

# STATISTICS

## Classification of Data.

(1) Grouped Data — which can be organised.

0-10 -	3
eg. 10-20 -	4
20-30 -	5

(2) Continuous data — upper limit of 1st interval is equal to the lower limit of second data interval.

(3) Discontinuous data — Not equal.

8-9 +	0-9 -	3
	10-19 -	2
	20-29 -	1

(2) Ungrouped data, cannot be organised into classes or it is just a list of nos.

eg,  
= 3, 5, 11, 15, ...

## frequency distribution

(1) discrete if data is represented in such a way that exact measurements of the units are clearly shown.

(ii) Continuous & in which data are arranged in classes groups which are not exactly measurable.

## # Cumulative frequency distribution -

Sum of all the frequencies i.e., first class + second class + third class and so on is so obtained is known as cf (cumulative frequency).

~~\* two types :-~~

- ① less than
- ② greater than.

① less than :- For less than cumulative frequency, we add up the frequency from above.

② greater than :- For greater than, cumulative frequency, we add up frequencies from below.

wages	cf (no. of workers)	wages	cf (no. of workers)
less than 1099.5	125	greater than 999.5	1000
" " 1199.5	275	" " 1099.5	875
" " 1299.5	475	" " 1199.5	726
" " 1399.5	726	" " 1299.5	528
" " 1499.5	925	" " 1399.5	277
" " 1599.5	1000	" " 1499.5	100

#

## Relative frequency Distribution :-

Very useful for the comparison of two or more frequency distribution.

$$\text{Relative frequency } R_i = \left[ \frac{\text{class frequency}}{\text{Total frequency}} \times 100 \right]$$

#

### (i) Pie diagrams :-

$$\text{Central angle} = \left[ \frac{\text{frequency} \times 360^\circ}{\text{Total frequency}} \right]$$

### (ii) Histogram :-

Set of adjacent rectangles whose area is proportional to the frequencies at a given continuous frequency distribution.

no gap b/w any two successive rectangles.

(iii)

### frequency polygon :- To draw the frequency

polygon of an ungrouped frequency distribution, we plot the points with abscissae as the mid-point values and the ordinates as the corresponding frequencies.

These plotted points are joined by straight lines to obtain the frequency polygon.

(iv) Ogive (of curve) + When we plot the upper class limits along X-axis and cumulative frequencies along Y-axis, And on joining them, we get a curve called an Ogive.

\* two types :-

- (i) Less than Ogive ; the rising curve.
- (ii) More than Ogive ; falling curve.

less

(i) Less than type

the upper limits along the X-axis and the corresponding cumulative frequencies along Y-axis.

To get a rising curve.

(ii) More than type -

Subtract the frequency of each class

Now, mark the lower class limits along X-axis and

the corresponding Y-axis. (of).

To get a declining curve.

## Measures of Central Tendency.

following are the 5 measures of central tendency :-

### (1) Mathematical Averages :-

- (a) Arithmetic Mean or Mean } Also Known
- (b) Geometric Mean } as
- (c) Harmonic Mean. } Measures of Location,

### (2) Positional Averages :-

- (a) Median
- (b) Mode.

## # ARITHMETIC MEAN (AM) :-

Sum of all the numbers in the series is divided by the total no. of series .

## # ARITHMETIC MEAN OF UNGROUPED OR INDIVIDUAL OBSERVATIONS :-

If  $x_1, x_2, x_3, \dots, x_n$  are n observations,  
If a variable  $X$ , then the AM is

① Direct Method :-

$$\boxed{\bar{X} = \frac{n_1 + n_2 + n_3 + \dots + n_n}{n} \quad \text{or} \quad \bar{X} = \frac{1}{n} \left( \sum_{i=1}^n n_i \right)}$$

② Shortcut Method or (Assumed Mean Method)

$$\boxed{\bar{X} = A + \frac{1}{n} \sum_{i=1}^n d_i}$$

where,

$A$  = assumed mean

$$d_i = n_i - A$$

# WEIGHTED ARITHMETIC MEAN :-

If  $w_1, w_2, w_3, \dots, w_n$  are the weights assigned to the  $n$  values of  $n_1, n_2, \dots, n_n$ , respectively,

then the weighted average of AM is given by

$$\boxed{\bar{X} = \frac{w_1 n_1 + w_2 n_2 + \dots + w_n n_n}{w_1 + w_2 + \dots + w_n} \Rightarrow \frac{\sum_{i=1}^n w_i n_i}{\sum_{i=1}^n w_i}}$$

OR

$$\bar{x} = \frac{\sum c o n}{\sum w}$$

#

## ARITHMETIC MEAN OF A DISCRETE FREQUENCY DISTRIBUTION :-

① Direct Method :- If a variable  $X$  takes values

$x_1, x_2, x_3, \dots, x_n$  - which corresponding frequencies  $f_1, f_2, f_3, \dots, f_n$ , respectively, then the arithmetic mean of the values is.

$$\bar{x} = \frac{f_1x_1 + f_2x_2 + \dots + f_nx_n}{f_1 + f_2 + f_3 + \dots + f_n}$$

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

where,  $N = f_1 + f_2 + \dots + f_n$ .

so,

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

(ii) Shortest Method :- If the number of  $n$  or (and)  $f$  is large, the calculation of AM by the formula used above is quite tedious and time consuming.

In such a case, we take the deviation from assumed mean  $A$  which is in the middle or just close to it in the data.

Then,

$$\bar{X} = A + \frac{1}{N} \sum_{i=1}^n f_i d_i$$

where,  $d_i = n_i - A$

$A$  = Assumed Mean.

NOTE

This method is nothing but shifting of origin from zero to the Assumed Mean  $A$  ~~on the~~ on the no. line.

(iii) Step Deviation Method :-

Sometimes, during the application of shortest method of finding the mean, the deviations  $d_i$  are divided by a common number  $b$  (say).

In such a case, the

arithmetic mean will be reduced to a great extent by taking.

$$U_i = \frac{u_i - A}{b}, \quad i = 1, 2, \dots, n$$

$$\bar{x} = A + b \left( \frac{1}{N} \sum_{i=1}^n f_i u_i \right)$$

NOTE

This process is called change of scale on the no. line.

## # ARITHMETIC MEAN OF A GROUPED OR CONTINUOUS FREQUENCY DISTRIBUTION :-

We need to compute the mid-point of class intervals ( $m_i$ ).

The mid-points are multiplied by the corresponding frequencies ( $f_i$ ).

The sum of this product is obtained and is divided by the sum of frequencies.

The AM may be computed by applying any of the methods used in a discrete frequency.

#

## COMBINED ARITHMETIC MEAN :-

If we are given the AM of the two data sets and their sizes, then the combined AM of two data sets can be obtained by

$$\bar{X}_{12} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$$

where,

$\bar{X}_{12}$  = Combined mean of two data sets 1 & 2.

$\bar{x}_1$  = Mean of 1st data

$\bar{x}_2$  = Mean of 2nd data

$n_1$  = size of the 1st data

$n_2$  = " " " 2nd data.

#

If  $\bar{x}_1, \bar{x}_2, \dots, \bar{x}_K$  are the means of  $K$  series of sizes  $n_1, n_2, \dots, n_K$ , respectively then the mean  $\bar{x}$  of the composite series is given by

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + \dots + n_K \bar{x}_K}{n_1 + n_2 + \dots + n_K}$$

## # PROPERTIES OF AM :-

- (a) If  $\bar{n}$  is the mean of  $n_1, n_2, \dots, n_n$ , then mean of  $an_1+b, an_2+b, \dots, an_n+b$ , is  $a\bar{n}+b$ .
- (b) Mean is dependent of change of origin but it is independent of change of scale.
- (c) Algebraic sum of the deviations of a set of values from their AM is zero.
- (d) The sum of the squares of the deviations of a set of values is minimum, when taken about mean.

## # GEOMETRIC MEAN (GM) :-

The  $n$ th root of the product of the values is called geometric mean.

## ① # GEOMETRIC MEAN FOR UNGROUPED DATA :-

If  $n_1, n_2, \dots, n_n$  are  $n$  non-zero values of a variate  $X$ , then geometric mean is

$$GM = \sqrt[n]{(n_1 \cdot n_2 \cdot \dots \cdot n_n)}$$

$$\log GM = \frac{1}{n} (\log n_1 + \log n_2 + \dots + \log n_n)$$

$$\log GM = \frac{1}{n} \sum_{i=1}^n \log n_i$$

$$\text{Go anti log } \left( \frac{1}{n} \sum_{i=1}^n \log n_i \right)$$

## # GEOMETRICAL MEAN FOR GROUPED - DATA :-

If  $n_1, n_2, \dots, n_n$  are  $n$  observations whose corresponding frequencies are  $f_1, f_2, \dots, f_n$ , then geometric mean is given by

$$GM = \sqrt[N]{(n_1 f_1 \cdot n_2 f_2 \cdot \dots \cdot n_n f_n)}$$

$$\log GM = \frac{1}{N} (f_1 \log n_1 + f_2 \log n_2 + \dots + f_n \log n_n)$$

$$\log GM = \frac{1}{n} \sum_{i=1}^n f_i \log n_i$$

$$\text{Go anti log } \left( \frac{1}{N} \sum_{i=1}^n f_i \log n_i \right)$$

#

## COMBINED GEOMETRIC MEAN :-

If  $G_1$  and  $G_2$  are the geometric means of two series of sizes  $n_1$  and  $n_2$  respectively,

then the GM of the combined series is given by

$$\log GM = \frac{n_1 \log G_1 + n_2 \log G_2}{n_1 + n_2}$$

as median divides a distribution into <sup>two</sup> equal parts,

$$N = f_1 + f_2 + \dots + f_n$$

#

## HARMONIC MEAN :-

It is the reciprocal of the AM of the reciprocals of the observations of any series.

#

## HARMONIC MEAN FOR UNGROUPED DATA :-

If  $n_1, n_2, \dots, n_n$  are  $n$  non-zero values of a variable  $X$ , then harmonic mean is

$$HM = \frac{n}{\frac{1}{n_1} + \frac{1}{n_2} + \dots + \frac{1}{n_n}} = \frac{n}{\sum_{i=1}^n \left(\frac{1}{n_i}\right)}$$

## # HARMONIC MEAN FOR GROUPED DATA :-

If  $n_1, n_2, \dots, n_n$  are  $n$  observations, whose corresponding frequencies of each variate is  $f_1, f_2, \dots, f_n$ , then

$$HM = \frac{\frac{f_1 + f_2 + \dots + f_n}{f_1 + f_2 + \dots + f_n}}{\frac{n_1}{f_1} + \frac{n_2}{f_2} + \dots + \frac{n_n}{f_n}} = \frac{N}{\sum_{i=1}^n \left(\frac{n_i}{f_i}\right)}$$

$$\frac{\sum_{i=1}^n f_i}{\sum_{i=1}^n \frac{f_i}{n_i}}$$

where,  $N = f_1 + f_2 + \dots + f_n$

NOTE:- If AM is to be find, and if GM and HM of a statistical sequence is known, then AM can be obtained as

$$AM = \frac{(GM)^2}{(HM)}$$

or

$$(GM)^2 = (AM) \times (HM)$$

#

### MEDIAN :-

It is the middle most or the central values of the variable in a set of observations when the observations are arranged either in ascending or in descending order of their magnitudes.

It divides the arranged series in two equal parts.

#

### MEDIAN OF AN UNGROUPED INDIVIDUAL OBSERVATIONS :-

Let  $n_1, n_2, n_3, \dots, n_n$  be set of observations.

- (a) arrange the observations in ascending or descending order.

If  $n$  is odd :-

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

If  $n$  is even :-

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

2

✓ #

## MEDIAN OF THE DISCRETE FREQUENCY DISTRIBUTION :-

In case of a discrete frequency distribution,

$x_i, f_i, i \geq 1, 2, \dots, n$ ,

- (a) Arrange the data in ascending or descending order and then find the cumulative frequency ( $f$ ).

(b)

find  $\left(\frac{N}{2}\right)$ , where  $N = \sum_{i=1}^n f_i$

① See the (cf) just greater than  $\frac{N}{2}$ .

The corresponding value of  $x_c$  is median.

~~ANOTHER :-~~

Find the cumulative frequency (CF)

~~If  $N$  is odd :-~~

Median ( $M$ ) = Value of  $\left(\frac{N+1}{2}\right)$ th observation

~~If  $N$  is even :-~~

Median ( $M$ ) = Value of

$\left(\frac{N}{2}\right)$ th observation +  $\left(\frac{N}{2} + 1\right)$ th observation

where,

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

## ~~#~~ MEDIAN OF CONTINUOUS FREQUENCY DISTRIBUTION :-

Let the Number of observation be  $N$ .

(a) prepare the cumulative frequency (cf) column and obtain  $N = \sum f_i$

and find  $\frac{N}{2}$ .

(b) Find the median class i.e., the class in which the observation whose cf. is equal to or just greater than  $\frac{N}{2}$  lies.

This class is known as the median class.

(c) Use the formula,

(M) Median =

$$M = l + \left[ \frac{\frac{N}{2} - c}{f} \right] \times h$$

where,  
 $l$  = lower limit of the median class  
 $f$  = frequency of the median class  
 $h$  = width (size) of the " "  
 $c$  = cf of the class preceding the median class.

#

## MODE :-

It is that value in a series which occurs most frequently. It has the maximum frequency.

It is the observation with maximum frequency, whenever the other observations have less frequencies.

#

### MODE OF INDIVIDUAL OBSERVATIONS :-

The value which is repeated more no. of times is called mode of the series.

#

### MODE OF A DISCRETE SERIES :-

It is the value of variable consisting highest frequency.

#

### MODE OF A CONTINUOUS SERIES :-

- ① first find the modal class i.e., which has the maximum frequency,

⑤ Use the formula,

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

where,

$l$  = lower limit of modal class

$h$  = width / size of the modal class

$f_1$  = frequency of .. .. ..

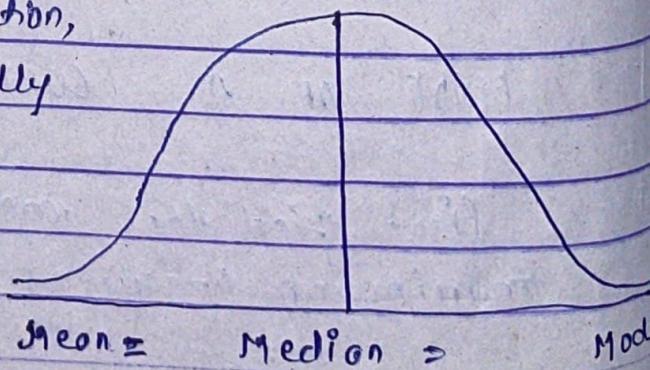
$f_0$  = " " " class preceding  
modal class

$f_2$  = frequency of the class succeeding  
modal class.

## # Symmetric Distribution :-

A distribution is a symmetric distribution,  
if the values of mean, median, mode coincide.

In a symmetric distribution,  
frequencies are symmetrically  
distributed on both sides of  
the centre point of the  
frequency curve.



## # Relationship b/w Mean, Median and Mode :-

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

$$\Rightarrow \boxed{\text{Mode} = 3\text{Median} - 2\text{Mean}}$$

## Measures of Dispersion

Dispersion is the measure of the variations.

The degree to which numerical data tend to spread about an average value is dispersion of the data.

The measures of dispersion commonly used are -

- (a) Range
- (b) Mean Deviation
- (c) Standard Deviation

## # RANGE :-

diff. b/w max. and min. values of variate

i.e.,  $\boxed{\text{Range} = L - S}$  where,

$L$  = largest value

$S$  = smallest value

NOTE

$$\text{Coefficient of Range} = \frac{L-S}{L+S}$$

#

### MEAN DEVIATION :-

Mean deviation about a central value is defined as the AM of the absolute deviations of all the values taken about that central value.

#

### M.D. OF INDIVIDUAL OBSERVATIONS :-

If  $x_1, x_2, \dots, x_n$  are  $n$  values of a variable  $x$ , then the mean deviation from an average  $A$  is given by -

$$\boxed{M.D.(A) = \frac{1}{n} \sum_{i=1}^n |x_i - A| = \frac{1}{n} \sum |d_i|}$$

where,  $d_i = x_i - A$

## # M.D. OF DISCRETE FREQUENCY DISTRIBUTION :-

If  $x_1, x_2, \dots, x_n$  are  $n$  observations with frequencies  $f_1, f_2, f_3, \dots, f_n$ , then mean deviation from an average  $A$  is given by

$$\text{M.D.}(A) = \frac{1}{N} \sum f_i |x_i - A|$$

where,

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

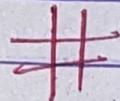
## # M.D. OF CONTINUOUS FREQUENCY DISTRIBUTION :-

The process is same as for a discrete frequency distribution. The only diff. is that here we have to obtain the midpoints of the continuous classes and take the deviations of these mid-points from the given average  $A$ .

NOTE:-

$$\text{Coefficient of M.D.} = \frac{\text{M.D.}}{\text{AM}} = \frac{\text{M.D.}}{\bar{x}}$$

**Note**  
 $\bar{x}$  = Mean deviation of a given set of observations  
 i.e. Least deviation about their median.



## VARIANCE AND STANDARD DEVIATION :-

The variance of a variable  $n$  is the AM of the squares of all deviations of  $n$  from the AM of the observations and is denoted by

$$\text{Var}(x) \text{ or } \sigma^2.$$

The positive sq. root of the variance of a variable  $n$  is known as standard deviation.

i.e.,

$$\text{Standard Deviation (S.D.)} = \sqrt{\text{Var}(n)}$$

$$= \sqrt{\sigma^2} = \sigma$$



## VARIANCE OF INDIVIDUAL OBSERVATIONS :-

If  $x_1, x_2, \dots, x_n$  are  $n$  values of a variable  $x$ ,

then

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

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

or

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

If the values of variable  $x$  are large, calculation of variance from the above formulae is quite tedious and time consuming.

In that case, we take deviation from an arbitrary point  $A$  (say), then

$$Var(x) = \frac{1}{n} \left[ \sum_{i=1}^n (d_i - \bar{d})^2 \right] = \frac{1}{n} \sum_{i=1}^n d_i^2 - \left( \frac{1}{n} \sum_{i=1}^n d_i \right)^2$$

where,  $d_i = x_i - A$

and  $\rightarrow$  let  $A$  be any Assumed Mean. and  $d_i = x_i - A$

$$S.D. = \sigma =$$

$$\sqrt{\frac{1}{n} \sum_{i=1}^n d_i^2}$$

$$S.D. = \sigma =$$

$$\sqrt{\frac{\sum_{i=1}^n d_i^2}{n} - \left( \frac{\sum_{i=1}^n d_i}{n} \right)^2}$$

Also,

$$\text{Var}(n) = \frac{h^2}{n} \left[ \sum_{i=1}^n (u_i - \bar{u})^2 \right]$$

$$\text{Var}(n) = h^2 \left[ \frac{1}{n} \sum_{i=1}^n u_i^2 - \left( \frac{1}{n} \sum_{i=1}^n u_i \right)^2 \right]$$

where,

$$u_i = \frac{n_i - A}{h}$$

#

## VARIANCE OF DISCRETE FREQUENCY

DISTRIBUTION:-

If  $n_1, n_2, \dots, n_n$  are  $n$  observations with frequencies  $f_1, f_2, \dots, f_n$ ; then

$$\text{Var}(n) = \frac{1}{N} \left\{ \sum_{i=1}^n f_i (n_i - \bar{n})^2 \right\} = \left( \frac{1}{N} \sum_{i=1}^n f_i n_i^2 \right) - \bar{n}^2$$

where,

$$\text{No. } \sum_{i=1}^n f_i$$

If the value of  $n$  or  $f$  are large, we take the deviations of the values of variable  $n$ , from an arbitrary point  $A$  (say).

$$\therefore d_i = n_i - A \quad ; \quad i = 1, 2, 3, \dots, N$$

$$\text{S.D. } (\sigma) = \sqrt{\left[ \frac{1}{N} \sum_{i=1}^n f_i (n_i - \bar{n})^2 \right]}$$

Variance  $\sigma^2$

$$\therefore \text{Var}(n) = \frac{1}{N} \left[ \sum_{i=1}^n f_i (d_i - \bar{d})^2 \right]$$

$$\text{Var}(n) = \frac{1}{N} \left( \sum_{i=1}^n f_i d_i^2 \right) - \left( \frac{1}{N} \sum_{i=1}^n f_i d_i \right)^2$$

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

$$S.D. (\sigma) = \sqrt{\frac{\sum_{i=1}^n f_i d_i^2}{N} - \left( \frac{\sum_{i=1}^n f_i d_i}{N} \right)^2}$$

where,

$$N = \sum_{i=1}^n f_i, \quad \text{and} \quad \bar{x} = \frac{\sum_{i=1}^n f_i d_i}{N}$$

Sometimes,  $d_i = n_i - A$  are divisible by a common number (h) say,

then

$$u_i = \frac{n_i - A}{h}, \quad i = 1, 2, 3, \dots, N$$

then

$$Var(n) = \frac{h^2}{N} \left[ \sum_{i=1}^n f_i (u_i - \bar{u})^2 \right]$$

$$\Rightarrow Var(n) = h^2 \left[ \frac{1}{N} \sum_{i=1}^n f_i u_i^2 - \left( \frac{1}{N} \sum_{i=1}^n f_i u_i \right)^2 \right]$$

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

## # VARIANCE OF A GROUPED OR CONTINUOUS FREQUENCY DISTRIBUTION :-

Some formulae, as used in discrete frequency distribution can be used.

$$\text{Variance} = \frac{\sum_{i=1}^n f_i d_i^2}{N} - \left( \frac{\sum_{i=1}^n f_i d_i}{N} \right)^2$$

where,

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

where,

$$N = \sum f_i$$

$$S.D. = \sigma$$

$$\sqrt{\frac{\sum_{i=1}^n f_i d_i^2}{N} - \left( \frac{\sum_{i=1}^n f_i d_i}{N} \right)^2}$$

NOTE:

Variance is independent of change of origin but dependent on change of scale. Adding or subtracting a positive number from each observation of a

group does not affect the variance. If each observation is multiplied by a constant  $b$ , then variance of the resulting group becomes  $b^2$  times the original variance.

- While calculating S.D., the deviations are to be taken about arithmetic mean only.

$$\text{Coefficient of dispersion} = \frac{\text{S.D.}}{\text{Mean}} = \frac{\sigma}{\bar{x}}$$

## # COEFFICIENT OF VARIATION :-

The M.D. and S.D. have the same units in which the data is given.

$$\text{C.V.} = \frac{\text{S.D.}}{\text{Mean}} \times 100 = \frac{\sigma}{\bar{x}} \times 100$$

The series having greater C.V. is said to be more variable and less consistent than the other.

## # STANDARD DEVIATION OF COMBINED GROUP OR SERIES $\sigma$ -

Let  $\bar{x}_1, \bar{x}_2$  be the respective means (A.M) and  $\sigma_1, \sigma_2$  be the respective S.D of the two given series having no. of observations as  $n_1$  and  $n_2$ , respectively, then

Combined standard deviation  $\sigma_{12}$  of all the observations taken together is given by

$$\boxed{\text{Combined S.D} \quad \sigma_{12} = \sqrt{\frac{n_1(\sigma_1^2 + D_1^2) + n_2(\sigma_2^2 + D_2^2)}{n_1 + n_2}}}$$

$$\boxed{\text{Variance } (\sigma^2) = \frac{1}{n_1 + n_2} [n_1(\sigma_1^2 + D_1^2) + n_2(\sigma_2^2 + D_2^2)]}$$

where,

$$D_1 = \bar{x}_1 - \bar{x}$$

$$D_2 = \bar{x}_2 - \bar{x}$$

and

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

#

## CORRELATION :-

The tendency of simultaneous variation b/w two variable is called correlation or covariation.

It denotes the degree of inter-dependence b/w variables.

#

## CORRELATION COEFFICIENT :-

The no. showing the degree or extent to which x are related to each other is called correlation coefficient.

It is denoted by  $r_{xy}$  or  $r$  or  $P(x,y)$ .

$r_{xy}$  or simply  $r$ .

#

## METHODS OF CALCULATING CORRELATION COEFFICIENT

(1)

### KARL PEARSON'S COEFFICIENT OF CORRELATION :-

$$\text{Covariance } (r_{xy}) = \text{cov}(x,y)$$

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

let  $\sigma_x$  and  $\sigma_y$  be the SD of variables  $x$  and  $y$ , respectively.

Then, coefficient of correlation

$$r(x,y) = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} \Rightarrow \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

$$\Rightarrow \frac{n \sum x_i y_i - (\sum x_i)(\sum y_i)}{\sqrt{n \sum x_i^2 - (\sum x_i)^2} \sqrt{n \sum y_i^2 - (\sum y_i)^2}}$$

## (2) RANK CORRELATION (SPEARMAN'S)

Let  $d$  be the difference b/w paired ranks and  $n$  be the no. of items ranked.

Then,  $r$  the coefficient of rank correlation is given by

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

NOTE

The rank correlation lies b/w -1 and 1.

## # PROPERTIES OF CORRELATION

- (i)  $-1 \leq r_{(n,y)} \leq 1$
- (ii) If  $r=1$ , then the coefficient of correlation is perfectly positive.
- (iii) If  $r=-1$ , the correlation is perfectly negative.
- (iv) The correlation coefficient is a pure no. independent of the unit of measurement.
- (v) The coefficient of correlation is independent of the change in origin and scale.
- (vi) If  $-1 < r < 1$ , it indicates the degree of linear relationship b/w  $x$  and  $y$ , whereas its sign tells about the direction of relationship.
- (vii) If  $x$  and  $y$  are two independent variables,  $r=0$
- (viii) If  $r=0$ ,  $x$  and  $y$  are said to be uncorrelated. It

does not imply that the two variables were independent.

$$r(x,y) > 0$$

(iv) If  $x$  and  $y$  are random variables and  $a, b, c$  and  $d$  are any nos. such that  $a \neq 0, c \neq 0$ , then

$$r(ax+b, cy+d) = \frac{ac}{\sqrt{a^2 + b^2} \sqrt{c^2 + d^2}} r(x, y).$$

## # LINES OF REGRESSION :-

$y_t$  is the line which gives the best estimate to the value of one variable for any specific value of the other variable.

Therefore, the line of regression is the line of best fit and is obtained by the principle of least squares.

## # REGRESSION ANALYSIS :-

① The line of regression of  $y$  on  $x$  or regression line of  $y$  on  $x$  is given by

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

(ii) The line of regression of  $x$  on  $y$  or regression line of  $x$  on  $y$  is given by

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

(iii) Regression coefficient of  $x$  on  $y$ , is denoted by  $b_{xy}$ ,

$$b_{xy} = r \frac{\sigma_x}{\sigma_y} \cdot \frac{\text{cov}(x, y)}{\sigma_x^2}$$

(iv) Regression coefficient of  $y$  on  $x$ , is denoted by  $b_{yx}$ ,

$$b_{yx} = r \frac{\sigma_x}{\sigma_y} \cdot \frac{\text{cov}(x, y)}{\sigma_y^2}$$

(v) If  $\theta$  is the angle between the two regression lines, then

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

$$\text{where, } \tan \theta = \frac{M_2 - M_1}{1 + M_1 M_2}$$

- (a) If  $r = 0, \theta_1 = \frac{\pi}{2}$ , then the two regression lines are perpendicular to each other.
- (b) If  $r = 1, \theta_1 = 0, \theta_2 = \pi$ , then the regression lines coincide.

### # PROPERTIES OF THE REGRESSION COEFFICIENTS :-

- (i) Both regression coefficients and correlation coefficient  $r$  have the same sign.
- (ii) Coefficient of correlation is the geometric mean between the regression coefficients.
- (iii) If one of the regression coefficient is greater than unity, the other must be less than unity,

$$0 < |b_{xy} \cdot b_{yx}| < 1, \text{ if } r \neq 0$$

i.e., if  $|b_{xy}| > 1$ ,

$$|b_{yx}| < 1.$$

- (iv) Regression coefficients are independent of the change of origin but not of scale.
- (v) Arithmetic Mean of the regression coefficient is greater than the correlation coefficient.
- (vi) The two lines of regression cut each other at the point  $(\bar{x}, \bar{y})$ .  
Thus, on solving the two lines of regression, we get the values of means of variables in the bivariate distribution.