

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 item.
- Statistics deals with quantitative data.
- Statistics fails to reveal the entire story.
- Statistics is liable to be misused.
- Statistics data should be uniform and homogeneous.

It's the data collected by a particular person or organization for his own use from the primary sources.

Primary data

- 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 the 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.

It is the data collected by some other person or organization for their own use but the investigator also gets it for his own use.

Secondary data

- 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 IGC.
- 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.

MODE OF PRESENTATION OF DATA

- Tabular presentation or Tabulation
- Diagrammatic presentation:
 - Line diagram
 - Bar diagram
 - Pie chart

MEASURES OF CENTRAL TENDENCY

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 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.

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 in 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 in 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.

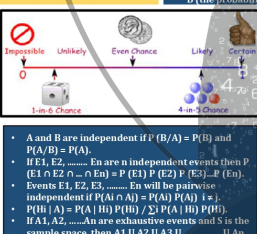
The second quartile also called median and denoted by Q2 has 50% of the items below it and 50% of the items above it.

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.

Formulae of calculating mode in continuous series Mode =

$$L_1 + \frac{f_1 - f_0}{f_1 - f_0 + f_1 - f_2} \times (L_2 - L_1)$$

- Where L1 = Lower limit of modal class
- f0 = Frequency of the group preceding the modal class
- f1 = Frequency of the modal class
- f2 = Frequency of the group succeeding the modal class
- C = Class interval of the modal class



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 an event not occurring is one minus the probability of it occurring.
- The complement of an event E is denoted as E' and is written as P(E') = 1 - P(E).
- P(A or B) is written as P(A ∪ B) and P(A and B) is written as P(A ∩ B).
- If A and B are mutually exclusive events, P(A or B) = P(A) + P(B).
- When two events A and B are independent, then the probability of both occurring together is simply the probability of event A, that is P(A).
- 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 and B) = P(A ∩ B).

The probability distribution of a count variable X is said to be the binomial distribution with parameter n and abbreviated B(n, p) if it satisfies the following conditions:

- The total number of observations is fixed.
- The observations are independent.
- Each outcome represents either a success or a failure.
- The probability of success i.e. p is same for every outcome.

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 Arithmetic mean:**
- Simplicity
 - Certainty
 - Based on all values.
 - Algebraic treatment possible.
 - Basis of comparison
 - Accuracy test possible
 - No scope for estimated value

- Merits of Median:**
- Simple measure of central tendency.
 - It is not affected by extreme observations.
 - Possible even when data is incomplete.
 - Median can be determined by graphical 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.
 - Less effect of marginal values.
 - No need of 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 | Measur | 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 | $\bar{x} = \frac{\sum fx}{\sum f}$ | $\bar{x} = A + \frac{\sum fd}{N}$ | $\bar{x} = A + \frac{\sum fd}{N} \times C$ | First Item | Item | Item | Item |
| Continuous Series | $\bar{x} = \frac{\sum fm}{\sum f}$ | $\bar{x} = A + \frac{\sum fd}{N}$ | $\bar{x} = A + \frac{\sum fd}{N} \times C$ | Third Quartile | Item | Item | Item |
| Combines Mean $\bar{x} = \frac{\sum x_1 f_1 + \sum x_2 f_2}{\sum f}$ | | | | Weighted Mean $\bar{x} = \frac{\sum fx}{\sum w}$ | | | |

Instagram.com/jatinDembla1

DISPERSTION

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

- Absolute measure of dispersion:
 - 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.
- Relative measure of dispersion:
 - We use a relative measure of dispersion for comparing distributions of two or more data set and for limit test comparison. Here 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 $\sigma^2 = \frac{n_1 \sigma_1^2 + n_2 \sigma_2^2}{n_1 + n_2} + \frac{n_1 (\bar{y}_1 - \bar{y})^2 + n_2 (\bar{y}_2 - \bar{y})^2}{n_1 + n_2}$ where, $\bar{y} = \frac{n_1 \bar{y}_1 + n_2 \bar{y}_2}{n_1 + n_2}$ and $\bar{y} = \frac{n_1 \bar{y}_1 + n_2 \bar{y}_2}{n_1 + n_2}$

Range - Range = $X_{max} - X_{min}$
Quartile Deviation - $Q_3 - Q_1 \div 2$
Mean Deviation
 Mean deviation is the average of the absolute deviations of the observations from a measure of central tendency. If x_1, x_2, \dots, x_n are the set of observations, then the mean deviation of about the average A (mean, median, or mode) is

Mean deviation from average $A = \frac{1}{N} \sum |x_i - A|$
 For a grouped frequency, it is calculated as
 Mean deviation from average $A = \frac{1}{N} \sum f_i |x_i - A| = \frac{\sum f_i |x_i - A|}{N}$
 Here, x_i and f_i are respectively the mid value and the frequency of the i^{th} class interval.

Standard Deviation
 A 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 $\sigma = \sqrt{\frac{\sum (y_i - \bar{y})^2}{N}} = \sqrt{\frac{\sum (y_i^2 - 2\bar{y}y_i + \bar{y}^2)}{N}}$
 For a grouped frequency distribution, it is $\sigma = \sqrt{\frac{\sum (f_i y_i^2 / N) - (\sum (f_i y_i / N))^2 / N}{N}}$
 The square of the standard deviation is the variance. It is also a measure of dispersion.
 $\sigma^2 = \frac{\sum (y_i - \bar{y})^2}{N} = \frac{\sum (y_i^2 - 2\bar{y}y_i + \bar{y}^2)}{N}$
 For a grouped frequency distribution, it is $\sigma^2 = \frac{\sum (f_i y_i^2 / N) - (\sum (f_i y_i / N))^2 / N}{N}$
 Instead of a mean, we choose any other arbitrary number, say A, the standard deviation becomes the root mean deviation.

Coefficient of Variation
 CV measures the coefficient of dispersion based on standard deviation is the coefficient of variation (C.V.).
 $C.V. = 100 \times \frac{S.D.}{\text{Mean}} = \frac{\sigma}{\bar{y}} \times 100$

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 attempts 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, is considered to be a function of one or more independent variables, X1, X2, and so on. The values of both the dependent and independent variables are assumed as being determined in an error-free random manner.

INDEX NUMBER AND TIME SERIES

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 suggests, 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 p0 be the base period price, and p1 be the price at the selected or given period. Thus, the simple price index is given by: $P = \frac{p_1}{p_0} \times 100$

Weighted index
Weighted Index Numbers = $(\sum \text{index number} \times \text{weight}) / \sum \text{weight}$
 There are two types of weighted indexes, they are Lapeyre's index

THEORETICAL DISTRIBUTION

Binomial Distribution:-

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

Poisson Distribution:-

- Mean = m
- Variance = m
- Mode = 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 Coincide
- 1st & 3rd Quartile
- a. $Q1 = \mu - 0.675 \sigma$
- b. $Q3 = \mu + 0.675 \sigma$
- 2 Points of Inflexion

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

- $P(X=x) = {}^n C_x (q)^x (p)^{n-x}$, where $q = 1 - p$ can be considered as the probability of a success, and q the probability of a failure.
- Note: nC_x ("n choose x") is more commonly written as $\binom{n}{x}$ and p and q 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 have 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:

- $E(X) = np$
- $\text{Var}(X) = npq$

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

Regression

In a simple regression analysis, one dependent variable is measured in relation to only one independent variable. The analysis is designed to develop an equation for the line that best models the relationship between the dependent and independent variables. This equation has the mathematical form:

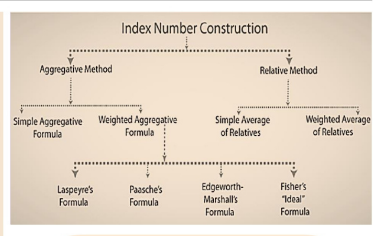
$Y = a + bX$
 In above equation, Y is the value of the dependent variable, X is the value of the independent variable is the intercept of the regression line on the Y axis when X = 0, and b is the slope of the regression line.

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 estimated.

Regression
 In regression analysis, the nature of relationship is studied. If value of variable is known, the value of other variable can be estimated using the functional relationships.

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 (XY) regression coefficient can be obtained

Regression coefficient is independent of change of origin but not of scale



Time Series Component Analysis

- Used primarily for forecasting
- Observed value in time series in the sum or product of components
- Additive model $Tt = Tt + St + Ct + It$
- Multiplicative model (Linear in Log Form) $Tt = TStCIt$

Kinshuk Institute