

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 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.

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

- 1. Direct personal observation In this method, the investigator collects the data personally and therefore, it gives reliable and correct information.
2. 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.
3. Estimates from the local sources and the correspondence.
4. Data through Questionnaire The data can be collected by preparing a Questionnaire and getting it filled by the persons concerned.
5. Investigations through enumerators. This method generally employed by the Government for population census etc.

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

- 1. Information collected through newspapers and periodicals.
2. Information obtained from the publications of trade associations.
3. Information obtained from the research papers published by university departments or research bureaus or IGC.
4. 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.
5. Information obtained from the official publications of the foreign governments for international organizations.

MODE OF PRESENTATION OF DATA

- Textual presentation
• Tabular presentation or Tabulation
• Diagrammatic presentation:
Types of diagrams
I. Line diagram
II. Bar diagram
III. 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.

Median is the middle value of the series when arranged in order of the magnitude.

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'.

Formula of calculating mode in continuous series Mode =

- Where L = 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 = Magnitude or class interval of the modal class

Deciles distribute the series into ten equal parts and generally expressed as D. Percentiles divide the series into hundred equal parts and generally expressed as P.

Mode is the value which occurs most frequently in the series, that modal value has the highest frequency in the series.

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.
(i) To represent a brief picture of data.
(ii) Comparison.
(iii) Formulation of policies.
(iv) Basis of statistical analysis.
(v) 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:
(i) Simple measure of central tendency.
(ii) It is not affected by extreme observations.
(iii) Possible even when data is incomplete.
(iv) Median can be determined by graphical presentation of data.
(v) It has a definite value.
(vi) Simple to calculate and understand
(vii) It is a positional value not a calculated value.

- Merits of mode:
(i) Simple and popular measure of central tendency.
(ii) It can be located graphically with the help of histogram.
(iii) Less effect of marginal values.
(iv) No need of knowing all the items of series.
(v) It is the most representative value in the given series.
(vi) It is less effected by extreme values.

Table with columns: Types of series, Direct Method, Shortcut Method, Step deviation method, Measures, Individual Series, Discrete Series, Continues Series. Includes formulas for Mean, Median, and Mode for various series types.

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
The nature of dispersion is categorized as:
(i) Absolute measure of dispersion:
• Absolute measures 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.
(ii) A 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 sigma_1, sigma_2 are two standard deviations of two series of sizes n_1 and n_2 with means y_1 and y_2. The variance of the two series of sizes n_1 + n_2 is sigma^2 = [(n_1 * sigma_1^2 + n_2 * sigma_2^2) + n_1 * (sigma_1^2 + d_1^2) + n_2 * (sigma_2^2 + d_2^2)] / (n_1 + n_2) where, d_1 = y_1 - y_bar, d_2 = y_2 - y_bar and y_bar = (n_1 * y_1 + n_2 * 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

Range - Range = X_max - X_min
Quartile Deviation - Q = 1/2 * (Q3 - Q1)
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, ..., 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 = 1/n * sum(|x_i - A|)
For a grouped frequency, it is calculated as:
Mean deviation from average A = 1/N * sum(f_i * |x_i - A|) = 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, sigma. It is also referred to as root mean square deviation. The standard deviation is given as:
sigma = sqrt([(x_1 - y_bar)^2 / n] + [(x_2 - y_bar)^2 / n] + ... + [(x_n - y_bar)^2 / n])
For a grouped frequency distribution, it is:
sigma = sqrt([(f_1 * (y_1 - y_bar)^2) / N] + [(f_2 * (y_2 - y_bar)^2) / N] + ... + [(f_n * (y_n - y_bar)^2) / N])
The square of the standard deviation is the variance. It is also a measure of dispersion.
sigma^2 = [(x_1 - y_bar)^2 / n] + [(x_2 - y_bar)^2 / n] + ... + [(x_n - y_bar)^2 / n]
For a grouped frequency distribution, it is:
sigma^2 = [(f_1 * (y_1 - y_bar)^2) / N] + [(f_2 * (y_2 - y_bar)^2) / N] + ... + [(f_n * (y_n - y_bar)^2) / N]
Instead of a mean, we choose any other arbitrary number, say A, the standard deviation becomes the root mean deviation.

Coefficient of Variation
100 times the coefficient of dispersion based on standard deviation is the coefficient of variation (C.V.).
C.V. = 100 * (S.D. / Mean) = (sigma / y_bar) * 100.

Correlation coefficient is the Statistic showing the degree of relation between two variables. The simple correlation coefficient, denoted by 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 concerns 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 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 = p1 / p0 * 100

Weighted index: Weighted Index Numbers = (sum(index number * weight)) / sum(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=q=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
Same as Binomial Distribution
Prob Of Success is Very Small
It is Uniparametric Distribution
2 Points of Inflexion

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

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) = nCx (q)^n-x (p)^x, where q = 1 - p
p can be considered as the probability of a success, and q the probability of a failure.
Note: nCn (1 choose 1) is more commonly written as 1, 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 ~ B(n, p) (n means the 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 have may have different parameters (different values for n and p).

Expectation and Variance
If X ~ 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 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:
1. Rectangular coordinate
2. Two quantitative variables
3. One variable is called independent (X) and the second is called dependent (Y)
4. Points are not joined
5. No frequency table

Correlation coefficient is the Statistic showing the degree of relation between two variables. The simple correlation coefficient, denoted by 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.

In scatter plot, the pattern of data is indicative of the type of relationship between two variables. It may be:
1. Positive relationship
2. Negative relationship
3. No relationship
The sign of the correlation coefficient indicates the direction of the association. The magnitude of the correlation coefficient indicates the strength of the association.

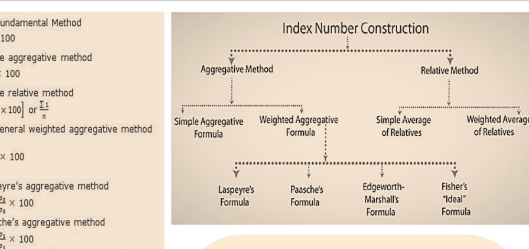
In correlation analysis, one dependent variable is measured in relation to only one independent variable. The analysis is designed to develop an equation for the line that best models the relationship between the dependent and independent variables. This equation has the mathematical form:
Y = a + bX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

In regression analysis, the degree and direction of relationship between the variables are studied. If the value of one variable is known, the value of other variable can be estimated.
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

In regression analysis, the nature of relationship is studied. If the value of 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

INDEX NUMBER AND TIME SERIES

- 1. According to the Fundamental Method Q11 = (I1 / I0) * 100
2. According to simple aggregative method Q11 = (sum I1 / sum I0) * 100
3. According to simple relative method Q11 = (sum (I1 / I0)) * 100 or sum (I1 / I0)
4. According of the general weighted aggregative method Q11 = (sum (I1 * W)) / (sum (I0 * W)) * 100
5. According to Laspeyre's aggregative method Q11 = (sum (I1 * P0)) / (sum (I0 * P0)) * 100
6. According to Paasche's aggregative method Q11 = (sum (I1 * P1)) / (sum (I0 * P1)) * 100
7. According to Drobnah and Bowley's aggregative method Q11 = 1/2 * ((sum (I1 * P0) / sum (I0 * P0)) + (sum (I1 * P1) / sum (I0 * P1))) * 100
8. According to Fisher's ideal method: Q11 = (sum (I1 * P0) / sum (I0 * P0)) * (sum (I1 * P1) / sum (I0 * P1))^0.5
9. According to Fisher's ideal method: Q11 = (sum (I1 * P0) / sum (I0 * P0)) * (sum (I1 * P1) / sum (I0 * P1))^0.5



Time Series Component Analysis
• Used primarily for forecasting
• Observed value in time series in the sum or product of components
• Additive model Xt = Tt + St + Ct + It
• Multiplicative model (Linear in Log Form) Xt = TStCIt

Kinshuk Institute

