

CA FOUNDATION - STATISTICS CHARTS BY - PROF. JATIN DEMBLA

STATISTICAL DESCRIPTION OF DATA

The word statistics has been derived either of the following Latin word "Status" which means a political State. Italian word "Statista" German word "Statistik" French word "Statistique"

CHARACTERISTICS:

- Statistics are aggregate of facts.
- Statistics are affected by a large number of causes.
- Statistics are always numerically expressed.
- Statistics should be enumerated or estimated.
- Statistics should be collected in systematic manner.
- Statistics should be collected for pre-determined purpose.
- Statistics should be placed in relation to each other.

Important limitations

- Statistics does not deal with individual items.
- Statistics deals with quantitative data.
- Statistics laws are true only on averages.
- Statistics does not reveal the entire story.
- Statistics is liable to be misused.
- Statistics data should be uniform and homogeneous.

MEASURES OF CENTRAL TENDENCY

Primary data	Secondary data
It's the data collected by a particular person or organization for his own use from the primary sources.	It's the data collected by some other person or organization for their own use but the investigator also gets it for his use.
<ol style="list-style-type: none"> Direct personal observation: In this method, the investigator collects the data personally and therefore, it gives reliable and correct information. USE: In case of natural calamity data can be collected more quickly and accurately by applying this method. Indirect oral investigation: In this method, a third person is contacted who is expected to know the necessary details about the persons for whom the enquiry is meant. USE: If there are some practical problems in reaching the respondents directly, as in case of rail accident. Estimates from the local sources and correspondence: Here the investigator appoints agents and correspondents to collect the data. Data through Questionnaire: The data can be collected by preparing a Questionnaire and getting it filled by the persons concerned. Investigations through enumerators: This method generally employed by the Government for population census etc. 	<ol style="list-style-type: none"> Information collected through newspapers and periodicals. Information obtained from the publications of trade associations. Information obtained from the research papers published by university departments or research bureaus or I.G.C. Information obtained from the official publications of the Central, State and the local governments dealing with crop statistics, Industrial statistics, Trade and transport statistics etc. Information obtained from the official publications of the foreign governments for international organizations

MEASURES OF DISPERSION

Arithmetic mean or mean is the number which is obtained by adding the values of all the items of a series and dividing the total by the number of items. When all items of a series are given equal importance then it is called simple arithmetic mean and when different items of a series are given different weights according to their relative importance is known as weighted arithmetic mean.

Median is the middle value of the series when arranged in order of the magnitude. When a series is divided into more than two parts, the dividing values are called Partition values.

Quartiles are the measures which divide the data into four equal parts, each portion contains equal number of observation

There are three quartiles

If a statistical series is divided into four equal parts, the end value of each part is called a quartile and denoted by 'Q'.

The lower half of a data set is the set of all values that are to the left of the median value when the data has been put into increasing order.

The upper half of a data set is the set of all values that are to the right of the median value when the data has been put into increasing order.

The first quartile, denoted by Q1 is the median of the lower half of the data set. This means that about 25% of the numbers in the data set lie below Q1 and about 75% lie above Q1

The second quartile also called median and denoted by Q2 has 50% of the items below it and 50% of the items above it.

The third quartile, denoted by Q3, is the median of the upper half of the data set. This means that about 75% of the numbers in the data set lie below Q3 and about 25% lie above Q3

Main purposes and functions of averages.

- To represent a brief picture of data.
- Comparison.
- Formulation of policies.
- Basis of statistical analysis.
- One value for all the group or series.

Merits of Median:

- Simple measure of central tendency.
- It is not affected by extreme observations.
- Possible even when data is incomplete.
- Graphic can be determined by median presentation of data.
- It has a definite value.
- Simple to calculate and understand
- It is a positional value not a calculated value.

Merits of mode:

- Simple and popular measure of central tendency.
- It can be located graphically with the help of histogram.
- Useful in knowing all the items of series.
- It is the most representative value in the given series.
- It is less affected by extreme values.

Types of series	Direct Method	Shortcut Method	Step deviation method	Measure	Individual Series	Discrete Series	Continues Series
Individual series	$\bar{X} = \frac{\sum X}{N}$	$\bar{X} = A + \frac{\sum d}{N} \times C$	$\bar{X} = A + \frac{\sum fd}{N} \times C$	Size of item	Size of item	Size of item	Size of item
Discrete Series	$\bar{X} = \frac{\sum fx}{N}$	$\bar{X} = A + \frac{\sum fd}{N} \times C$	$\bar{X} = A + \frac{\sum fd}{N} \times C$	First Quartile	Item	Item	Item
Continuous Series	$\bar{X} = \frac{\sum fm}{N}$	$\bar{X} = A + \frac{\sum fd}{N} \times C$	$\bar{X} = A + \frac{\sum fd}{N} \times C$	Third Quartile	Item	Item	Item

Combines Mean $\bar{x}_{12} = \frac{x_1 n_1 + x_2 n_2}{n_1 + n_2}$

Weighted Mean $\bar{X} = \frac{\sum Wx}{\sum W}$

DISPERSION

Characteristics of Measures of Dispersion

- A measure of dispersion should be rigidly defined.
- It must be easy to calculate and understand.
- Not affected much by the fluctuations of observations Based on all observations

Classification of Measures of Dispersion

The measure of dispersion is categorized as:

- An absolute measure of dispersion:
 - The measures which express the scattering of observation in terms of distances i.e., range, quartile deviation.
 - The measure which expresses the variations in terms of the average of deviations of observations like mean deviation and standard deviation.
- A relative measure of dispersion:
 - We use a relative measure of dispersion for comparing distributions of two or more data set and for unit free comparison. They are the coefficient of range, the coefficient of mean deviation, the coefficient of quartile deviation, the coefficient of variation, and the coefficient of standard deviation.

Variance of the Combined Series

If σ_1, σ_2 are two standard deviations of two series of sizes n_1 and n_2 with means \bar{y}_1 and \bar{y}_2 , the variance of the two series of sizes $n_1 + n_2$ is:

$$\sigma^2 = \frac{1}{n_1 + n_2} [n_1(\sigma_1^2 + d_1^2) + n_2(\sigma_2^2 + d_2^2)]$$

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

Coefficient of Dispersion

Whenever we want to compare the variability of the two series which differ widely in their averages. Also, when the unit of measurement is different. We need to calculate the coefficients of dispersion along with the measure of dispersion. The coefficients of dispersion (C.D.) based on different measures of dispersion are:

- Based on Range = $(X_{max} - X_{min}) / (X_{max} + X_{min})$
- C.D. based on quartile deviation = $(Q3 - Q1) / (Q3 + Q1)$
- Based on mean deviation = Mean deviation / average from which it is calculated.
- For Standard Deviation = S.D. / Mean

Standard Deviation

Standard deviation is the positive square root of the arithmetic mean of the squares of the deviations of the given values from their arithmetic mean. It is denoted by a Greek letter sigma, σ . It is also referred to as root mean square deviation. The standard deviation is given as:

$$\sigma = \sqrt{\frac{\sum (y_i - \bar{y})^2}{n}} = \sqrt{\frac{\sum y_i^2/n - \bar{y}^2}{n}}$$

For a grouped frequency distribution, it is:

$$\sigma = \sqrt{\frac{\sum f_i (y_i - \bar{y})^2}{N}} = \sqrt{\frac{\sum f_i y_i^2/n - \bar{y}^2}{N}}$$

The square of the standard deviation is the variance. It is also a measure of dispersion.

Coefficient of Variation

100 times the coefficient of dispersion based on standard deviation is the coefficient of variation (C.V.).

$$C.V. = 100 \times (S.D. / \text{Mean}) = (\sigma / \bar{y}) \times 100$$

PROBABILITY

The sum of all the probabilities in the sample space is 1.

The probability of an event which cannot occur is 0.

The probability of any event which is not in the sample space is zero.

The probability of an event which must occur is 1.

The probability of the sample space is 1.

The probability of an event E is denoted as P(E) and is written as $P(E) = \frac{n(E)}{n(S)}$

The probability of the complement of an event E is denoted as $P(\bar{E}) = 1 - P(E)$

P(A|B) is written as $P(A \cap B)$ and P(A|B) is written as P(A|B)

If A and B are mutually exclusive events, $P(A \cup B) = P(A) + P(B)$

When two events A and B are independent i.e. when event A has no effect on the probability of event B, the conditional probability of event B given event A is simply the probability of event B, that is P(B).

If events A and B are not independent then the probability of the intersection of A and B (the probability that both events occur) is defined by $P(A \cap B) = P(A)P(B|A)$

Regression

Regression analysis encompasses to identify the relationship between a dependent variable and one or more independent variables. Regression calculates the "best fit" line for a certain set of data. The regression line makes the sum of the squares of the residuals smaller than for any other line. In regression analysis, a single dependent variable, Y_i is considered to be a function of one or more independent variables, X_1, X_2 , and so on. The values of both the dependent and independent variables are assumed as being determined in an error-free random manner.

Index Numbers: An index number is an economic data figure reflecting price or quantity compared with a standard or base value. OF INDEX NUMBERS index numbers are names after the activity they measure. Their types are as under:

Price Index: Measure changes in price over a specified period of time. It is basically the ratio of the price of a certain number of commodities at the present year as against base year.

Quantity Index: As the name suggest, these indices pertain to measuring changes in volumes of commodities like goods produced or goods consumed, etc.

Value Index: These pertain to compare changes in the monetary value of imports, exports, production or consumption of commodities.

Simple index numbers: A simple index number index numbers is a number that expresses the relative change in price, quantity, or value from one period to another. Let p_0 be the base period price, and p_1 be the price at the selected or given period. Thus, the simple price index is given by: $P = p_1 / p_0 \times 100$

Weighted index

Weighted Index Numbers = $(\sum \text{index number} \times \text{weight}) / \sum \text{weight}$

There are two types of weighted indices, they are Lapeyre's index

INDEX NUMBER AND TIME SERIES

Index Number Construction

- According to the Fundamental Method $Q_{11} = \frac{\sum p_1 q_0}{\sum p_0 q_0} \times 100$
- According to simple aggregative method $Q_{11} = \frac{\sum p_1}{\sum p_0} \times 100$
- According to simple relative method $Q_{11} = \frac{\sum \frac{p_1 q_0}{p_0 q_0}}{\sum \frac{p_0 q_0}{p_0 q_0}} \times 100$ or $\frac{\sum \frac{p_1 q_0}{p_0 q_0}}{\sum 1} \times 100$
- According of the general weighted aggregative method $Q_{11} = \frac{\sum w \frac{p_1 q_0}{p_0 q_0}}{\sum w} \times 100$
- According to Laspeyre's method $Q_{11(p)} = \frac{\sum p_1 q_0}{\sum p_0 q_0} \times 100$
- According to Paasche's aggregative method $Q_{11(p)} = \frac{\sum p_1 q_1}{\sum p_0 q_1} \times 100$
- According to Drobish and Bowley's aggregative method $Q_{11} = \frac{\sum p_1 q_0 + \sum p_0 q_1}{2 \sum p_0 q_0} \times 100$
- According to Fisher's ideal method: $Q_{11(F)} = \sqrt{\frac{\sum p_1 q_0}{\sum p_0 q_0} \times \frac{\sum p_1 q_1}{\sum p_0 q_1}} \times 100$
- According to Fisher's ideal method: $Q_{11(F)} = \sqrt{\frac{\sum p_1 q_0}{\sum p_0 q_0} \times \frac{\sum p_1 q_1}{\sum p_0 q_1}} \times 100$

Time Series Component Analysis

- Used primarily for forecasting
- Observed value in time series in the sum or product of components
- Additive model $X_t = T_t + S_t + C_t + I_t$
- Multiplicative model (Linear in log Form) $X_t = T_t S_t C_t I_t$

THEORETICAL DISTRIBUTION

Binomial Distribution:-

Mean = np
Variance = npq
• Variance is Always Less than Mean
• Note - Variance Will be Highest When $p = q = 0.5$
• Mode = (n+1)P (If Mode is Integer then it will be Bi-Modal & If Non Integer than will be Uni-Modal)

Poisson Distribution:-

Mean = m
Variance = m
Mode = As Same as Binomial Distribution
• Prob of Success is Very Small
• It is Uniparametric Distribution

Normal Distribution Or Bell Curve Or Gaussian curve

• Mean, Median & Mode are Concise
• 1st & 3rd Quartile
a. $Q1 = \mu - 0.675 \sigma$
b. $Q3 = \mu + 0.675 \sigma$
• 2 Points of Inflection

If a discrete random variable X has the following probability density function (p.d.f), it is said to have a binomial distribution:

$$P(X = x) = nC_x q^n p^x$$

Note: nCr ("n choose r") is more commonly written, but I shall use the former because it is easier to write on a computer. It means the number of ways of choosing r objects from a collection of n objects (see permutations and combinations).

If a random variable X has a binomial distribution, we write $X \sim B(n, p)$ ("means" has distribution...").

n and p are known as the parameters of the distribution (n can be any integer greater than 0 and p can be any number between 0 and 1).

All random variables with a binomial distribution may have different parameters (different values for n and p).

Expectation and Variance

If $X \sim B(n, p)$, then the expectation and variance is given by:

- $E(X) = np$
- $Var(X) = npq$

Binomial & Normal Distribution Is Known As Biparametric Distribution

CORRELATION AND REGRESSION

Correlation

Correlation is a statistical method used to determine the extent to which two variables are related. Correlation analysis measures the degree of association between two or more variables. Parametric methods of correlation analysis assume that for any pair or set of values taken under a given set of conditions, variation in each of the variables is random and follows a normal distribution pattern. In scattered diagram, following elements are represented:

- Rectangular coordinate
- Two quantitative variables
- One variable is called independent (X) and the second is called dependent (Y)
- Points are not joined
- No frequency table

Correlation coefficient is the Statistic showing the degree of relation between two variables. The simple correlation coefficient, denoted r , ranges between -1 and +1 and quantifies the direction and strength of the linear association between the two variables. It is also called Pearson's correlation or product moment correlation coefficient. It measures the nature and strength between two variables of the quantitative type.

Regression

Regression analysis encompasses to identify the relationship between a dependent variable and one or more independent variables. Regression calculates the "best fit" line for a certain set of data. The regression line makes the sum of the squares of the residuals smaller than for any other line. In regression analysis, a single dependent variable, Y_i is considered to be a function of one or more independent variables, X_1, X_2 , and so on. The values of both the dependent and independent variables are assumed as being determined in an error-free random manner.

Correlation

- In correlation analysis the degree and direction of relationship between the variables are studied
- If value of one variable is known, the value of other variable cannot be estimates
- Correlation coefficient lies between -1 & 1
- Correlation coefficient is independent of change of origin and scale
- With the help of correlation coefficient and standard deviation of two random variable (X, Y) regression coefficient can be obtained

Regression

- In regression analysis, the nature of relationship is studied
- If value of one variable is known, the value of other variable can be estimated using the functional relationships
- Only one relation coefficient can be greater than 1
- Regression coefficient is independent of change of origin but not of scale