

UNIT - IV

MEASURES OF CENTRAL TENDENCY

DEFINITION

Average is a value which is typical or representation of a set of data given by MURRY. R. SPEIGAL.

Average is an attempt to find one single figure to describe the whole of figures. - CLARK & SEKKODS

Average is sometime described as number which is typical of the whole group - LEABO.

FUNCTIONS OF AVERAGES

- (1) To facilitate quick understanding of complex Data.
- (2) To facilitate comparison
- (3) To know about the universe from a sample.
- (4) To help in Decision making.
- (5) To estimate or establish mathematical Relationship.

CHARACTERISTICS OF A GOOD AVERAGES

- * It should rigidly defined (well defined)
- * It should be easy to understand.
- * It should be simple to complete.
- * Its definition should be in the form of a mathematical formula.
- * It should be based on all items in the form of data.
- * It should not be influenced by any single item (or)

Group of item.

- * It should be capable of further algebraic treatment.
- * It should be capable of being used in further statistical computation.
- * It should have sampling stability.

TYPES OF AVERAGES

The following are four types of averages:

- Arithmetic mean
 - Simple
 - Weighted
- Median
- Mode
- Geometric mean
- Harmonic mean

(i) ARITHMETIC MEAN

Average average is also called as 'Mean'. Arithmetic average of a series is the figure obtained dividing the total value of the various items by their number. There are two types:

a) SIMPLE : * INDIVIDUAL SERIES

- Add up all the values of the variables (x) and find Σx .
- Divide Σx by their no. of observation (n)

$$\therefore \bar{x} = \frac{\Sigma x}{n}$$

b) WEIGHTED : INDIVIDUAL SERIES

The weighted is a number standing for the relative importance of the items. Weighted average can be defined as an average where components items are multiplied by their respective weights and the aggregate of the product are divided by total of arithmetic average.

$$\therefore \bar{X}_W = \frac{\sum WX}{\sum W}$$

ARITHMETIC DISCRETE AVERAGE:

- Multiply each size of item (variable by its frequency) that is F_x .
- Add all the F_x to $\sum f_x$ divided $\sum f_x$ by total frequency

$$\therefore \bar{X} = \frac{\sum f_x}{N}$$

ARITHMETIC CONTINUOUS AVERAGE:

- The midpoint value must be find.
- Assume any one of midpoint as a assumed mean.
- Find out the deviation of mid value of each class.
- Deviation are divided by a common pattern to obtained 'd'
- Multiply d of each class by its frequency to get fd .
- Add up the product value to get summation fd .

$$\therefore \bar{X} = A + \left(\frac{\sum fd}{N} \right) \times c$$

A - ASSUMED MEAN

N = $\sum f$

d - $(M-A)/c$

c - WIDTH OF THE CLASS-INTERVAL OF A COMMON FACTOR.

MEDIAN

The median is that value of the variable which divides the group.

INDIVIDUAL SERIES:

Median refers to the middle value in the distribution.

If the total no. of items is an odd figure, the median is the middle value. If the total is even figure, the average of items in the center of the distribution.

DISCRETE SERIES:

- Arrange the data in Ascending order.
- Find the cumulative frequency
- Apply the formula

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

$$(N = \sum f)$$

CONTINUOUS SERIES:

- Find out the Median by using $\left(\frac{N}{2}\right)$.
- Find out the class in which Median lies.
- The formula,

$$\text{Median} = L + \frac{\left(\frac{N}{2}\right) - cf}{f} \times C$$

L - lower limit

f - frequency

c.f - cumulative frequency

C - class interval

PARTITIONED VALUES:

Median divides the distribution into two equal parts. There are other values also which divide the series into equal parts and they are called the partition values.

ONE POINT divides the series (HALVES)

THREE POINT divides the series (QUARTILES)

NINE POINT divides the series (DECILES)

99 POINT divides the series (PERCENTILES)

CALCULATION OF QUANTILES (INDIVIDUAL & DISCRETE SERIES)

- Find out the cumulative frequency.
- Then apply the formula.

FIRST QUANTILE (Q_1) = size of $\left(\frac{N+1}{4}\right)^{\text{th}}$ item

THIRD QUANTILE (Q_3) = size of $\left(\frac{3(N+1)}{4}\right)^{\text{th}}$ item

FOR CONTINUOUS SERIES

$$Q_1 = L + \left(\frac{N/4 - c.f}{f} \right) \times c$$

$$Q_3 = L + \left(\frac{\frac{3N}{4} - c.f}{f} \right) \times c$$

DECILES:

$$D_2 = L + \left(\frac{\frac{2N}{10} - c.f}{f} \right) \times c$$

PERCENTILES:

$$P_{25} = L + \left(\frac{\frac{25N}{100} - c.f}{f} \right) \times c$$

MODE

The mode of the distribution is the value at the point around which the items tend to most nearly concentrated.

INDIVIDUAL :

Mode can often be found out by which occur more no. of times than the other.

DISCRETE :

Mode is the variable value which corresponds to the maximum frequency.

CONTINUOUS :

First we have to find Modal class and then value of mode by formula,

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

ASYMPTOTIC MODE :

$$\text{MODE} = 3 \text{ MEDIAN} - 2 \text{ MEAN}$$

$$\text{Symmetrical} = \bar{x} = \text{Median} = \text{Mode}$$

$$\text{Asymmetrical} = \bar{x} \neq \text{Median} = \text{Mode}$$

GEOMETRIC MEAN

It is defined as the n^{th} root of the product of n times.

$$\therefore \text{G.M.} = \sqrt[n]{x_1 x_2 \dots x_n}$$

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

DISCRETE :

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

CONTINUOUS :

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

HARMONIC MEAN

Harmonic mean is the reciprocal of the arithmetic average of values of various items in the variable.

INDIVIDUAL :

$$H.M. = \frac{N}{\sum \frac{1}{x}}$$

DISCRETE :

$$H.M. = \frac{N}{\sum f \left(\frac{1}{x} \right)}$$

CONTINUOUS :

$$H.M. = \frac{N}{\sum f \frac{1}{M}}$$

MEASURES OF DISPERSION

DEFINITION

Dispersion is the measure of variation of the items.

- Bowley

Dispersion is a measure of extent to which the individual items - CONNER.

INTER-QUARTILE RANGE AND QUARTILE DEVIATION:

The distance between first & third quartile is called Inter-quartile range.

$$\text{Inter Quartile range} = Q_3 - Q_1$$

Semi Inter Quartile range (or) Quartile deviation

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

$$\text{co-efficient of Quartile deviation} = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

MEAN DEVIATIONS (OR) AVERAGE DEVIATIONS:

Average deviations is the average amount of scatter of the items in a distribution from either the mean or the median, ignoring the sign of the deviation. The relative measure of deviation of mean is called co-efficient of mean deviation.

$$\text{co-efficient} = \frac{\text{Mean deviation}}{\text{Mean (or) Median}}$$

INDIVIDUAL:

$$M.D. = \frac{\sum |D|}{n}$$

DISCRETE:

$$M.D. = \frac{\sum f |D|}{N}$$

CONTINUOUS:

$$M.D. = \frac{\sum f |D|}{N}$$

STANDARD DEVIATION:

Standard deviation is also called ROOT MEAN SQUARE DEVIATION OR mean error OR Mean square error. It is denoted by greek letter σ (sigma).

INDIVIDUAL:

* Deviation taken from actual mean and assumed mean.

by ACTUAL MEAN:

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

$$\bullet \sigma = \sqrt{\frac{\sum x^2}{N}} \quad (\text{or}) \quad \sqrt{\frac{\sum (x - \bar{x})^2}{N}}$$

by ASSUMED MEAN:

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

Step deviation Method:

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

COMBINED STANDARD DEVIATION

$$\bar{x}_{12} = \frac{N_1 \bar{x}_1 + N_2 \bar{x}_2}{N_1 + N_2}$$

$$\sigma_{12} = \sqrt{\frac{N_1 \sigma_1^2 + N_2 \sigma_2^2 + N_1 d_1^2 + N_2 d_2^2}{N_1 + N_2}}$$

$$\sigma_{12} = \sqrt{\frac{N_1 (\sigma_1^2 + d_1^2) + N_2 (\sigma_2^2 + d_2^2)}{N_1 + N_2}}$$

For, 'n' natural numbers, the standard deviation can be computed by the formula,

$$\sigma = \sqrt{\frac{1}{12} (n^2 - 1)}$$

VARIANCE : square of standard deviation is called variance

$$\sigma = \sqrt{\text{Variance}}$$

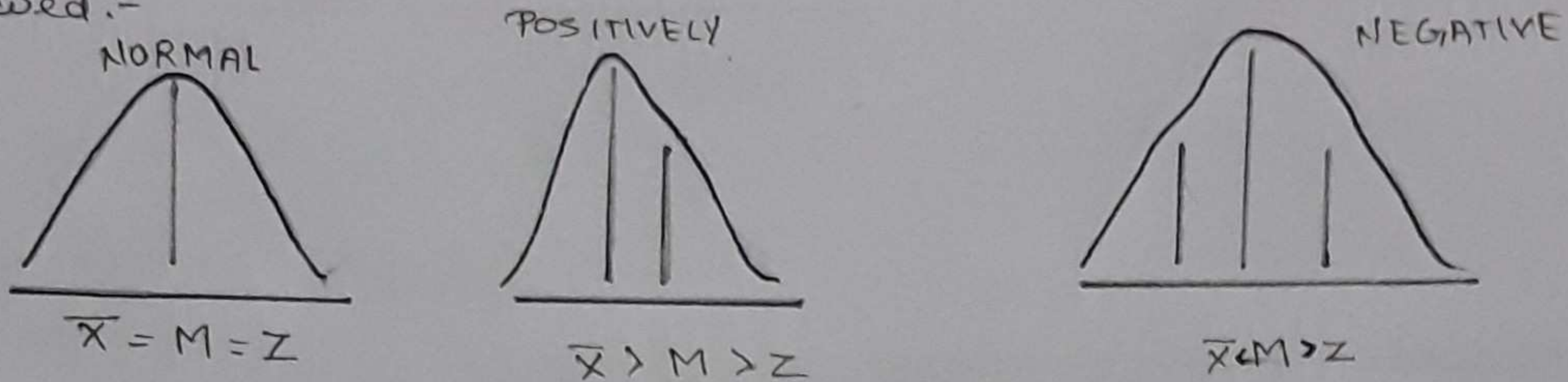
$$\text{Co-efficient of Variation (C.V.)} = \frac{\sigma}{\bar{x}} \times 100.$$

SKEWNESS AND KURTOSIS

DEFINITION:

Skewness or asymmetry is the attribute of a frequency distribution that extends further on one side of the class with the highest frequency than on the other - SIMPSON & KALFA.

When a series is not symmetrical, it is said to be asymmetrical or skewed:-



CHARACTERISTICS

- * $\bar{x} \neq \text{Median} \neq Z$
- * $Q_3 - \text{Median} \neq \text{Median} - Q_1$
- * sum of +ve deviation \neq sum of -ve deviation.
- * The frequencies are unequal on all sides.
- * The plotted graph have unequal halves.

CO-EFFICIENT OF SKEWNESS:

1. KARL PEARSON $SK_p = \frac{\bar{x} - \text{Mode}}{\sigma}$

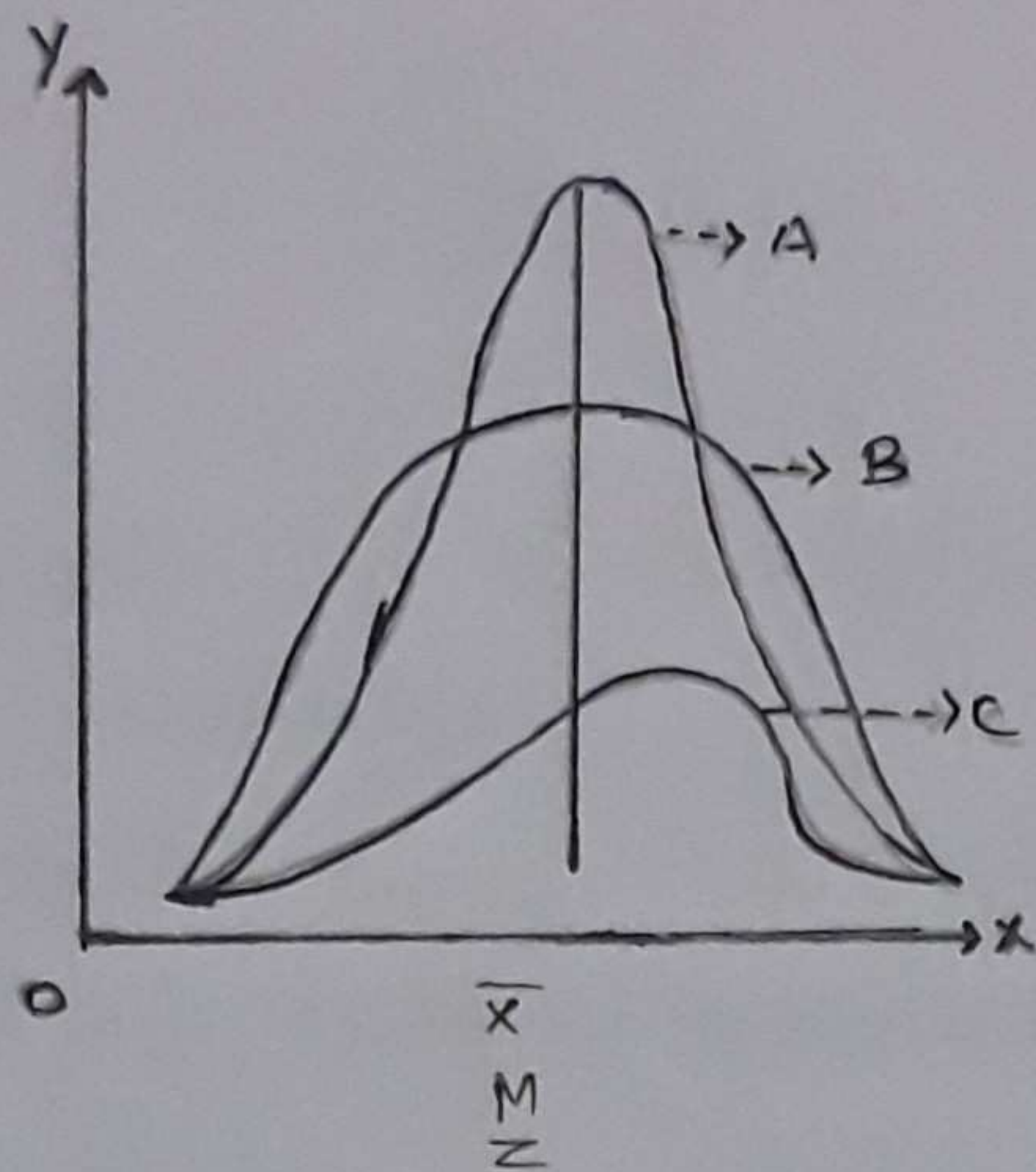
$$SK_p = \frac{3(\text{Mean} - \text{Median})}{\sigma}$$

2. BOWLEY $SK = \frac{Q_3 + Q_1 - 2\text{Median}}{Q_3 - Q_1}$

KURTOSIS

DEFINITION:

The degree of kurtosis of a distribution is measured relative to the peakedness of a normal - SIMPSON & KALFA



- A - PEAKED (LEPTO)
- B - NORMAL (MESO)
- C - FLATTOPPED (PLATY)

MEASURES OF KURTOSIS:

It is based on the fourth moments about the mean of distribution.

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

- If $\beta_2 = 3$, the distribution is normal
- $\beta_2 > 3$, the curve is leptokurtic
- $\beta_2 < 3$, the curve is platykurtic

MOMENTS

DEFINITION:

Moment is a familiar term for the measure of a force with reference to its tendency to produce rotation. The strength of this tendency depends, obviously, upon the amount of the force and the distance from the origin of the point at which the force is exerted. — FEDERIC MILLS.

UNIT - IV
CORRELATION

DEFINITION:

Correlation is an analysis of co-variation between two or more variables.

Correlation Analysis attempts to determine the degree of relationship between variables. — YA KUN CHOU

TYPES OF CORRELATION

There are 3 important types of correlation.

- (i) Positive and Negative
- (ii) Simple, Partial & Multiple
- (iii) Linear & Non-linear

POSITIVE & NEGATIVE:

This type of correlation depends upon the direction of change of variables.

If two variables tend to move together in the same direction, there is an increase in those two variables or a decrease in those variables is called positive correlation.

If two variables tend to move in opposite directions together, there is an increase in one variable and a decrease in another and vice versa is called negative.

SIMPLE, PARTIAL, MULTIPLE:

This type depends on the number of variables.

If only two variables, it is called simple.

eg: Demand & price.

In multiple, we study more than two variables simultaneously.

Eg: Price, Demand and commodity.

If we study two variables excluding the effect of other variables is called partial.

eg: The study of Price and demand eliminating the supply.

LINEAR OR NON-LINEAR

If the ratio of change between two variables is uniform then it will be linear. If we plot the graph, it will be the straight line.

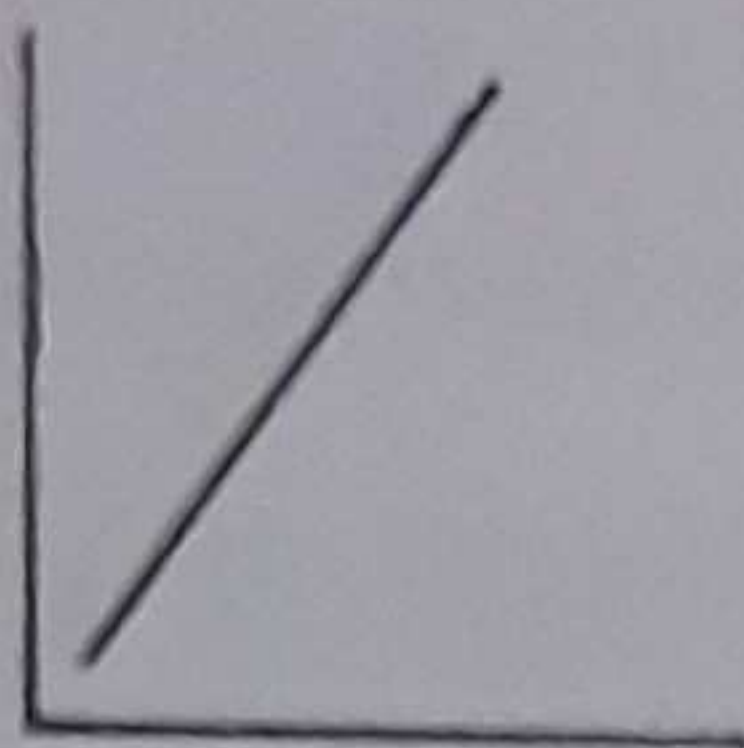
In Non-linear, we get a curve.

METHODS OF MEASURING CORRELATION:

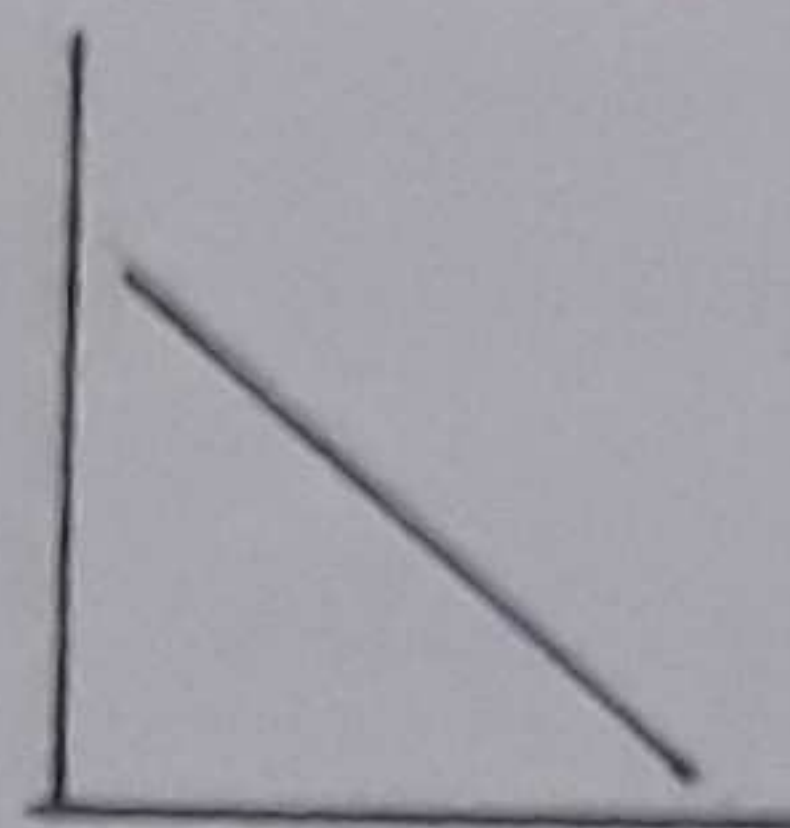
- * scattered diagram
- * Karl Pearson co-efficient of correlation
- * Spearman's correlation
- * RANK
- * Co-efficient of concurrent Deviation

SCATTERED DIAGRAM:

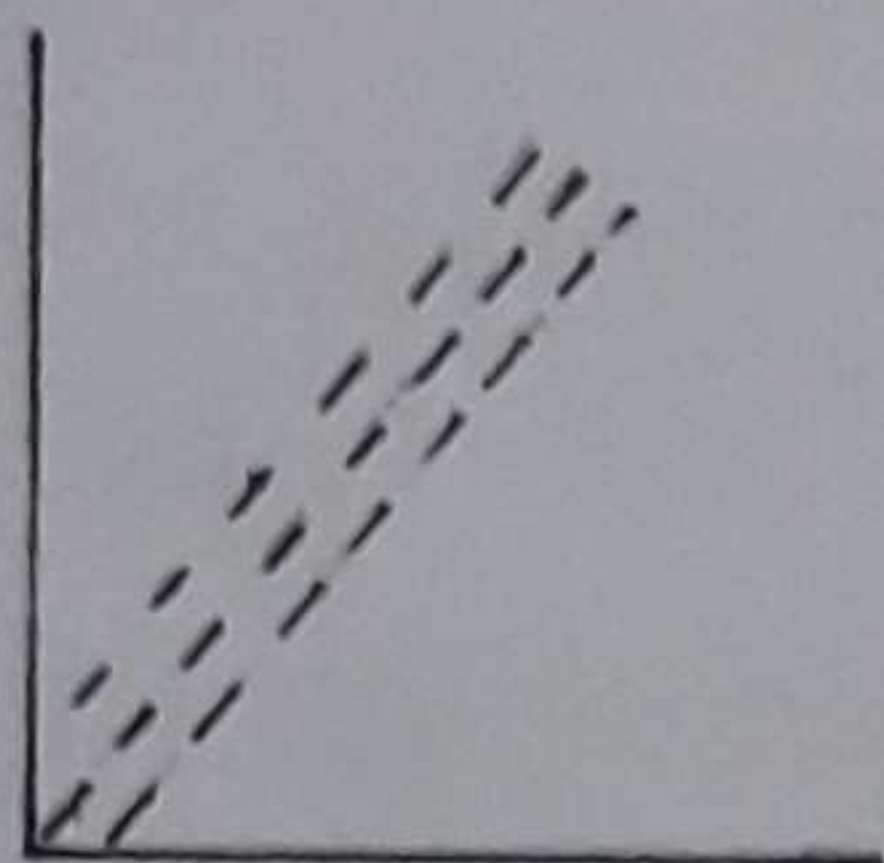
This is the simplest method of finding out whether there is any relationship present between 2 variables by plotting the values on a chart, known as scatter diagram. In this method, the given data are plotted on a graphed paper in the form of dots. This will show the types of correlation.



Perfect Positive correlation
(ie, $r = +1$ diagram)



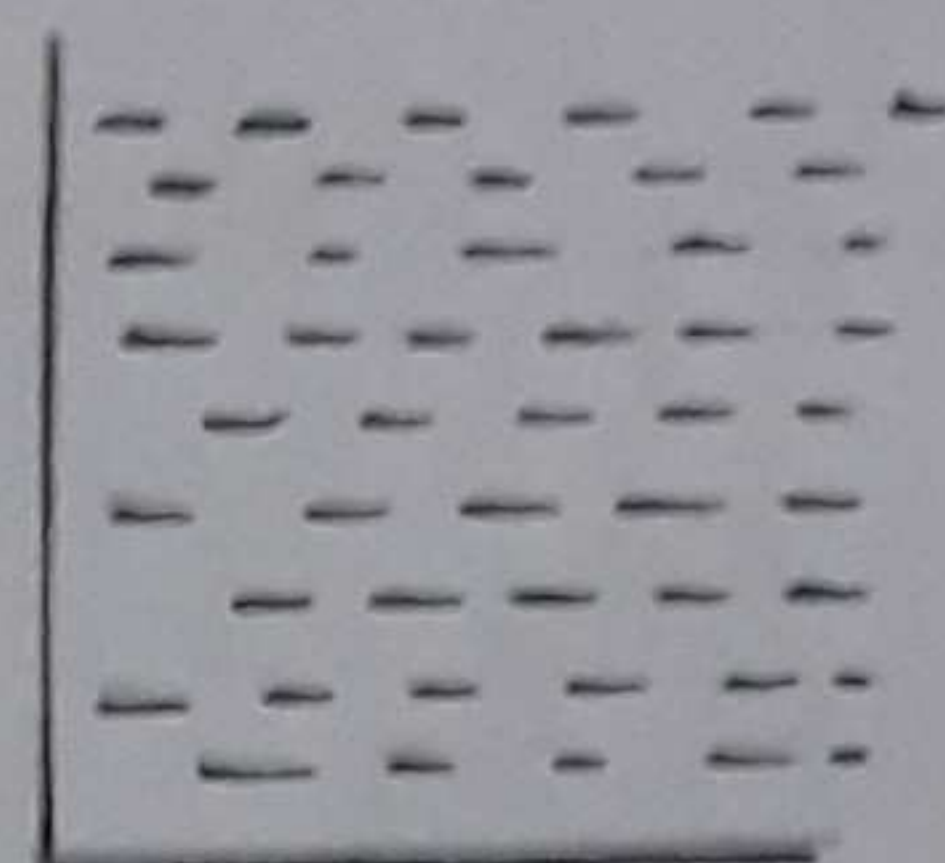
Perfect Negative Correlation
(ie, $r = -1$ diagram)



High degrees of +ve diagram correlation



High degrees of -ve correlation



No correlation

KARL PEARSON'S CO-EFFICIENT OF CORRELATION:

This method is the most widely used method in Practice and is known as Pearsonian co-efficient of correlation

It is denoted by the symbol ' r ', the formula,

$$(1) r = \frac{\text{Covariance of } xy}{\sigma_x \times \sigma_y}$$

$$(2) r = \frac{\sum xy}{N \sigma_x \times \sigma_y}$$

$$(3) r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

σ_x = standard deviation of series x .

σ_y = standard deviation of series y .

The value of co-efficient will always lie between -1 and $+1$.

$r = +1$ \rightarrow Perfect +ve correlation

$r = -1$ \rightarrow Perfect -ve correlation

$r = 0$ \rightarrow No correlation.

RANK CORRELATION CO-EFFICIENT:

This method is based on rank. This measure is useful in dealing with qualitative characteristic.

RANK ARE NOT GIVE:

When no rank is given but actual data are given then we must assign ranks. We can give ranks by the highest value has rank 1 and next value has rank 2 for both variable separately.

$$\text{The correlation coefficient} = 1 - \frac{6 \sum D^2}{n^3 - n}$$

TIED RANKS:

When two or more items have equal values it is difficult to give ranks to the in that case, the items are given the average of the ranks. they could received if they are not tied.

REGRESSION

DEFINITION:

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

REGRESSION EQUATION:

There are two regression equation namely,

(i) R.E. of X on Y

(ii) R.E. of Y on X

There are two methods of finding regression co-efficient

DEVIATION TAKEN FROM ASSUMED MEAN:

When the actual is not a whole number then deviation are taken from the assumed mean and correlation is found out by using,

$$r = \frac{n \sum dx dy - (\sum dx)(\sum dy)}{\sqrt{n \sum d^2 - (\sum dx)^2} \sqrt{n \sum d^2 - (\sum dy)^2}}$$

where,

$$dx = (x - A) \quad ; \quad dy = (y - B) \quad ; \quad A, B \text{ are assumed means of } x \text{ \& } y.$$

SPEARSMAN RANK CORRELATION:

- Take the deviation of x series and find $\sum dx$
- Similarly, take the deviation of y series and find $\sum dy$
- Square dx and dy and get total $\sum dx^2$ and $\sum dy^2$ respectively.
- Multiply dx and dy and get $\sum dx \cdot dy$.
- Then the formula,

$$= \frac{\sum dx dy - \frac{\sum dx \sum dy}{N}}{\sqrt{\left[\sum dx^2 - \frac{(\sum dx)^2}{N} \right]} \sqrt{\left[\sum dy^2 - \frac{(\sum dy)^2}{N} \right]}}$$

MATHEMATICAL PROPERTIES:

1) The co-efficient lies between -1 & $+1$.

$$-1 \leq r \leq +1$$

Moreover, this provides check on our calculations.

2) co-efficient is independent of change of origin and scale.

In change in scale, all the variables are multiplied or divided by same constant.

namely deviations from original mean of x and y and deviations from assumed mean:

DEVIATION FROM ORIGINAL MEAN:

Regression of x on y

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

where $b_{xy} = r \frac{\sigma_x}{\sigma_y} = \frac{\sum xy}{\sum y^2}$

DEVIATION FROM ASSUMED MEAN:

If the actual mean is in fraction, this method can be used.

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

We can find out the value of $r \frac{\sigma_x}{\sigma_y}$ (by applying formula)

$$b_{xy} = \frac{n \sum dxdy - (\sum dx)(\sum dy)}{n \sum dy^2 - \sum dy^2}$$

MATHEMATICAL PROPERTIES:

(i) The Geometric means between two regression co-efficient is the co-efficient of correlation.

$$r = \pm \sqrt{b_{xy} \times b_{yx}}$$

(ii) Arithmetic mean of b_{xy} and b_{yx} is equal to or greater than r .

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

(iii) Regression co-efficients are independent of change of origin but not to side.

(iv) If one of the regression co-efficient is greater than unity the other one will be less than unity because correlation co-efficient lie between -1 & $+1$.

(v) In regression co-efficient if any of the three values are given the fourth one can be found out.

DIFFERENCE BETWEEN CORRELATION AND REGRESSION

CORRELATION	REGRESSION
<ul style="list-style-type: none">• Correlation is the relationship between two or more variables, which vary in sympathy with the other in the same (or) the opposite direction.	<ul style="list-style-type: none">• Regression means going back and it is a mathematical measure showing the average relationship between two variables.
<ul style="list-style-type: none">• Both the variables X and Y are random variables.	<ul style="list-style-type: none">• Here X is a random variable and Y is a fixed variable. Sometimes both variables may be random variables.
<ul style="list-style-type: none">• It finds out the degree of relationship between two variables are not their cause and effect of the variables.	<ul style="list-style-type: none">• It indicates the cause and effect relationship the variables and establishes a functional relationship.
<ul style="list-style-type: none">• There may be non-sense correlation b/n two variables.	<ul style="list-style-type: none">• There is no such non-sense regression.

- It is used for Testing and Verifying the relationship between two variables and series limited information.

- Besides, Verifications it is used for the prediction of one value in relationship to the other value.

- The co-efficient of correlation is a relative measure. The range of relationship lies b/n the variable.

- The regression co-efficient is an absolute figure. If we know the value of the independent variable, we can find the value of the dependent variable.

- It is not veryful for further mathematical treatment.

- It is widely used for further mathematical treatment.

- If the co-efficient of correlation is positive, then their two variables are positively correlated and vice versa.

- Regression co-efficient explicit that decrease in one variable is associated with the increase in the other variables.

- It has limited application because, it is confined only to linear relationship b/n the variables.

- It has wider application as it studies linear and non-linear relationship b/n the variables.

- It is immediately whether X depends upon Y (or) Y depends upon X.

- There is a functional relationship b/n the two variables. so, that we may identify b/n the independent and dependent variables.

UNIT - V

ASSOCIATION OF ATTRIBUTE

NOTATIONS AND DEFINITIONS:

(i) POSITIVE AND NEGATIVE CLASSES

The attributes may be positive or negative. If the attribute is present it is termed as positive class.

They are denoted by A, B, C . . .

If the attribute is not present, it is termed as negative class.

They are denoted by greek letters α, β, γ . . .

(ii) ULTIMATE CLASS FREQUENCIES:

The number of values which possess a particular attribute is called the class frequencies.

Those classes which specify the attributes of the highest order are known as the ultimate class and their frequencies are known as the ultimate class frequencies.

In the study of Attribute (A) and (B) are the ultimate class frequency (AB) ($\alpha\beta$) (AB) ($\alpha\beta$) formula 2^n .

In the study formula, the number of ultimate class frequency $= 2^n$.

(iii) ORDER OF CLASSES:

The order of a class depends upon the no. of

attributes under study. A class having one attribute is known as the class of 1st order; A class having two attributes is known as the second order class. N denotes zero order class.

FIRST ORDER - A, B, α , β

SECOND ORDER - AB, αB , $\alpha \beta$, $A\beta$

N - ZERO ORDER

CONTINGENCY TABLE (OR) NINE SQUARE TABLE:

The class frequency in a study of two attributes

A and B are represented by following table.

	A	α	
B	(AB)	(αB)	(B)
β	(A β)	($\alpha \beta$)	(β)
	(A)	(α)	N

Since it has parts 3^2 which represent the class frequency which also called nine square table.

RELATIONSHIP:

(i) $(A) = (AB) + (A\beta)$

(ii) $(\alpha) = (\alpha B) + (\alpha \beta)$

(iii) $(B) = (AB) + (\alpha B)$

(iv) $(\beta) = (A\beta) + (\alpha \beta)$

(v) $N = (A) + (\alpha)$ (or) $N = (B) + (\beta)$ (or) $N = (AB) + (\alpha B) + (A\beta) + (\alpha \beta)$

PROBLEM

Calculate the missing value and fill up the contingency table:

$$(AB) = 35, (A) = 55, N = 100 \text{ \& } (B) = 65$$

	A	α	
B	35	30	65
β	20	15	35
	55	45	100

$$(\alpha) = (N - A) = 100 - 55 = 45$$

$$(\beta) = (N - B) = 100 - 65 = 35$$

$$(A\beta) = A - (AB) = 55 - 35 = 20$$

$$(\alpha B) = (B) - (AB) = 65 - 35 = 30$$

$$(\alpha\beta) = (\alpha) - (\alpha B) = 45 - 30 = 15$$

TYPES OF ASSOCIATION:

There are 3 types,

- i) Positive Association
- ii) Negative Association
- iii) Independent Association

POSITIVE ASSOCIATION

when two attributes are present (or) absent together in the data, and actual frequency is more than

the expected frequency is called positive association.

$$\text{i.e., } (AB) > \frac{(A) \times (B)}{N}$$

(actual) (expected)

NEGATIVE ASSOCIATION

when the existence of one attribute causes absence of another attribute and actual frequency is less than the expected frequency is called negative association.

$$\text{i.e., } (AB) < \frac{(A) \times (B)}{N}$$

(actual) (expected)

INDEPENDENT ASSOCIATION

when there exist no association two attributes (or) when they have no tendency to be present together (or) the presence of one attribute are set to be independent actual frequency is equal to the expected frequency.

$$\text{i.e., } AB = \frac{(A) \times (B)}{N}$$

(actual) (expected)

METHODS OF DETERMINING ASSOCIATION:

YULE'S CO-EFFICIENT OF ASSOCIATION:

The above mentioned methods will give us as a

rough idea about their association but the association cannot be found out. Prof. Yule, has suggested a formula to measure the association (Q) false between ± 1 .

If $Q=0$, no association; If $Q=+1$, there is perfect positive association; If $Q=-1$, there is perfect negative association.

PROBLEM:

Find association between literacy and unemployment from the following frequencies.

$$\text{TOTAL ADULTS} = 10,000 (N)$$

$$\text{LITERATE} = 1290 (A)$$

$$\text{UNEMPLOYED} = 1390 (B)$$

$$\text{LITERATE UNEMPLOYED} = 820 (AB)$$

soln:

	A	α	
B	820		1390
β			
	1290		10,000

$$\alpha = N - A = 10,000 - 1290 = 8710$$

$$\beta = N - B = 10,000 - 1390 = 8610$$

$$(A\bar{B}) = A - AB = 1290 - 820 = 470$$

$$(\bar{\alpha}B) = B - AB = 1390 - 820 = 570$$

$$(\alpha\bar{\beta}) = \alpha - \alpha B = 8710 - 570 = 8140$$

Since it is a positive, and consistent.

$$R = \frac{(820)(8140) - (470)(570)}{(820)(8140) + (470)(570)}$$

$$R = \frac{6,406,900}{69,42,700} = 0.9228$$

There is a positive association b/w the attributes of A and B.

	A	α	
B	820	570	1390
\bar{B}	470	8140	8610
	1290	8710	10,000