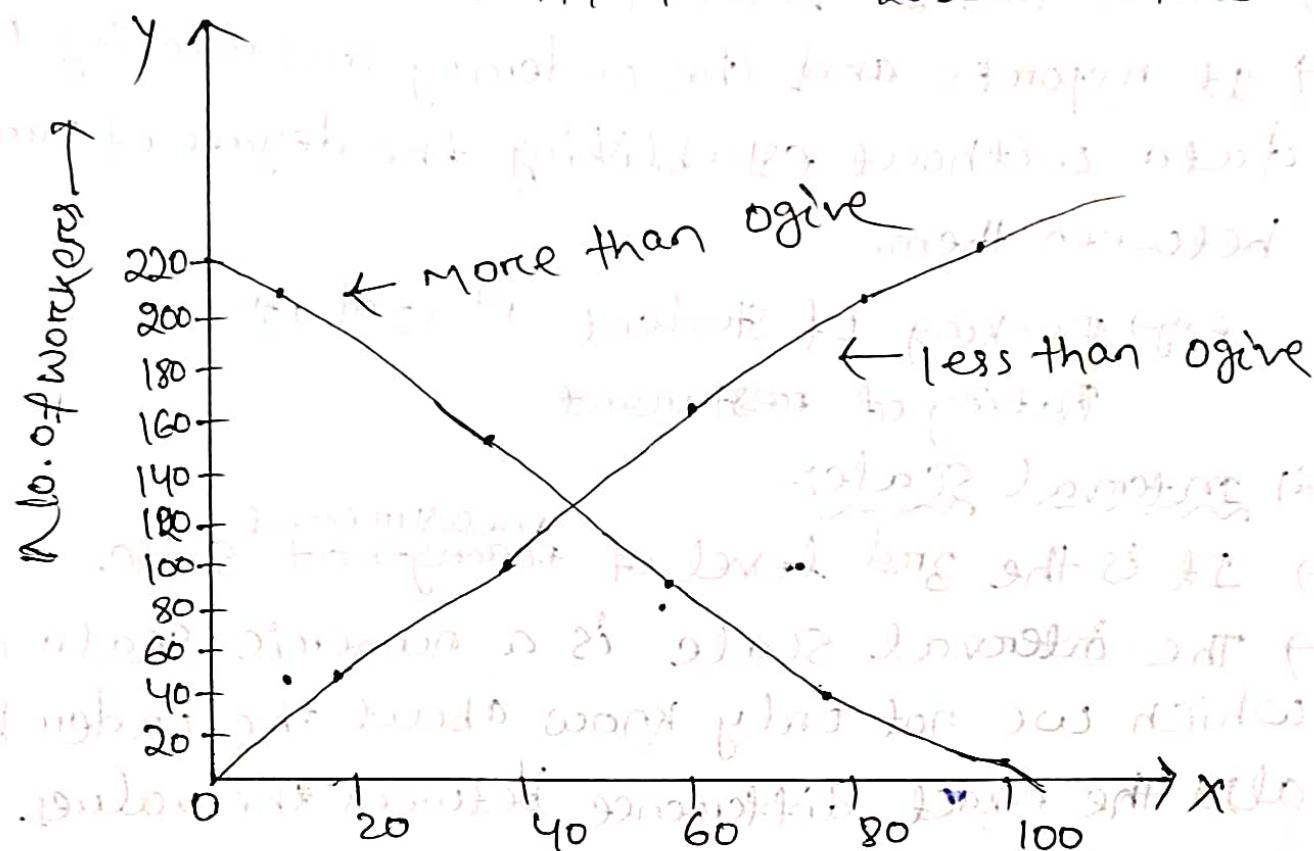


② More than Ogive: In more than ogive or more than cumulative frequency are plotted against the lower class boundaries of that respective classes. The points are obtained and are joined by a smooth hand curve or free hand

→ More than ogive is a decreasing curve and slopes downward from left to right.

Ex:— Draw less than and more than ogive and find the value of median.

<u>Weekly wages</u>	<u>No. of workers</u>	<u>less than Cf</u>	<u>More than Cf</u>
0-20	41	41	201
20-40	$92 - 41 = 51$	92	$201 - 41 = 160$
40-60	$156 - 92 = 64$	156	$160 - 51 = 109$
60-80	$194 - 156 = 38$	194	$109 - 64 = 45$
80-100	$201 - 194 = 7$	201	$45 - 38 = 7$



Scales of measurement:

The Scales of Measurement is a classification of that describe the nature of information within the numbers assign to the variable.

→ There are four types of scale.

- Measurement Scale
- (i) Nominal Scale
 - (ii) Ordinal Scale
 - (iii) Interval Scale
 - (iv) Ratio Scale

(i) Nominal Scale:-

- It is the 1st level of measurement scale.
- It deals with non-numeric variable or the numbers that do not have any value.

e.g. → Sex, Region, Cast etc.

(ii) Ordinal Scale:-

- It is the 2nd level of measurement scale.
- It reports ~~and~~ the ordering and ranking of data without establishing the degree of variation between them.

e.g. → Ranking of student 1st, 2nd, 3rd

Rating of restaurant

(iii) Interval Scale:-

- It is the 3rd level of ~~measurement~~ scale.
- The interval scale is a numeric scale on which we not only know about the order but also the exact difference between the values.

e.g. → Celsius, temperature

(iv) Ratio Scale:-

- It is the 4th level of measurement scale which is quantitative.
- It allows the researcher to compare the differences between unit.

→ One of the most important property that is zero position indicates the absence of quantity being measured.

e.g) Height, weight etc.

Concepts of statistical Population & Sample:-

Population:- Group of individual under study is known as population.

Ex:- All people living in India indicates the population of India

Sample:- It indicates one or more observations that are drawn from the population.

Ex:- All Sikhs people living in India is the sample of population.

Tabulation:-

After classification, the next step in the purpose of summarization of data is to put the classified data in rows and columns having special characteristics on a piece of paper. Such representation of data is called as tabulation.

Graphical:-

The graphs provide an alternative method of representing data in a condensed & summary form. A graph is a scale-dependent geometrical figures and provides a visual representation of statistical data.