

24BA1101

STATISTICS TOOL FOR MANAGEMENT

LT P C

4 0 0 4

## COURSE OBJECTIVES

- To introduce the concepts of basic statistical techniques.
- To acquaint the knowledge of hypothesis testing, using tools such as Z test, F test, ANOVA, chi-square tests and non-parametric tests.
- To introduce the concepts solving regression & correlation problems.

### UNIT I INTRODUCTION

12

Measures of central tendency – Mean – Median – Mode – Measures of Dispersion– Range – Quartile Deviation – Mean Deviation – Standard Deviation and Co-efficient of Variation – Skewness – Bowley's Co-efficient – Kelley's Co-efficient.

### UNIT II SAMPLING DISTRIBUTION AND ESTIMATION

12

Introduction to sampling distributions – Sampling distribution of mean and proportion – Sampling techniques – Estimation: Point and Interval estimates for population parameters of large sample and small samples.

### UNIT III TESTING OF HYPOTHESIS – PARAMETRIC TESTS

12

Hypothesis testing: one sample and two sample tests for means and proportions of large samples (z-test), one sample and two sample tests for means of small samples (t-test) – F-test for two sample standard deviations – ANOVA one and two way.

### UNIT IV NON-PARAMETRIC TESTS

12

Chi-square test for single sample standard deviation – Chi-square tests for independence of attributes and goodness of fit – Rank sum test – Comparing two populations – Mann – Whitney U test and Kruskal Wallis test.

### UNIT V CORRELATION AND REGRESSION

12

Correlation – Coefficient of Determination – Rank Correlation – Regression – Estimation of Regression line – Standard Error of estimate.

**TOTAL: 60 PERIODS**

## TEXT BOOKS

1. Richard I. Levin, David S. Rubin, Masood H. Siddiqui, Sanjay Rastogi, Statistics for Management, Pearson Education, 8th Edition, 2017.
2. G C Beri, "Business Statistics", Tata Mc Graw Hill Publishing Company Ltd., 3rd Edition, 2017.
3. Gupta S.P., "Statistical Method", Sultan Chand & Sons, New Delhi, 46th Edition, 2019.

## REFERENCE BOOKS

1. Prem. S. Mann, Introductory Statistics, Wiley Publications, 9th Edition, 2015.
2. T N Srivastava and Shailaja Rego, Statistics for Management, Tata McGraw Hill, 3rd Edition 2017.
3. David R. Anderson, Dennis J. Sweeney, Thomas A. Williams, Jeffrey D. Camm, James J. Cochran, Statistics for business and economics, 13th Edition, Thomson (South – Western) Asia, Singapore, 2016.
4. N. D. Vohra, Business Statistics, Tata McGraw Hill, 2017.
5. Johnson, R.A., Miller.I and Freund J., "Miller and Freund's Probability and Statistics for Engineers", Pearson Education, Asia, 8th Edition, 2015.

Mapping of COs and POs COs			POs		
PO1	PO2	PO3	PO4	PO5	
CO1	3	3	1	2	2
CO2	3	3	2	2	2
CO3	2	3	2	2	2
CO4	3	3	3	2	3
CO5	3	3	3	2	2
AVG	2.80	3.00	2.20	2.00	2.20

## UNIT-I

### INTRODUCTION TO STATISTICS

#### DEFINITION OF STATISTICS

Branch of mathematics concerned with collection, classification, analysis, and interpretation of numerical facts, for drawing inferences on the basis of their quantifiable likelihood (probability). Statistics can interpret aggregates of data too large to be intelligible by ordinary observation because such data (unlike individual quantities) tend to behave in regular, predictable manner. It is subdivided into descriptive statistics and inferential statistics.

#### HISTORY OF STATISTICS

The Word statistics have been derived from Latin word —Status or the Italian word —Statistall, meaning of these words is —Political Statell or a Government. Shakespeare used a word Statist in his drama Hamlet (1602). In the past, the statistics was used by rulers. The application of statistics was very limited but rulers and kings needed information about lands, agriculture, commerce, population of their states to assess their military potential, their wealth, taxation and other aspects of government.

Gottfried Achenwall used the word statistik at a German University in 1749 which means that political science of different countries. In 1771 W. Hooper (Englishman) used the word statistics in his translation of Elements of Universal Erudition written by Baron B.F Bieford, in his book statistics has been defined as the science that teaches us what is the political arrangement of all the modern states of the known world. There is a big gap between the old statistics and the modern statistics, but old statistics also used as a part of the present statistics.

During the 18th century the English writer have used the word statistics in their works, so statistics has evolved gradually during last few centuries. A lot of work has been done in the end of the nineteenth century.

At the beginning of the 20th century, William S Gosset was developed the methods for decision making based on small set of data. During the 20th century several statistician are active in developing new methods, theories and application of statistics. Now these days the availability of electronics computers is certainly a major factor in the modern development of statistics.

## Descriptive Statistics and Inferential Statistics

Every student of statistics should know about the different branches of statistics to correctly understand statistics from a more holistic point of view. Often, the kind of job or work one is involved in hides the other aspects of statistics, but it is very important to know the overall idea behind statistical analysis to fully appreciate its importance and beauty.

The two main branches of statistics are descriptive statistics and inferential statistics. Both of these are employed in scientific analysis of data and both are equally important for the student of statistics.

## Descriptive Statistics

Descriptive statistics deals with the presentation and collection of data. This is usually the first part of a statistical analysis. It is usually not as simple as it sounds, and the statistician needs to be aware of designing experiments, choosing the right focus group and avoid biases that are so easy to creep into the experiment.

Different areas of study require different kinds of analysis using descriptive statistics. For example, a physicist studying turbulence in the laboratory needs the average quantities that vary over small intervals of time. The nature of this problem requires that physical quantities be averaged from a host of data collected through the experiment.

## MANAGERIAL APPLICATIONS OF STATISTICS

Statistics is a mathematical science involving the collection, analysis and interpretation of data. A number of specialties have evolved to apply statistical theory and methods to various disciplines. Certain topics have "statistical" in their name but relate to manipulations of probability distributions rather than to statistical analysis.

- **Actuarial science** is the discipline that applies mathematical and statistical methods to assess risk in the insurance and finance industries.
- **Astrostatistics** is the discipline that applies statistical analysis to the understanding of astronomical data.
- **Biostatistics** is a branch of biology that studies biological phenomena and observations by means of statistical analysis, and includes medical statistics.

- **Business analytics** is a rapidly developing business process that applies statistical methods to data sets (often very large) to develop new insights and understanding of business performance & opportunities
- **Chemometrics** is the science of relating measurements made on a chemical system or process to the state of the system via application of mathematical or statistical methods.
- **Demography** is the statistical study of all populations. It can be a very general science that can be applied to any kind of dynamic population, that is, one that changes over time or space.
- **Econometrics** is a branch of economics that applies statistical methods to the empirical study of economic theories and relationships.
- **Environmental statistics** is the application of statistical methods to environmental science. Weather, climate, air and water quality are included, as are studies of plant and animal populations.
- **Geostatistics** is a branch of geography that deals with the analysis of data from disciplines such as petroleum geology, hydrogeology, hydrology, meteorology, oceanography, geochemistry, geography.
- **Operations research** (or Operational Research) is an interdisciplinary branch of applied mathematics and formal science that uses methods such as mathematical modeling, statistics, and algorithms to arrive at optimal or near optimal solutions to complex problems.
- **Population ecology** is a sub-field of ecology that deals with the dynamics of species populations and how these populations interact with the environment.
- **Quality control** reviews the factors involved in manufacturing and production; it can make use of statistical sampling of product items to aid decisions in process control or in accepting deliveries.
- **Quantitative psychology** is the science of statistically explaining and changing mental processes and behaviors in humans.
- **Statistical finance**, an area of econophysics, is an empirical attempt to shift finance from its normative roots to a positivist framework using exemplars from statistical physics with an emphasis on emergent or collective properties of financial markets.



- **Statistical mechanics** is the application of probability theory, which includes mathematical tools for dealing with large populations, to the field of mechanics, which is concerned with the motion of particles or objects when subjected to a force.
- **Statistical physics** is one of the fundamental theories of physics, and uses methods of probability theory in solving physical problems.

## STATISTICS AND COMPUTERS

Crunch numbers to the nth degree — and see what happens. When you study computer science and mathematics, you'll use algorithms and computational theory to create mathematical models or define formulas that solve mathematical problems. In other words, you'll design new tools that can predict the future.

The Computer Applications option gives students the flexibility to combine a traditional computer science degree with a non-traditional field. Our state-of-the-art labs for high- performance computing, networks and artificial intelligence will give you experience with the tools you'll use in the field. Through labs, lectures and projects, you'll also:

1. Investigate the computational limits of the algorithms and data structures that support complex software systems
2. Develop new applications and tools in multi-disciplinary areas of science and research
3. Explore opportunities for advanced computer modeling and simulation

## IMPORTANCE OF STATISTICS IN DIFFERENT FIELDS

Statistics plays a vital role in every fields of human activity. Statistics has important role in determining the existing position of per capita income, unemployment, population growth rate, housing, schooling medical facilities etc...in a country. Now statistics holds a central position in almost every field like Industry, Commerce, Trade, Physics, Chemistry, Economics, Mathematics, Biology, Botany, Psychology, Astronomy etc..., so application of statistics is very wide. Now we discuss some important fields in which statistics is commonly applied.

## 1. Business:

Statistics play an important role in business. A successful businessman must be very quick and accurate in decision making. He knows that what his customers wants, he should therefore, know what to produce and sell and in what quantities. Statistics helps businessman to plan production according to the taste of the costumers, the quality of the products can also be checked more efficiently by using statistical methods. So all the activities of the businessman based on statistical information. He can make correct decision about the location of business, marketing of the products, financial resources etc...

## 2. In Economics:

Statistics play an important role in economics. Economics largely depends upon statistics. National income accounts are multipurpose indicators for the economists and administrators. Statistical methods are used for preparation of these accounts. In economics research statistical methods are used for collecting and analysis the data and testing hypothesis. The relationship between supply and demands is studies by statistical methods, the imports and exports, the inflation rate, the per capita income are the problems which require good knowledge of statistics.

## 4. In Mathematics:

Statistical plays a central role in almost all natural and social sciences. The methods of natural sciences are most reliable but conclusions draw from them are only probable, because they are based on incomplete evidence. Statistical helps in describing these measurements more precisely. Statistics is branch of applied mathematics. The large number of statistical methods like probability averages, dispersions, estimation etc... is used in mathematics and different techniques of pure mathematics like integration, differentiation and algebra are used in statistics.

## 5. In Banking:

Statistics play an important role in banking. The banks make use of statistics for a number of purposes. The banks work on the principle that

all the people who deposit their money with the banks do not withdraw it at the same time. The bank earns profits out of these deposits by lending to others on interest. The bankers use statistical approaches based on probability to estimate the numbers of depositors and their claims for a certain day.

**6 In State Management (Administration):**

Statistics is essential for a country. Different policies of the government are based on statistics. Statistical data are now widely used in taking all administrative decisions. Suppose if the government wants to revise the pay scales of employees in view of an increase in the living cost, statistical methods will be used to determine the rise in the cost of living. Preparation of federal and provincial government budgets mainly depends upon statistics because it helps in estimating the expected expenditures and revenue from different sources. So statistics are the eyes of administration of the state.

**7 In Accounting and Auditing:**

Accounting is impossible without exactness. But for decision making purpose, so much precision is not essential the decision may be taken on the basis of approximation, known as statistics. The correction of the values of current assets is made on the basis of the purchasing power of money or the current value of it. In auditing sampling techniques are commonly used. An auditor determines the sample size of the book to be audited on the basis of error.

**8 In Natural and Social Sciences:**

Statistics plays a vital role in almost all the natural and social sciences. Statistical methods are commonly used for analyzing the experiments results, testing their significance in Biology, Physics, Chemistry, Mathematics, Meteorology, Research chambers of commerce, Sociology, Business, Public Administration, Communication and Information Technology etc...

**9 In Astronomy:**

Astronomy is one of the oldest branches of statistical study; it deals with the measurement of distance, sizes, masses and densities of



heavenly bodies by means of observations. During these measurements errors are unavoidable so most probable measurements are founded by using statistical methods.

## UNIT-II

### MEASURES OF CENTRAL TENDENCY

#### MEASURES OF CENTRAL TENDENCY:

The term **central tendency** refers to the "middle" value or perhaps a typical value of the data, and is measured using the **mean**, **median** or **mode**. Each of these measures is calculated differently, and the one that is best to use depends upon the situation.

In the study of a population with respect to one in which we are interested we may get a large number of observations. It is not possible to grasp any idea about the characteristic when we look at all the observations. So it is better to get one number for one group. That number must be a good representative one for all the observations to give a clear picture of that characteristic. Such representative number can be a central value for all these observations. This central value is called a measure of central tendency or an average or a measure of locations. There are five averages. Among them mean, median and mode are called simple averages and the other two averages geometric mean and harmonic mean are called special averages.

#### Arithmetic mean or mean

Arithmetic mean or simply the mean of a variable is defined as the sum of the observations divided by the number of observations. It is denoted by the symbol  $\mu$ . If the variable  $x$  assumes  $n$  values  $x_1, x_2 \dots x_n$  then the mean is given by

The arithmetic mean is the most common measure of central tendency. It is simply the sum of the numbers divided by the number of numbers. The symbol " $\mu$ " is used for the mean of a population. The symbol " $M$ " is used for the mean of a sample. The formula for  $\mu$  is shown below:

$$\mu = \Sigma X / N$$

where  $\Sigma X$  is the sum of all the numbers in the population and

$N$  is the number of numbers in the population.

The formula for  $M$  is essentially identical:

$$M = \Sigma X / N$$

where  $\Sigma X$  is the sum of all the numbers in the sample and

N is the number of numbers in the sample.

As an example, the mean of the numbers 1, 2, 3, 6, 8 is  $20/5 = 4$  regardless of whether the numbers constitute the entire population or just a sample from the population.

### Example 1

Calculate the mean for pH levels of soil 6.8, 6.6, 5.2, 5.6, 5.8

$$\bar{x} = \frac{6.8 + 6.6 + 5.2 + 5.6 + 5.8}{5} = \frac{30}{5} = 6$$

### Grouped Data

The mean for grouped data is obtained from the following formula:

$$\bar{x} = \frac{\sum f_x}{n}$$

Where  $x$  = the mid-point of individual class

$f$  = the frequency of individual class

$n$  = the sum of the frequencies or total frequencies in a sample.

### Short-cut method

$$\bar{x} = A + \frac{\sum fd}{n} \times c$$

Where  $d = \frac{x - A}{c}$

A = any value in  $x$

$n$  = total frequency

$c$  = width of the class interval

### Example 2

Given the following frequency distribution, calculate the arithmetic mean

Marks	: 64	63	62	61	60	59
Number of Students	: 8	18	12	9	7	6

### Solution

X	F	F <sub>x</sub>	D=x-A	Fd
64	8	512	2	16
63	18	1134	1	18
<b>62</b>	12	744	0	0
61	9	549	-1	-9
60	7	420	-2	-14
59	6	354	-3	-18
	60	3713		-7

### Short-cut method

Here A = 62

$$\bar{x} = 62 - \frac{7}{60} \times 1 = 61.88$$

$$\bar{x} = A + \frac{\sum fd}{n} \times c$$

### Example 3

For the frequency distribution of seed yield of plot given in table, calculate the mean yield per plot.

Yield per plot in(ing)	64.5-84.5	84.5-104.5	104.5-124.5	124.5-144.5
No of plots	3	5	7	20

### Solution

Yield ( in g)	No of Plots (f)	Mid X	$d = \frac{x - A}{c}$	Fd
64.5-84.5	3	74.5	-1	-3
84.5-104.5	5	94.5	0	0
104.5-124.5	7	114.5	1	7
124.5-144.5	20	134.5	2	40
<b>Total</b>	<b>35</b>			<b>44</b>

$$A=94.5$$

The mean yield per plot is

Direct method:

$$\begin{aligned} \bar{x} &= \frac{\sum fx}{n} = \frac{(74.5 \times 3) + (94.5 \times 5) + (114.5 \times 7) + (134.5 \times 20)}{35} \\ &= \frac{4187.5}{35} = 119.64 \text{ gms} \end{aligned}$$

### Shortcut method

$$\bar{x} = A + \frac{\sum fd}{n} \times c$$

$$\bar{x} = 94.5 + \frac{44}{35} \times 20 = 119.64 \text{ g}$$

### Merits and demerits of Arithmetic mean

#### Merits

1. It is rigidly defined.
2. It is easy to understand and easy to calculate.
3. If the number of items is sufficiently large, it is more accurate and more reliable.
4. It is a calculated value and is not based on its position in the series.
5. It is possible to calculate even if some of the details of the data are lacking.
6. Of all averages, it is affected least by fluctuations of sampling.
7. It provides a good basis for comparison.

#### Demerits

1. It cannot be obtained by inspection nor located through a frequency graph.
2. It cannot be in the study of qualitative phenomena not capable of numerical measurement i.e.  
Intelligence, beauty, honesty etc.,
3. It can ignore any single item only at the risk of losing its accuracy.
4. It is affected very much by extreme values.
5. It cannot be calculated for open-end classes.
6. It may lead to fallacious conclusions, if the details of the data from which it is computed are not given.

## Median

The median is the middle most item that divides the group into two equal parts, one part comprising all values greater, and the other, all values less than that item.

### Ungrouped or Raw data

Arrange the given values in the ascending order. If the number of values are odd, median is the middle value. If the number of values are even, median is the mean of middle two values.

By formula  $\left(\frac{n+1}{2}\right)^{th}$

When n is odd, Median = Md *value*

When n is even, Average of  $\left(\frac{n}{2}\right)$  and  $\left(\frac{n}{2} + 1\right)^{th}$  *value*

### Example 4

If the weights of sorghum ear heads are 45, 60, 48, 100, 65 gms, calculate the median

#### Solution

Here n = 5

First arrange it in ascending order

45, 48, 60, 65, 100

Median =  $\left(\frac{n+1}{2}\right)^{th}$  *value*

$$= \left(\frac{5+1}{2}\right) = 3^{rd} \text{ value} = 60$$



### Example 5

If the sorghum ear- heads are 5,48, 60, 65, 65, 100 gms, calculate the median.

### Solution

Here  $n = 6$

Median = Average of  $\left(\frac{n}{2}\right)$  and  $\left(\frac{n}{2} + 1\right)^{th}$  value

$$\left(\frac{n}{2}\right) = \frac{6}{2} = 3^{rd} \text{ value} = 60 \quad \text{and} \quad \left(\frac{n}{2} + 1\right) = \frac{6}{2} + 1 = 4^{th} \text{ value} = 65$$

$$\text{Median} = \frac{60 + 65}{2} = 62.5 \text{ g}$$

### Grouped data

In a grouped distribution, values are associated with frequencies. Grouping can be in the form of a discrete frequency distribution or a continuous frequency distribution. Whatever may

be the type of distribution, cumulative frequencies have to be calculated to know the total number of items.

### Cumulative frequency (cf)

Cumulative frequency of each class is the sum of the frequency of the class and the frequencies of the previous classes, ie adding the frequencies successively, so that the last cumulative frequency gives the total number of items.

### Discrete Series

Step1: Find cumulative frequencies.

Step2: Find  $\left(\frac{n}{2} + 1\right)$

Step3: See in the cumulative frequencies the value just greater than  $\left(\frac{n}{2} + 1\right)$

Step4: Then the corresponding value of x is median.

### Example 6

The following data pertaining to the number of insects per plant. Find median number of insects per plant.

Number of insects per plant (x)	1	2	3	4	5	6	7	8	9	10	11	12
No. of plants(f)	1	3	5	6	10	13	9	5	3	2	2	1

### Solution

Form the cumulative frequency table

x	f	cf
1	1	1
2	3	4
3	5	9
4	6	15
5	10	25
6	13	38
7	9	47
8	5	52
9	3	55
10	2	57
11	2	59
12	1	60
	60	

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

Here the number of observations is even. Therefore median = average of (n/2)th item and (n/2+1)th item.

$$= (30^{\text{th}} \text{ item} + 31^{\text{st}} \text{ item}) / 2 = (6+6)/2 = 6$$

Hence the median size is 6 insects per plant.

### Continuous Series

The steps given below are followed for the calculation of median in continuous series.

Step1: Find cumulative frequencies.

Step2: Find

Step3: See in the cumulative frequency the value first greater than  $\left(\frac{n}{2}\right)$ , Then the corresponding

class interval is called the Median class. Then apply the formula

$$\text{Median} = l + \frac{\frac{n}{2} - m}{f} \times c$$

where  $l$  = Lower limit of the median class

$m$  = cumulative frequency preceding the median class

$c$  = width of the class

$f$  = frequency in the median class.

$n$  = Total frequency.

### Example 7

For the frequency distribution of weights of sorghum ear-heads given in table below. Calculate the median.

Weights of ear heads ( in g )	No of ear heads (f)	Less than class	Cumulative frequency (m)
60-80	22	<80	22
80-100	38	<100	60
100-120	45	<120	105
120-140	35	<140	140
140-160	24	<160	164
Total	164		

### Solution

$$\text{Median} = l + \frac{\frac{n}{2} - m}{f} \times c$$

$$\left(\frac{n}{2}\right) = \left(\frac{164}{2}\right) = 82$$

It lies between 60 and 105. Corresponding to 60 the less than class is 100 and corresponding to 105 the less than class is 120. Therefore the median class is 100-120. Its lower limit is 100.

Here  $l = 100$ ,  $n=164$ ,  $f = 45$ ,  $c = 20$ ,  $m = 60$

$$\text{Median} = 100 + \frac{82 - 60}{45} \times 20 = 109.78 \text{ gms}$$

### Merits of Median

1. Median is not influenced by extreme values because it is a positional average.
2. Median can be calculated in case of distribution with open-end intervals.
3. Median can be located even if the data are incomplete.

### Demerits of Median

1. A slight change in the series may bring drastic change in median value.
2. In case of even number of items or continuous series, median is an estimated value other than any value in the series.
3. It is not suitable for further mathematical treatment except its use in calculating mean deviation.
4. It does not take into account all the observations.

### Mode

The mode refers to that value in a distribution, which occur most frequently. It is an actual value, which has the highest concentration of items in and around it. It shows the centre of concentration of the frequency in around a given value. Therefore, where the purpose is to know the point of the highest concentration it is preferred. It is, thus, a positional measure.

Its importance is very great in agriculture like to find typical height of a crop variety, maximum source of irrigation in a region, maximum disease prone paddy variety. Thus the mode is an important measure in case of qualitative data.

### Computation of the mode Ungrouped or Raw Data

For ungrouped data or a series of individual observations, mode is often found by mere inspection.

#### Example 8

Find the mode for the following seed weight

2, 7, 10, 15, 10, 17, 8, 10, 2 gms

∴ Mode = 10

In some cases the mode may be absent while in some cases there may be more than one mode.

#### Example 9

(1) 12, 10, 15, 24, 30 (no mode)

(2) 7, 10, 15, 12, 7, 14, 24, 10, 7, 20, 10

the modal values are 7 and 10 as both occur 3 times each.

## Grouped Data

For Discrete distribution, see the highest frequency and corresponding value of x is mode.

Example:

Find the mode for the following

Weight of sorghum in gms (x)	No. of ear head(f)
50	4
65	6
75	16
80	8
95	7
100	4

## Solution

The maximum frequency is 16. The corresponding x value is 75.

∴ mode = 75 gms.

## Continuous distribution

Locate the highest frequency the class corresponding to that frequency is called the modal class.

Then apply the formula.

$$\text{Mode} = l + \frac{f_s}{f_p + f_s} \times c$$

Where  $l$  = lower limit of the modal class

$f_p$  = the frequency of the class preceding the modal class

$f_s$  = the frequency of the class succeeding the modal class

and  $c$  = class interval

## Example 10

For the frequency distribution of weights of sorghum ear-heads given in table below. Calculate the mode



Weights of ear heads (g)	No of ear heads (f)	
60-80	22	
80-100	38	$f_p$
100-120	45	f
120-140	35	$f_s$
140-160	20	
<b>Total</b>	<b>160</b>	

### Solution

$$\text{Mode} = l + \frac{f_s}{f_p + f_s} \times c$$

Here  $l = 100$ ,  $f = 45$ ,  $c = 20$ ,  $m = 60$ ,  $f_p = 38$ ,  $f_s = 35$

$$\begin{aligned} \text{Mode} &= 100 + \frac{35_s}{38 + 35} \times 20 \\ &= 100 + \frac{35_s}{73} \times 20 \\ &= 109.589 \end{aligned}$$

### Geometric mean

The geometric mean of a series containing  $n$  observations is the  $n$ th root of the product of the values. If  $x_1, x_2, \dots, x_n$  are observations then

$$\begin{aligned} \text{G.M} &= \sqrt[n]{x_1, x_2, \dots, x_n} \\ &= (x_1, x_2, \dots, x_n)^{1/n} \end{aligned}$$

$$\text{Log GM} = \frac{1}{n} \log(x_1, x_2, \dots, x_n)$$

$$\begin{aligned} &= \frac{1}{n} (\log x_1 + \log x_2 + \dots + \log x_n) \\ &= \frac{\sum \log x_i}{n} \end{aligned}$$

$$GM = \text{Antilog } \frac{\sum \log x_i}{n}$$

For grouped data

$$\left[ \frac{\sum f \log x_i}{n} \right]$$

GM = Antilog

GM is used in studies like bacterial growth, cell division, etc.

### Example 11

If the weights of sorghum ear heads are 45, 60, 48, 100, 65 gms. Find the Geometric mean for the following data

Weight of ear head x (g)	Log x
45	1.653
60	1.778
48	1.681
100	2.000
65	1.813
<b>Total</b>	<b>8.925</b>

### Solution

$$\begin{aligned}
 \text{Here } n &= 5 \\
 GM &= \text{Antilog } \frac{\sum \log x_i}{n} \\
 &= \text{Antilog } \frac{8.925}{5} \\
 &= \text{Antilog } 1.785
 \end{aligned}$$

## Grouped Data

### Example 12

Find the Geometric mean for the following

Weight of sorghum (x)	No. of ear head(f)
50	4
65	6
75	16
80	8
95	7
100	4

### Solution

Weight of sorghum (x)	No. of ear head(f)	Log x	f x log x
50	5	1.699	8.495
63	10	10.799	17.99
65	5	1.813	9.065
130	15	2.114	31.71
135	15	2.130	31.95
<b>Total</b>	<b>50</b>	<b>9.555</b>	<b>99.21</b>

Here n= 50

$$\begin{aligned}
 \text{GM} &= \text{Antilog} \left[ \frac{\sum f \log x_i}{n} \right] \\
 &= \text{Antilog} \left[ \frac{99.21}{50} \right] \\
 &= \text{Antilog } 1.9842 = 96.43
 \end{aligned}$$

## Continuous distribution

### Example 13

For the frequency distribution of weights of sorghum ear-heads given in table below.

Calculate the Geometric mean

Weights of ear heads ( in g)	No of ear heads (f)
60-80	22
80-100	38
100-120	45
120-140	35
140-160	20
<b>Total</b>	<b>160</b>

### Solution

Weights of ear heads ( in g)	No of ear heads (f)	Mid x	Log x	f log x
60-80	22	70	1.845	40.59
80-100	38	90	1.954	74.25
100-120	45	110	2.041	91.85
120-140	35	130	2.114	73.99
140-160	20	150	2.176	43.52
<b>Total</b>	<b>160</b>			<b>324.2</b>

Here  $n = 160$

$$\text{GM} = \text{Antilog} \left[ \frac{\sum f \log x_i}{n} \right]$$

$$= \text{Antilog} \left[ \frac{324.2}{160} \right]$$

$$= \text{Antilog} [2.02625]$$

$$= 106.23$$

**Harmonic mean (H.M)**

Harmonic mean of a set of observations is defined as the reciprocal of the arithmetic average of the reciprocal of the given values. If  $x_1, x_2, \dots, x_n$  are  $n$  observations,

$$\text{H.M} = \frac{n}{\sum_{i=1}^n \left( \frac{1}{x_i} \right)}$$

For a frequency distribution

$$\text{H.M} = \frac{n}{\sum_{i=1}^n f \left( \frac{1}{x_i} \right)}$$

H.M is used when we are dealing with speed, rates, etc.

### Example 13

From the given data 5, 10, 17, 24, 30 calculate H.M.

X	$\frac{1}{x}$
5	0.2000
10	0.1000
17	0.0588
24	0.0417
30	0.4338

$$H.M = \frac{5}{0.4338} = 11.526$$

### Example 14

Number of tomatoes per plant are given below. Calculate the harmonic mean.

Number of tomatoes per plant	20	21	22	23	24	25
Number of plants	4	2	7	1	3	1

### Solution

Number of tomatoes per plant (x)	No of plants(f)	$\frac{1}{x}$	$f\left(\frac{1}{x}\right)$
20	4	0.0500	0.2000
21	2	0.0476	0.0952
22	7	0.0454	0.3178
23	1	0.0435	0.0435
24	3	0.0417	0.1251
25	1	0.0400	0.0400
	18		0.8216



$$\sum f \left( \frac{1}{x_i} \right) = \frac{0.1968}{21.91}$$

### Merits of H.M

1. It is rigidly defined.
2. It is defined on all observations.
3. It is amenable to further algebraic treatment.
4. It is the most suitable average when it is desired to give greater weight to smaller observations and less weight to the larger ones.

### Demerits of H.M

1. It is not easily understood.
2. It is difficult to compute.
3. It is only a summary figure and may not be the actual item in the series
4. It gives greater importance to small items and is therefore, useful only when small items have to be given greater weightage.
5. It is rarely used in grouped data.

### Percentiles

The percentile values divide the distribution into 100 parts each containing 1 percent of the cases. The  $x^{\text{th}}$  percentile is that value below which  $x$  percent of values in the distribution fall. It may be noted that the median is the  $50^{\text{th}}$  percentile.

For raw data, first arrange the  $n$  observations in increasing order. Then the  $x^{\text{th}}$  percentile is given by

$$P_x = \left( \frac{x(n+1)}{100} \right)^{\text{th}} \text{ item}$$

For a frequency distribution the  $x^{\text{th}}$  percentile is given by

$$P_x = l + \left( \frac{(x.n/100) - cf}{f} \times C \right)$$

Where

$l$  = lower limit of the percentile class which contains the  $x^{\text{th}}$  percentile value ( $x = n/100$ )

$cf$  = cumulative frequency upto  $l$

$f$  = frequency of the percentile class

$C$  = class interval

$N$  = total number of observations

### Percentile for Raw Data or Ungrouped Data

#### Example 15

The following are the paddy yields (kg/plot) from 14 plots:

30, 32, 35, 38, 40, 42, 48, 49, 52, 55, 58, 60, 62, and 65 (after arranging in ascending order). The computation of 25<sup>th</sup> percentile ( $Q_1$ ) and 75<sup>th</sup> percentile ( $Q_3$ ) are given below:

$$P_{25} \text{ (or } Q_1) = \left( \frac{25(14+1)}{100} \right)^{\text{th}} \text{ item}$$

$$= \left( 3\frac{3}{4} \right)^{\text{th}} \text{ item}$$

$$= 3^{\text{rd}} \text{ item} + (4^{\text{th}} \text{ item} - 3^{\text{rd}} \text{ item}) \left( \frac{3}{4} \right)$$

$$= 35 + (38 - 35)$$

$$= 35 + 3 \left( \frac{3}{4} \right) = 37.25 \text{ kg}$$

$$P_{75} \text{ (or } Q_3) = \left( \frac{75(14+1)}{100} \right)^{\text{th}} \text{ item}$$

$$= \left( 11\frac{1}{4} \right)^{\text{th}} \text{ item}$$

$$= 11^{\text{th}} \text{ item} + (12^{\text{th}} \text{ item} - 11^{\text{th}} \text{ item}) \left( \frac{1}{4} \right)$$

$$= 55 + (58 - 55) \left( \frac{1}{4} \right)$$

$$= 55 + 3 \left( \frac{1}{4} \right) = 55.75 \text{ kg}$$

### Example 16

The frequency distribution of weights of 190 sorghum ear-heads are given below. Compute 25<sup>th</sup> percentile and 75<sup>th</sup> percentile.

Weight of ear-heads (in g)	No of ear heads
40-60	6
60-80	28
80-100	35
100-120	55
120-140	30
140-160	15
160-180	12
180-200	9
<b>Total</b>	<b>190</b>

### Solution

Weight of ear-heads (in g)	No of ear heads	Less than class	Cumulative frequency	
40-60	6	< 60	6	
60-80	28	< 80	34	
80-100	35	<100	69	→ 47.5
100-120	55	<120	124	
120-140	30	<140	154	→ 142.5
140-160	15	<160	169	
160-180	12	<180	181	
180-200	9	<200	190	
<b>Total</b>	<b>190</b>			

or  $P_{25}$ , first find out  $\left(\frac{25(190)}{100}\right)$ , and for  $P_{75}$ ,  $\left(\frac{75(190)}{100}\right)$ , and proceed as in the case of median.

For  $P_{25}$ , we have  $\left(\frac{25(190)}{100}\right) = 47.5$ .

The value 47.5 lies between 34 and 69. Therefore, the percentile class is 80-100. Hence,

$$\begin{aligned}
 P_{25} &= Q_1 \\
 &= 80 + \left( \frac{(47.5) - 34}{35} \times 20 \right) \\
 &= l + \left( \frac{(25 \cdot n / 100) - cf}{f} \times C \right)
 \end{aligned}$$

## Quartiles

The quartiles divide the distribution in four parts. There are three quartiles. The second quartile divides the distribution into two halves and therefore is the same as the median. The first (lower).quartile ( $Q_1$ ) marks off the first one-fourth, the third (upper) quartile ( $Q_3$ ) marks off the three-fourth. It may be noted that the second quartile is the value of the median and 50<sup>th</sup> percentile.

### Raw or ungrouped data

First arrange the given data in the increasing order and use the formula for  $Q_1$  and  $Q_3$  then quartile deviation, Q.D is given by

Where  $Q_1$ . item and item

$$Q.D = \frac{Q_3 - Q_1}{\left(\frac{n+1}{4}\right)^{th}} \quad Q_3 = 3\left(\frac{n+1}{4}\right)^{th}$$

### Example 18

Compute quartiles for the data given below (grains/panicles) 25, 18, 30, 8, 15, 5, 10, 35, 40, 45

### Solution

5, 8, 10, 15, 18, 25, 30, 35, 40, 45

$$\begin{aligned} Q_1 &= \left(\frac{n+1}{4}\right)^{th} \\ &= \left(\frac{10+1}{4}\right)^{th} \\ &= (2.75)^{th} \text{ item} \\ &= 2^{nd} \text{ item} + \left(\frac{3}{4}\right)(3^{rd} \text{ item} - 2^{nd} \text{ item}) \\ &= 8 + \frac{3}{4}(10-8) \\ &= 8 + \frac{3}{4} \times 2 \end{aligned}$$

$$= 8 + 1.5$$

$$= 9.5$$

$$Q_3 = 3 \left( \frac{n+1}{4} \right)^{th}$$

$$= 3 \times (2.75)^{th} \text{ item}$$

$$= (8.75)^{th} \text{ item}$$

$$= 8^{th} \text{ item} + \left( \frac{1}{4} \right) (9^{th} \text{ item} - 8^{th} \text{ item})$$

$$= 35 + \frac{1}{4} (40 - 35)$$

$$= 35 + 1.25$$

$$= 36.25$$

### Discrete Series

Step1: Find cumulative frequencies.

Step2: Find  $\left( \frac{n+1}{4} \right)$

Step3: See in the cumulative frequencies, the value just greater than  $\left( \frac{n+1}{4} \right)$ , then the

corresponding value of  $x$  is  $Q_1$

Step4: Find  $3 \left( \frac{n+1}{4} \right)$

Step5: See in the cumulative frequencies, the value just greater than  $3 \left( \frac{n+1}{4} \right)$ , then the

corresponding value of  $x$  is  $Q_3$

### Example 19

Compute quartiles for the data given below (insects/plant).

X	5	8	12	15	19	24	30
f	4	3	2	4	5	2	4



## Solution

x	f	cf
5	4	4
8	3	7
12	2	9
15	4	13
19	5	18
24	2	20

$$Q_1 = \left( \frac{n+1}{4} \right)^{th} item = \left( \frac{24+1}{4} \right) = \left( \frac{25}{4} \right) = 6.25^{th} item$$

$$Q_3 = 3 \left( \frac{n+1}{4} \right)^{th} item = 3 \left( \frac{24+1}{4} \right) = 18.75^{th} item \therefore Q_1 = 8; Q_3 = 24$$

## Continuous series

Step1: Find cumulative frequencies

Step2: Find  $\left( \frac{n}{4} \right)$

Step3: See in the cumulative frequencies, the value just greater than  $\left( \frac{n}{4} \right)$ , then the

corresponding class interval is called first quartile class.

Step4: Find  $3 \left( \frac{n}{4} \right)$  See in the cumulative frequencies the value just greater than  $3 \left( \frac{n}{4} \right)$  then the

corresponding class interval is called 3<sup>rd</sup> quartile class. Then apply the respective formulae

$$Q_1 = l_1 + \frac{\frac{n}{4} - m_1}{f_1} \times c_1$$

$Q_3$

Where  $l_1$  = lower limit of the first quartile class

$$= l_3 + \frac{3\left(\frac{n}{4}\right) - m_3}{f_3} \times c_3$$

$f_1$  = frequency of the first quartile class

$c_1$  = width of the first quartile class

$m_1$  = c.f. preceding the first quartile class

$l_3$  = lower limit of the 3<sup>rd</sup> quartile class  $f_3$  =  
frequency of the 3<sup>rd</sup> quartile class  $c_3$  = width of  
the 3<sup>rd</sup> quartile class  
 $m_3$  = c.f. preceding the 3<sup>rd</sup> quartile class

Table 1 shows the number of touchdown (TD) passes thrown by each of the 31 teams in the National Football League in the 2000 season. The mean number of touchdown passes thrown is 20.4516 as shown below.

$$\begin{aligned}\mu &= \Sigma X/N \\ &= 634/31 \\ &= 20.4516\end{aligned}$$

Table 1. Number of touchdown passes.

37 33 33 32 29 28 28 23 22 22 22 21 21 21 20 20 19 19 18 18 18 18 16 15 14 14 14 12 12 9 6

Although the arithmetic mean is not the only "mean" (there is also a geometric mean), it is by far the most commonly used. Therefore, if the term "mean" is used without specifying whether it is the arithmetic mean, the geometric mean, or some other mean, it is assumed to refer to the arithmetic mean.

## Median

The median is also a frequently used measure of central tendency. The median is the midpoint of a distribution: the same number of scores is above the median as below it. For the data in Table 1, there are 31 scores. The 16th highest score (which equals 20) is the median because there are 15 scores below the 16th score and 15 scores above the 16th score. The median can also be thought of as the 50th percentile.

## Computation of the Median

When there is an odd number of numbers, the median is simply the middle number. For example, the median of 2, 4, and 7 is 4. When there is an even number of numbers, the median is the mean of the two middle numbers. Thus, the median of the numbers 2, 4, 7, 12 is  $(4+7)/2 = 5.5$ . When there are numbers with the same values, then the formula for the third definition of the 50th percentile should be used.

## Mode

The mode is the most frequently occurring value. For the data in Table 1, the mode is 18 since more teams (4) had 18 touchdown passes than any other number of touchdown passes. With continuous data such as response time measured to many decimals, the frequency of each value is one since no two scores will be exactly the same (see discussion of continuous variables).

Therefore the mode of continuous data is normally computed from a grouped frequency distribution. Table 2 shows a grouped frequency distribution for the target response time data. Since the interval with the highest frequency is 600-700, the mode is the middle of that interval (650).

## GEOMETRIC MEAN

Geometric Mean is a special type of average where we multiply the numbers together and then take a square root (for two numbers), cube root (for three numbers) etc.

**Example: What is the Geometric Mean of 2 and 18?**

- First we multiply them:  $2 \times 18 = 36$
- Then (as there are two numbers) take the square root:  $\sqrt{36} = 6$

In one line:

$$\text{Geometric Mean of 2 and 18} = \sqrt{(2 \times 18)} = 6$$

It is like the area is the same!

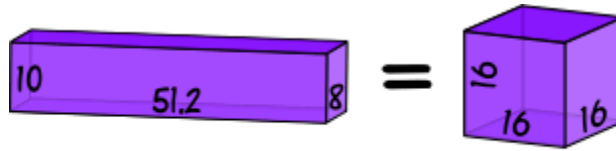
**Example: What is the Geometric Mean of 10, 51.2 and 8?**

- First we multiply them:  $10 \times 51.2 \times 8 = 4096$
- Then (as there are three numbers) take the cube root:  $\sqrt[3]{4096} = 16$

In one line:

$$\text{Geometric Mean} = \sqrt[3]{(10 \times 51.2 \times 8)} = 16$$

It is like the volume is the same:



$$10 \times 51.2 \times 8 = 16 \times 16 \times 16$$

**Example:** What is the Geometric Mean of 1,3,9,27 and 81?

- First we multiply them:  $1 \times 3 \times 9 \times 27 \times 81 = 59049$
- Then (as there are 5 numbers) take the 5th root:  $\sqrt[5]{59049} = 9$

In one line:

$$\text{Geometric Mean} = \sqrt[5]{(1 \times 3 \times 9 \times 27 \times 81)} = 9$$

I can't show you a nice picture of this, but it is still true that:

$$1 \times 3 \times 9 \times 27 \times 81 = 9 \times 9 \times 9 \times 9 \times 9$$

## Harmonic Mean

A kind of [average](#). To find the harmonic mean of a [set](#) of  $n$  numbers, add the [reciprocals](#) of the numbers in the set, divide the [sum](#) by  $n$ , then take the reciprocal of the result. The harmonic mean of  $\{a_1, a_2, a_3, a_4, \dots, a_n\}$  is given below.

Formula: Harmonic Mean = 
$$\frac{n}{\frac{1}{a_1} + \frac{1}{a_2} + \frac{1}{a_3} + \frac{1}{a_4} + \dots + \frac{1}{a_n}}$$

Example: For the numbers 4 and 9,  
Harmonic Mean = 
$$\frac{2}{\frac{1}{4} + \frac{1}{9}} = \frac{72}{13} = 5.54$$

## RANGE

The difference between the lowest and highest values.

In {4, 6, 9, 3, 7} the lowest value is 3, and the highest is 9, so the range is  $9 - 3 = 6$ .

Range can also mean all the output values of a function.

## Quartile Deviation :

In a distribution, partial variance between the upper quartile and lower quartile is known as 'quartile deviation'. Quartile Deviation is often regarded as semi inter quartile range.

**Formula :** (Upper quartile- lower quartile) / 2

upper quartile = 400, lower quartile = 200 then

Quartile deviation (QD) =  $(400-200)/2 = 200/2$

=100.

## Mean Deviation

The mean of the distances of each value from their mean.

Yes, we use "**mean**" twice: Find the mean ... use it to work out distances ... then find the mean of those distances!

Three steps:

- 1. Find the mean of all values
- 2. Find the **distance** of each value from that mean (subtract the mean from each value, ignore minus signs)
- 3. Then find the **mean of those distances** Example: the Mean Deviation of 3, 6, 6, 7, 8, 11, 15, 16

Step 1: Find the **mean**:

$$\text{Mean} = \frac{3 + 6 + 6 + 7 + 8 + 11 + 15 + 16}{8} = \frac{72}{8} = 9$$

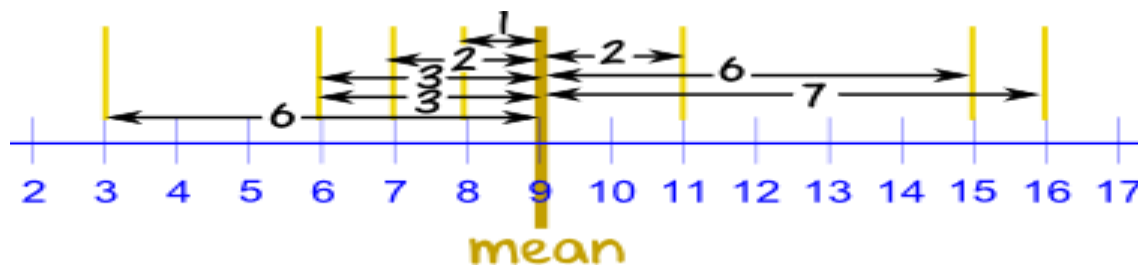


Step 2: Find the **distance** of each value from that mean:

Value	Distance from 9
3	6
6	3
6	3
7	2
8	1
11	2
15	6
16	7

Step 3. Find the **mean of those distances**:

$$6 + 3 + 3 + 2 + 1 + 2 + 6 + 7 = 30$$



$$\text{Mean Deviation} = \frac{30}{8} = 3.75$$

So, the **mean = 9**, and the **mean deviation = 3.75**

It tells us how far, on average, all values are from the middle.

In that example the values are, on average, 3.75 away from the middle.

For **deviation** just think **distance** Formula

The formula is:

$$\text{Mean Deviation} = \frac{\sum |x - \mu|}{N}$$

Let's learn more about those symbols!

Firstly:

- $\mu$  is the mean (in our example  $\mu = 9$ )
- $x$  is each value (such as 3 or 16)
- $N$  is the number of values (in our example  $N = 8$ )

### Standard Deviation

The Standard Deviation is a measure of how spread out numbers are.

Its symbol is  $\sigma$  (the greek letter sigma)

The formula is easy: it is the **square root** of the **Variance**. So now you ask, "What is the Variance?"

### Variance:

The Variance is defined as. The average of the squared differences from the Mean.

### Examples of Standard Deviation:

This tutorial is about some examples of standard deviation using all methods which are discussed in the previous tutorial.

### Example:

Calculate the standard deviation for the following sample data using all methods: 2, 4, 8, 6, 10, and 12.

### Method-I: Actual Mean Method

Marks	f	X	$F(X)$	$(X-X^-)^2$	$f(X-X^-)^2$
1-3	40	2	80	4	160
3-5	30	4	120	0	0
5-7	20	6	120	4	80
7-9	10	8	80	16	160
<b>Total</b>	100		400		400

$$X^- = \sum f(X) / \sum f = 400 / 100 = 4$$

### Method-II: Taking assumed mean as 2

Marks	f	X	$D=(X-2)$	$fD$	$fD^2$
1-3	40	2	0	0	0
3-5	30	4	2	60	120
5-7	20	6	4	80	320
7-9	10	8	6	60	160

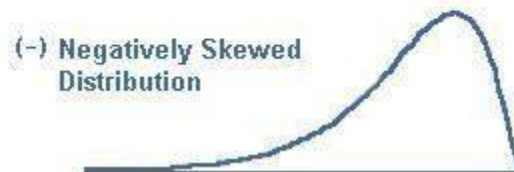
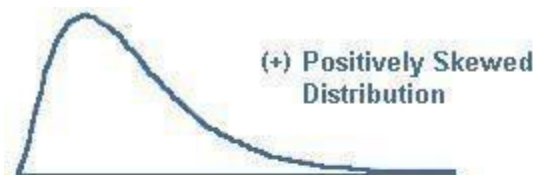
<b>Total</b>	100			200	800
--------------	-----	--	--	-----	-----

## SKEWNESS

Lack of symmetry is called Skewness. If a distribution is not symmetrical then it is called skewed distribution. So, mean, median and mode are different in values and one tail becomes longer than other. The skewness may be positive or negative.

### Positively skewed distribution:

If the frequency curve has longer tail to right the distribution is known as positively skewed distribution and  $Mean > Median > Mode$ .



### Negatively skewed distribution:

If the frequency curve has longer tail to left the distribution is known as negatively skewed distribution and  $Mean < Median < Mode$ .

### Measure of Skewness:

The difference between the mean and mode gives as absolute measure of skewness. If we divide this difference by standard deviation we obtain a relative measure of skewness known as coefficient and denoted by  $SK$ .

Karl Pearson coefficient of Skewness

$$SK = \frac{Mean - Mode}{S.D}$$

Sometimes the mode is difficult to find. So we use another formula

$$SK = 3(\text{Mean} - \text{Median}) / S.D$$

Bowley's coefficient of Skewness

$$SK = (Q_1 + Q_3 - 2\text{Median}) / (Q_3 - Q_1)$$

Kelly's Measure of Skewness is one of several ways to measure skewness in a data distribution. Bowley's skewness is based on the middle 50 percent of the observations in a data set. It leaves 25 percent of the observations in each tail of the distribution. Kelly suggested that leaving out fifty percent of data to calculate skewness was too extreme. He created a measure to find skewness with more data. Kelly's measure is based on  $P_{90}$  (the 90th percentile) and  $P_{10}$  (the 10th percentile). Only twenty percent of observations (ten percent in each tail) are excluded from the measure.

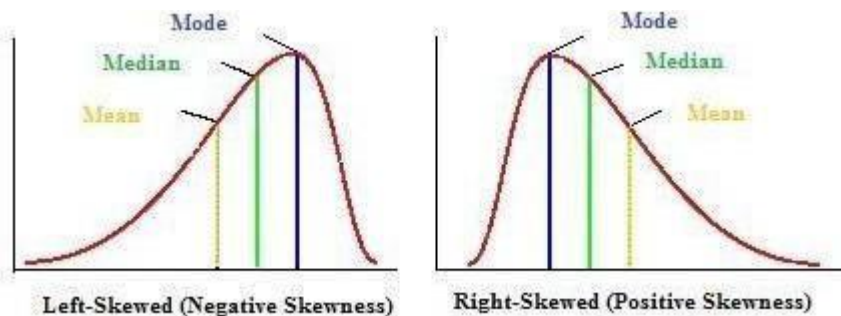
### Kelly's Measure Formula.

Kelley's measure of skewness is given in terms of percentiles and [deciles](#) (D). Kelley's absolute measure of skewness ( $S_k$ ) is:

$$S_k = P_{90} + P_{10} - 2 * P_{50} = D_9 + D_1 - 2 * D_5.$$

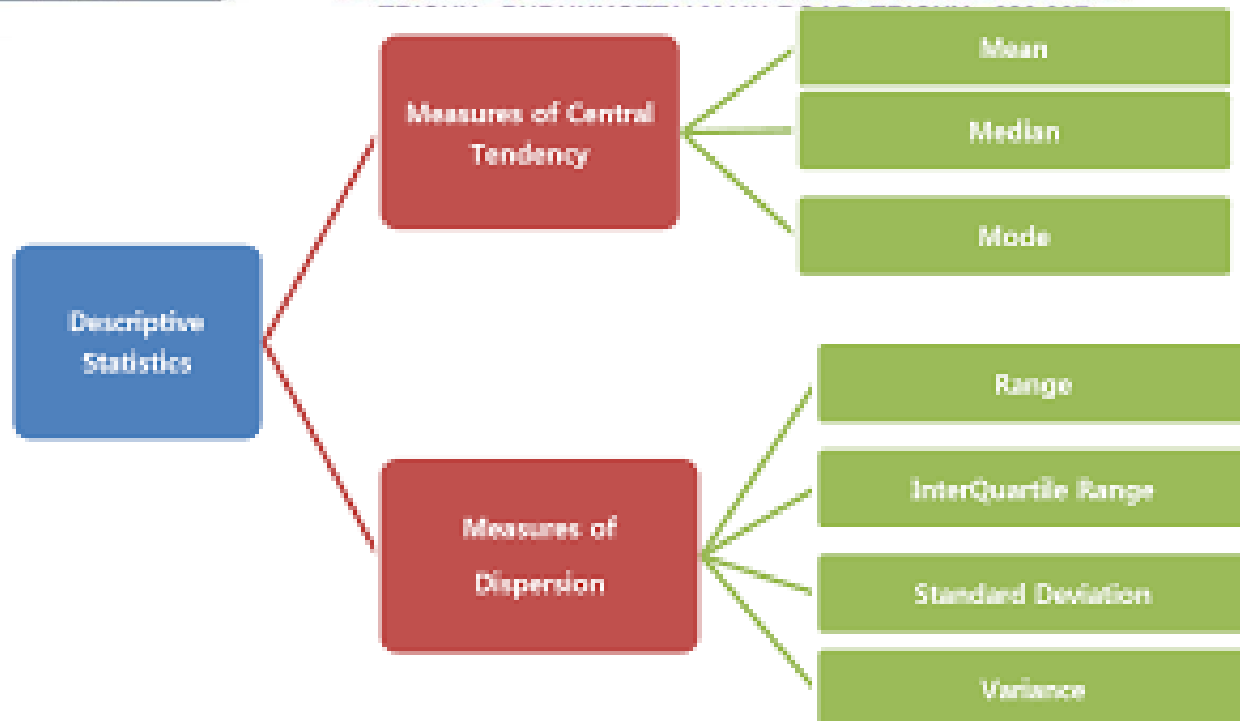
Kelly's Measure of Skewness gives you the same information about skewness as the other three types of skewness measures.

A measure of skewness = 0 means that the distribution is symmetrical.



A measure of skewness > 0 means a positive skewness.

A measure of skewness < means a negative skewness.



### UNIT –III

#### TABULATION OF UNIVARIATE

**Univariate data:** Univariate means "one variable" (one type of data): Example: Travel Time (minutes): 15, 29, 8, 42, 35, 21, 18, 42,

26 The variable is **Travel Time**

**Bivariate data:** Bivariate means "two variables", in other words there are two types of data with bivariate data you have two sets of related data that you want to compare: Example:

An ice cream shop keeps track of how much ice cream they sell versus the temperature on that day.

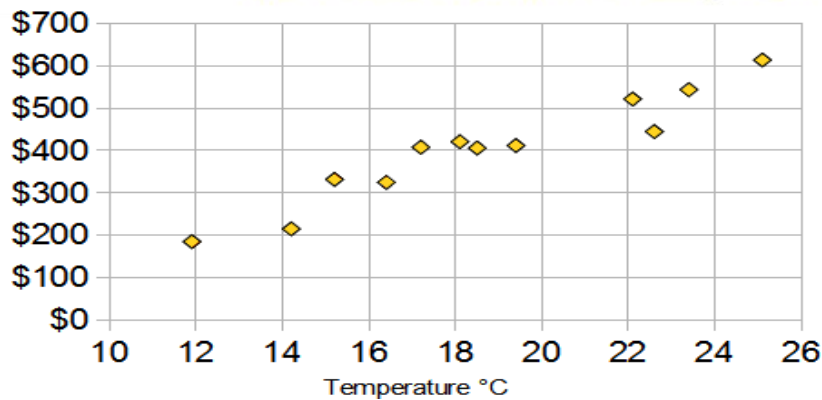
The two variables are **Ice Cream Sales** and **Temperature**.

Here are their figures for the last 12 days:

<i>Ice Cream Sales vs Temperature</i>	
Temperature °C	Ice Cream Sales
14.2°	\$215
16.4°	\$325
11.9°	\$185
15.2°	\$332
18.5°	\$406
22.1°	\$522
19.4°	\$412
25.1°	\$614
23.4°	\$544
18.1°	\$421
22.6°	\$445
17.2°	\$408

And here is the same data as a [Scatter Plot](#):





Now we can easily see that **warmer weather** and **more ice cream sales** are linked, but the relationship is not perfect.

### Multivariate data

**Multivariate Data** Analysis refers to any statistical technique used to analyze **data** that arises from more than one variable. This essentially models reality where each situation, product, or decision involves more than a single variable.

Univariate Data		Bivariate Data	
<ul style="list-style-type: none"> <li>involving a <b>single variable</b></li> <li>does not deal with causes or relationships</li> <li>the major purpose of univariate analysis is to describe</li> <li>central tendency - mean, mode, median</li> <li>dispersion - range, variance, max, min, quartiles, standard deviation.</li> <li>frequency distributions</li> <li>bar graph, histogram, pie chart, line graph, box-and-whisker plot</li> </ul>		<ul style="list-style-type: none"> <li>involving <b>two variables</b></li> <li>deals with causes or relationships</li> <li>the major purpose of bivariate analysis is to explain</li> <li>analysis of two variables simultaneously</li> <li>correlations</li> <li>comparisons, relationships, causes, explanations</li> <li>tables where one variable is contingent on the values of the other variable.</li> <li>independent and dependent variables</li> </ul>	
<b>Sample question:</b> Is there a relationship between the number of females in Computer Programming and their scores in Mathematics?		<b>Sample question:</b> Is there a relationship between the number of females in Computer Programming and their scores in Mathematics?	

## DIAGRAMATIC AND GRAPHICAL; REPRESENTATION OF DATA

Although tabulation is very good technique to present the data, but diagrams are an advanced technique to represent data. As a layman, one cannot understand the tabulated data easily but with only a single glance at the diagram, one gets complete picture of the data presented. According to M.J. Moroney, –diagrams register a meaningful impression almost before we think.

### Importance or utility of Diagrams

- Diagrams give a very clear picture of data. Even a layman can understand it very easily and in a short time.
- We can make comparison between different samples very easily. We don't have to use any statistical technique further to compare.
- This technique can be used universally at any place and at any time. This technique is used almost in all the subjects and other various fields.
- Diagrams have impressive value also. Tabulated data has not much impression as

compared to Diagrams. A common man is impressed easily by good diagrams.

- This technique can be used for numerical type of statistical analysis, e.g. to locate Mean, Mode, Median or other statistical values.
- It does not save only time and energy but also is economical. Not much money is needed to prepare even good diagrams.
- These give us much more information as compared to tabulation. Technique of tabulation has its own limits.
- This data is easily remembered. Diagrams which we see leave their lasting impression much more than other data techniques.
- Data can be condensed with diagrams. A simple diagram can present what even cannot be presented by 10000 words.

### **General Guidelines for Diagrammatic presentation**

- The diagram should be properly drawn at the outset. The pith and substance of the subject matter must be made clear under a broad heading which properly conveys the purpose of a diagram.
- The size of the scale should neither be too big nor too small. If it is too big, it may look ugly. If it is too small, it may not convey the meaning. In each diagram, the size of the paper must be taken note-of. It will help to determine the size of the diagram.
- For clarifying certain ambiguities some notes should be added at the foot of the diagram. This shall provide the visual insight of the diagram.
- Diagrams should be absolutely neat and clean. There should be no vagueness or overwriting on the diagram.
- Simplicity refers to love at first sight. It means that the diagram should convey the meaning clearly and easily.
- Scale must be presented along with the diagram.
- It must be Self-Explanatory. It must indicate nature, place and source of data presented.
- Different shades, colors can be used to make diagrams more easily understandable.
- Vertical diagram should be preferred to Horizontal diagrams.
- It must be accurate. Accuracy must not be done away with to make it attractive or impressive.

### **Limitations of Diagrammatic Presentation**

- Diagrams do not present the small differences properly.

- These can easily be misused.
- Only artist can draw multi-dimensional diagrams.
- In statistical analysis, diagrams are of no use.
- Diagrams are just supplement to tabulation.
- Only a limited set of data can be presented in the form of diagram.
- Diagrammatic presentation of data is a more time consuming process.
- Diagrams present preliminary conclusions.
- Diagrammatic presentation of data shows only an estimate of the actual behavior of the variables.

## Types of Diagrams

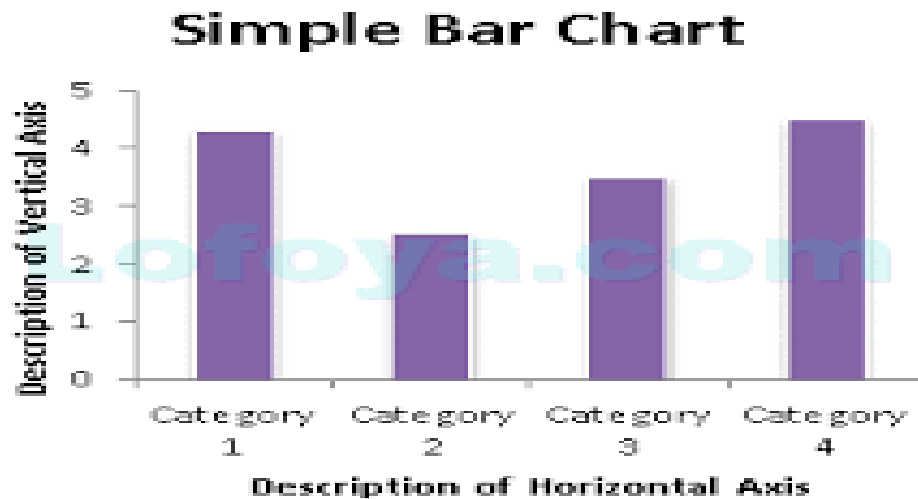
### (a) Line Diagrams



In these diagrams only line is drawn to represent one variable. These lines may be vertical or horizontal. The lines are drawn such that their length is the proportion to value of the terms or items so that comparison may be done easily.

## (b) Simple Bar Diagram

Like line diagrams these figures are also used where only single dimension i.e. length can present the data. Procedure is almost the same, only one thickness of lines is measured. These can also be drawn either vertically or horizontally. Breadth of these lines or bars should be equal. Similarly distance between these bars should be equal. The breadth and distance between them



should be taken according to space available on the paper.

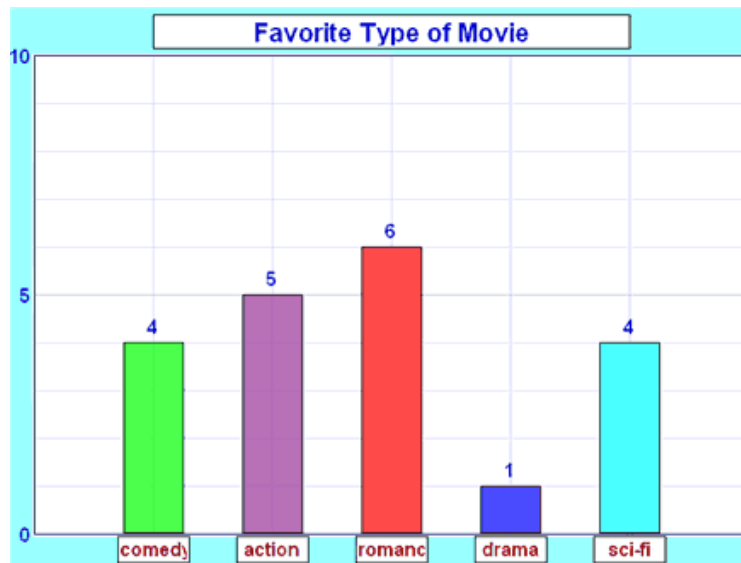
Imagine you just did a survey of your friends to find which kind of movie they liked best:

*Table: Favorite Type of Movie*

Comedy	Action	Romance	Drama	SciFi
4	5	6	1	4

We can show that on a bar graph like this:

It is a really good way to show relative sizes: we can see which types of movie are most liked, and which are least liked, at a glance.

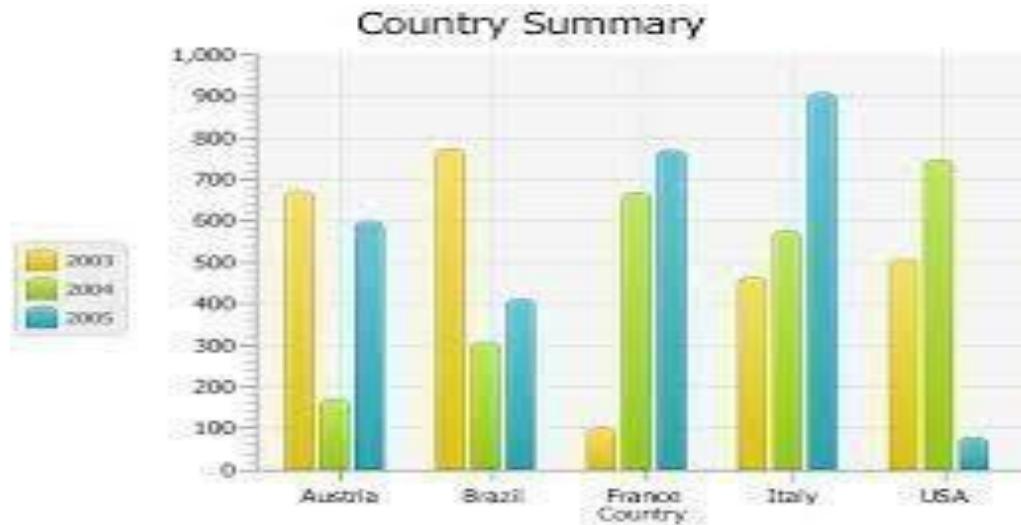


We can use bar graphs to show the relative sizes of many things, such as what type of car people have, how many customers a shop has on different days and so on.

### (c) Multiple Bar Diagrams

The diagram is used, when we have to make comparison between more than two variables. The number of variables may be 2, 3 or 4 or more. In case of 2 variables, pair of bars is drawn. Similarly, in case of 3 variables, we draw triple bars. The bars are drawn on the same proportionate basis as in case of simple bars. The same shade is given to the same item.





Draw a multiple bar chart to represent the import and export of Canada (values in \$) for the years 1991 to 1995.

Years	Imports	Exports
1991	7930	4260
1992	8850	5225
1993	9780	6150
1994	11720	7340
1995	12150	8145

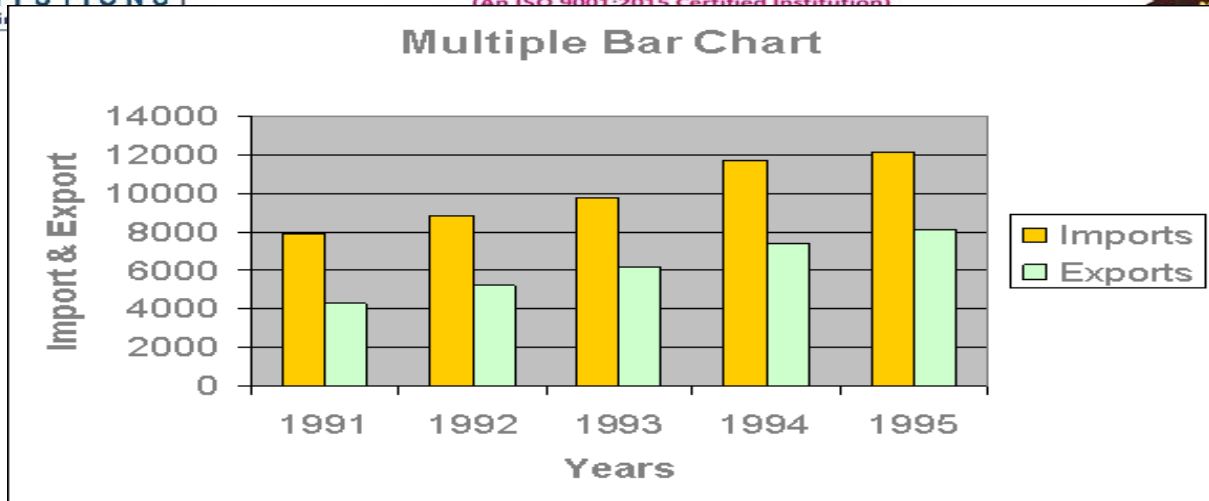
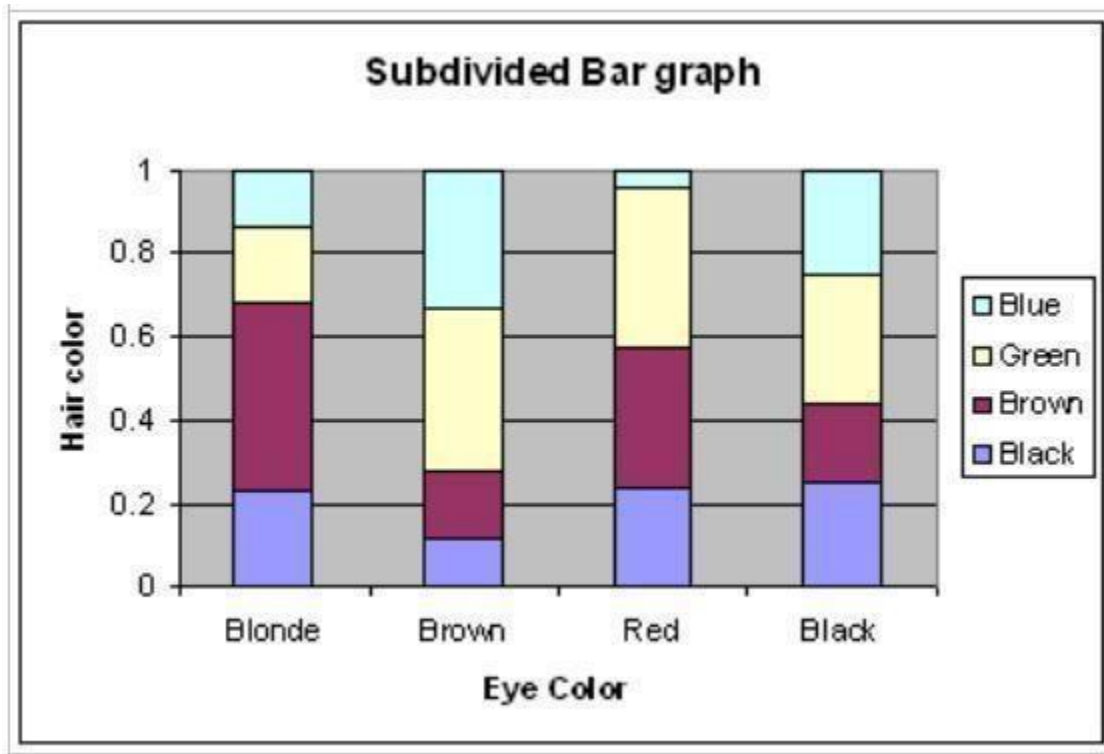


chart showing the import and export of Canada from 1991 – 1995.



#### (d) Sub-divided Bar Diagram

The data which is presented by multiple bar diagram can be presented by this diagram. In this case we add different variables for a period and draw it on a single bar as shown in the following examples. The components must be kept in same order in each bar. This diagram is more efficient if number of components is less i.e. 3 to 5.

#### (e) Percentage Bar Diagram

Like sub-divide bar diagram, in this case also data of one particular period or variable is put on single bar, but in terms of percentages. Components are kept in the same order in each bar for Easy comparison.

#### (f) Duo-directional Bar Diagram

In this case the diagram is on both the sides of base line i.e. to left and right or to above or below sides.

#### (g) Broken Bar Diagram

This diagram is used when value of some variable is very high or low as compared to others. In this case the bars with bigger terms or items may be shown broken.

### One dimensional diagrams:

A diagram in which the size of only one dimension i.e. length is fixed in proportion to the value of the data is called one dimensional diagram. Such diagrams are also popularly called bar diagrams. These diagrams can be drawn in both vertical and horizontal manner. The related different bar diagrams differ from each other only in respect of their length dimension, while they remain the same in respect of their other two dimensions i.e. breadth and thickness. The size of breadth of each of such diagrams is determined taking into consideration the number of diagrams to be drawn, and the size of the paper at one end. They may take the form of a line, or a thread if they are to be drawn in large numbers on the surface of a paper. However, their breadth should neither be too large nor too small for that in both the cases they look ugly. The dimension of thickness does not look prominent in such diagrams. The examples of such diagrams are given on the next page.

Techniques of drawing bar diagrams

The following are the techniques of drawing the bar diagrams :

- (i) First, draw the base line, preferably horizontally, and divide it into a number of equal parts keeping in view the number of diagrams to be drawn.
- (ii) Then, draw the scale line preferably vertically, and divide into a number of equal parts keeping in view the maximum value to be represented.
- (iii) Then, fix the width of the bar uniformly keeping in view the number of bars to be drawn and the gaps to be provided in between each two of them.
- (iv) Then, fix the size of gaps to be provided between each of the two bars uniformly.
- (v) Then, fix the lengths of the different bars in proportion of the value of the data.
- (vi) Then, draw the different bars in accordance with their length, and width thus fixed, and arranged in order of their length, or time of occurrence.
- (vii) Then, decorate the bars with similar or different colours, or shades according to the similarity or dissimilarity in the nature of the data respectively.
- (viii) Give a description of the data in short at the bottom of respective bars.
- (xiv) Put the respective figures at the top of each bar to read out the exact value at a glance without looking at the scale.

## Advantages

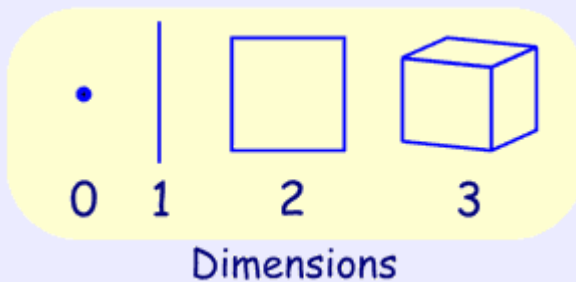
The chief advantages of a bar diagram can be outlined as under:

1. It is very simple to draw and read as well.
2. It is the only form of diagram which can represent a large number of data on a piece of paper.
3. It can be drawn both vertically and horizontally.
4. It gives a better look and facilitates comparison.

## Disadvantages

1. It cannot exhibit a large number of aspects of the data.
2. The which of the bars are fixed arbitrarily by a drawer.

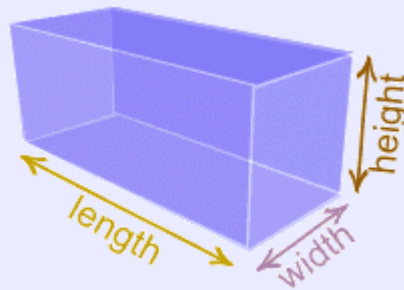
## Two-Dimensional



A shape that only has two dimensions (such as width and height) and no thickness. Squares, Circles, Triangles, etc are two dimensional objects.

Also known as "2D".

### 3. Three-Dimensional



An object that has height, width and depth, like any object in the real world. Example: your body is three-dimensional.

Also known as "3D".

**Pie Chart:** a special chart that uses "pie slices" to show relative sizes of data.

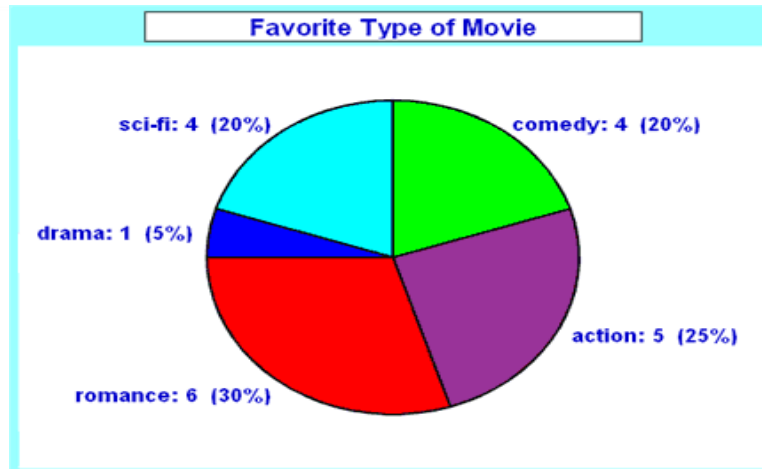
Imagine you survey your friends to find the kind of movie they like best:

*Table: Favorite Type of Movie*

Comedy	Action	Romance	Drama	SciFi
4	5	6	1	4

You can show the data by this Pie Chart:





It is a really good way to show relative sizes: it is easy to see which movie types are most liked, and which are least liked, at a glance.

## UNIT-IV

statistical Test for Population Mean (Small Sample)

In this section will adjust our statistical test for the population mean to apply to small sample situations. Fortunately (sic!), this will be easy (in fact, once you understand *one* statistical test, additional tests are easy since they all follow a similar procedure.

The only difference in performing a "small sample" statistical test for the mean as opposed to a "large sample test" is that we do not use the normal distribution as prescribed by the Central Limit theorem, but instead a more conservative distribution called the **T-Distribution**. The Central Limit theorem applies best when sample sizes are large so that we need to make some adjustment in computing probabilities for small sample sizes. The appropriate function in Excel is the TDIST function, defined as follows:

TDIST(T, N-1, TAILS), where

- T is the value for which we want to compute the probability
- N is the sample size (and N-1 is frequently called the "degrees of freedom")
- **TAILS** is either 1 (for a 1-tail test) or 2 (for a 2-tail test). Since we again consider 2-tailed tests only we always use 2

With that new Excel function our test procedure for a sample mean, small sample size, is as follows:

Statistical Test for the Mean (small sample size  $N < 30$ ):

Fix an error level you are comfortable with (something like 10%, 5%, or 1% is most common). Denote that "comfortable error level" by the letter "A".

Then setup the test as follows:

Null Hypothesis  $H_0$ :

mean = M, i.e. The mean is a known number M Alternative Hypothesis  $H_a$ :

mean  $\neq$  M, i.e. mean is different from M (*2-tailed test*)

Test Statistics:

Select a random sample of size N, compute its sample mean X and the standard deviation S. Then compute the corresponding t-score as follows:

$$T = (X - M) / (S / \sqrt{N})$$

Rejection Region (Conclusion)

$$\text{Compute } p = 2 * (1 - P(t > |T|)) = \text{TDIST}(\text{ABS}(T), N-1, 2)$$

If the probability p computed in the above step is less than A (the error level you were comfortable with initially, you reject the null hypothesis  $H_0$  and accept the alternative hypothesis. Otherwise you declare your test inconclusive.

Comments:

- The null and alternative hypothesis for this test are the same as before
- The calculation of the test statistics is the same as before, but the result is called T instead of Z (oh well -:-)
- The TDIST function is similar to the NORMSDIST function, but it does not work for negative values of T (a limitation of Excel), and it automatically gives a "tail" probability. Thus, the computation of the p-value had to be adjusted accordingly.
- The ABS function in the above formulas stands for the "absolute value" function. (In other words, just drop any minus signs ... -:-)

**Example 1:** A group of secondary education student teachers were given 2 1/2 days of training in interpersonal communication group work. The effect of such a training session on the dogmatic nature of the student teachers was measured by the difference of scores on the "Rokeach Dogmatism test" given before and after the training session. The difference "post minus pre score" was recorded as follows:

-16, -5, 4, 19, -40, -16, -29, 15, -2, 0, 5, -23, -3, 16, -8, 9, -14, -33, -64, -33

Can we conclude from this evidence that the training session makes student teachers less dogmatic (at the 5% level of significance)?

This is of course the same example as before, where we *incorrectly* used the normal distribution

to compute the probability in the last step. This time, we will do it correctly, which is fortunately almost identical to the previous case (except that we use TDIST instead of NORMDIST):

- **Null Hypothesis:** there is no difference in dogmatism, i.e. mean = 0
- **Alternative Hypothesis:** dogmatism is different, i.e. mean not equal to 0
- **Test statistics:** sample mean = -10.9, standard deviation = 21.33, sample size = 20.

Compute

$$T = (-10.9 - 0) / (21.33 / \sqrt{20}) = -2.28$$

- **Rejection Region:** We use Excel to compute  $p = \text{TDIST}(2.28, 19, 2) = 0.034$ , or 3.4%.

That probability is less than 0.05 so we reject the null hypothesis.

Note that in the previous section we (incorrectly) computed the probability  $p$  to be 2.2%, now it is 3.4%. The difference is small, but can be significant in special situations. Thus, to be safe:

- if  $N > 30$  use the Z-Test based on the standard normal distribution NORMSDIST as in the previous section
- if  $N < 30$  use the T-Test based on the T-Distribution TDIST as in this section

**Example 2:** Suppose GAP, the clothing store, wants to introduce their line of clothing for women to another country. But their clothing sizes are based on the assumption that the average size of a woman is 162 cm. To determine whether they can simply ship the clothes to the new country they select 5 women at random in the target country and determine their heights as follows: 149, 165, 150, 158, 153

should they adjust their line of clothing or they ship them without change? Make sure to decide at the 0.05-level. By now statistical testing is second-nature (I hope -:)

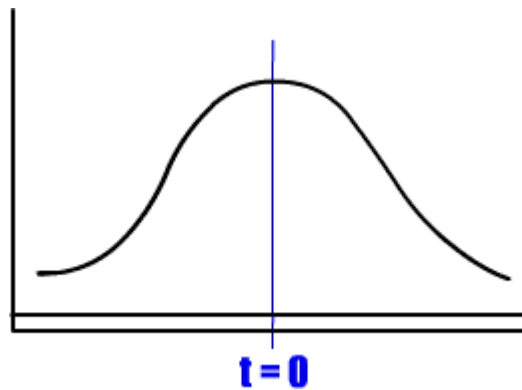
- **Null Hypothesis:** mean height in new country is the same as in old country, i.e.  $M = 162$
- **Alt. Hypothesis:** mean height in new country is different from old country, i.e.  $M$  not equal to 162 (*either too small or too tall would be bad for GAP*)
- **Test Statistics:** we can compute the sample mean = 155 and the sample standard deviation = 6.59 while the sample size is clearly  $N = 5$ .

- **Therefore:**

$$T = (155 - 162) / (6.59 / \sqrt{5}) = -2.37$$

- **Rejection Region:** We use Excel to compute  $p = \text{TDIST}(2.37, 4, 2) = 0.077$ . Thus, if we did decide to reject the null hypothesis the probability of that decision being wrong is 7.7%. That is larger than 0.05, thus we declare the test inconclusive.

#### t-Distribution.



#### ● Properties of the t-Distribution

If a population is essentially normal, then the distribution is:

$$t = \frac{\bar{X} - \mu}{S / \sqrt{n}}$$

This is the equation for the Student t-Distribution, or simply t-Distribution, for all samples of size  $n$  less than 30.

To find the Rejection Region, we can use the t-Distribution Table. This table merely requires knowledge of the sample size, which allows us to calculate the Degrees of Freedom, and the Significance Level.

As illustrated above, the t-distribution has many properties which differentiate it from the standard normal or z-distribution.

The distribution shares the bell curve of the z, but reflects the variability that is inherent with smaller sample sizes.

The shape of the t-distribution is dependent on the sample size  $n$ .

The standard deviation is greater than 1.

As the sample size  $n$  increases, the shape of the curve approaches the standard deviation.

### PAIRED T TEST

- Paired sample t-test is a statistical technique that is used to compare two population means in the case of two samples that are correlated. Paired sample t-test is used in before-after studies, or when the samples are the matched pairs, or when it is a case-control study. For example, if we give training to a company employee and we want to know whether or not the training had any impact on the efficiency of the employee, we could use the paired sample test. We collect data from the employee on a seven scale rating, before the training and after the training. By using the paired sample t-test, we can statistically conclude whether or not training has improved the efficiency of the employee. In medicine, by using the paired sample t-test, we can figure out whether or not a particular medicine will cure the illness.

#### Steps:

- Set up hypothesis:** We set up two hypotheses. The first is the null hypothesis, which assumes that the mean of two paired samples are equal. The second hypothesis will be an alternative hypothesis, which assumes that the means of two paired samples are not equal.
- Select the level of significance:** After making the hypothesis, we choose the level of significance. In most of the cases, significance level is 5%, (in medicine, the significance level is set at 1%).
- Calculate the parameter:** To calculate the parameter we will use the following



formula:

$$t = \frac{\bar{d}}{\sqrt{s^2/n}}$$

Where  $\bar{d}$  is the mean difference between two samples,  $s^2$  is the sample variance,  $n$  is the [sample size](#) and  $t$  is a paired sample t-test with  $n-1$  degrees of freedom. An alternate formula for paired sample t-test is:

$$t = \frac{\bar{d}}{\sqrt{\frac{n(\sum d^2) - (\sum d)^2}{n-1}}}$$

**4. Testing of hypothesis or decision making:** After calculating the parameter, we will compare the calculated value with the table value. If the calculated value is greater than the table value, then we will reject the null hypothesis for the paired sample t-test. If the calculated value is less than the table value, then we will accept the null hypothesis and say that there is no significant mean difference between the two paired samples.

#### Assumptions:

1. Only the matched pairs can be used to perform the test.
2. Normal distributions are assumed.
3. The variance of two samples is equal.
4. Cases must be independent of each other.

#### Administration, Analysis and Reporting

Statistics Solutions consists of a team of professional methodologists and statisticians that can assist the student or professional researcher in administering the survey instrument, collecting the data, conducting the analyses and explaining the results.

Let  $x$  = the difference in weight 3 months after the program starts. The null hypothesis is:

$H_0: \mu = 0$ ; i.e. any differences in weight is due to chance

We can make the following calculations using the difference column D:

$$\text{s.e.} = \text{std dev} / \sqrt{n} = 6.33 / \sqrt{15} = 1.6343534$$

$$t_{obs} = (\bar{x} - \mu) / \text{s.e.} = (10.93 - 0) / 1.63 = 6.6896995$$

$$t_{crit} = \text{TINV}(\alpha, df) = \text{TINV}(.05, 14) = 2.1447867$$

Since  $t_{obs} > t_{crit}$  we reject the null hypothesis and conclude with 95% confidence that the difference in weight before and after the program is not due solely to chance.

Alternatively we can use a type 1 TTEST to perform the analysis as follows:

$$\text{p-value} = \text{TTEST}(B4:B18, C4:C18, 2, 1) = 1.028\text{E-}05 < .05 = \alpha$$

and so once again we reject the null hypothesis.

## Chi-Square Test for Independence

This lesson explains how to conduct a **chi-square test for independence**. The test is applied when you have two [categorical variables](#) from a single population. It is used to determine whether there is a significant association between the two variables.

For example, in an election survey, voters might be classified by gender (male or female) and voting preference (Democrat, Republican, or Independent). We could use a chi-square test for independence to determine whether gender is related to voting preference. The [sample problem](#) at the end of the lesson considers this example.

## When to Use Chi-Square Test for Independence

The test procedure described in this lesson is appropriate when the following conditions are met:

- The sampling method is simple random sampling.
- The variables under study are each categorical.
- If sample data are displayed in a contingency table, the expected frequency count for each cell of the table is at least 5.

This approach consists of four steps: (1) state the hypotheses, (2) formulate an analysis plan, (3) analyze sample data, and (4) interpret results.

### **State the Hypotheses**

Suppose that Variable A has  $r$  levels, and Variable B has  $c$  levels. The null hypothesis states that knowing the level of Variable A does not help you predict the level of Variable B. That is, the variables are independent.

$H_0$ : Variable A and Variable B are independent.

$H_a$ : Variable A and Variable B are not independent.

The alternative hypothesis is that knowing the level of Variable A *can* help you predict the level of Variable B.

**Note:** Support for the alternative hypothesis suggests that the variables are related; but the relationship is not necessarily causal, in the sense that one variable "causes" the other.

### **Formulate an Analysis Plan**

The analysis plan describes how to use sample data to accept or reject the null hypothesis. The plan should specify the following elements.

- Significance level. Often, researchers choose significance levels equal to 0.01, 0.05, or 0.10; but any value between 0 and 1 can be used.
- Test method. Use the chi-square test for independence to determine whether there is a significant relationship between two categorical variables.

### Analyze Sample Data

Using sample data, find the degrees of freedom, expected frequencies, test statistic, and the P- value associated with the test statistic. The approach described in this section is illustrated in the sample problem at the end of this lesson.

- **Degrees of freedom.** The degrees of freedom (DF) is equal to:

$$DF = (r - 1) * (c - 1)$$

where  $r$  is the number of levels for one categorical variable, and  $c$  is the number of levels for the other categorical variable.

- **Expected frequencies.** The expected frequency counts are computed separately for each level of one categorical variable at each level of the other categorical variable. Compute  $r * c$  expected frequencies, according to the following formula.

$$E_{r,c} = (n_r * n_c) / n$$

where  $E_{r,c}$  is the expected frequency count for level  $r$  of Variable A and level  $c$  of Variable B,  $n_r$  is the total number of sample observations at level  $r$  of Variable A,  $n_c$  is the total number of sample observations at level  $c$  of Variable B, and  $n$  is the total sample size.

- **Test statistic.** The test statistic is a chi-square random variable ( $X^2$ ) defined by the following equation.

$$X^2 = \sum [ (O_{r,c} - E_{r,c})^2 / E_{r,c} ]$$

where  $O_{r,c}$  is the observed frequency count at level  $r$  of Variable A and level  $c$  of Variable B, and  $E_{r,c}$  is the expected frequency count at level  $r$  of Variable A and level  $c$  of Variable B.

- **P-value.** The P-value is the probability of observing a sample statistic as extreme as the test statistic. Since the test statistic is a chi-square, use the [Chi-Square Distribution Calculator](#) to assess the probability associated with the test statistic. Use the degrees of freedom computed above.

## Interpret Results

If the sample findings are unlikely, given the null hypothesis, the researcher rejects the null hypothesis. Typically, this involves comparing the P-value to the [significance level](#), and rejecting the null hypothesis when the P-value is less than the significance level.

## Test Your Understanding

### Problem

A public opinion poll surveyed a simple random sample of 1000 voters. Respondents were classified by gender (male or female) and by voting preference (Republican, Democrat, or Independent). Results are shown in the contingency table below.

		Voting Preferences			Row total
		Republican	Democrat	Independent	
Male	200	150	50		400
Female	250	300	50		600
Column total	450	450	100		1000

is there a gender gap? Do the men's voting preferences differ significantly from the women's preferences? Use a 0.05 level of significance.

### Solution

The solution to this problem takes four steps:

- (1) state the hypotheses,
- (2) formulate an analysis plan,
- (3) analyze sample data, and
- (4) interpret results.

We work through those steps below:

**State the hypotheses.** The first step is to state the null hypothesis and an alternative hypothesis.

$H_0$ : Gender and voting preferences are independent.

$H_a$ : Gender and voting preferences are not independent.

**Formulate an analysis plan.** For this analysis, the significance level is 0.05. Using sample data, we will conduct a chi-square test for independence.

- **Analyze sample data.** Applying the chi-square test for independence to sample data, we compute the degrees of freedom, the expected frequency counts, and the chi-square test statistic. Based on the chi-square statistic and the degrees of freedom, we determine the P-value.

$$DF = (r - 1) * (c - 1) = (2 - 1) * (3 - 1) = 2$$

$$E_{r,c} = (n_r * n_c) / n$$

$$E_{1,1} = (400 * 450) / 1000 = 180000/1000 = 180$$

$$E_{1,2} = (400 * 450) / 1000 = 180000/1000 = 180$$

$$E_{1,3} = (400 * 100) / 1000 = 40000/1000 = 40$$

$$E_{2,1} = (600 * 450) / 1000 = 270000/1000 = 270$$

$$E_{2,2} = (600 * 450) / 1000 = 270000/1000 = 270$$

$$E_{2,3} = (600 * 100) / 1000 = 60000/1000 = 60$$



$$X^2 = \sum [ (O_{r,c} - E_{r,c})^2 / E_{r,c} ]$$

$$X^2 = (200 - 180)^2/180 + (150 - 180)^2/180 + (50 - 40)^2/40 \\ + (250 - 270)^2/270 + (300 - 270)^2/270 + (50 - 60)^2/60$$

$$X^2 = 400/180 + 900/180 + 100/40 + 400/270 + 900/270 + 100/60 \\ X^2 = 2.22 + 5.00 + 2.50 + 1.48 + 3.33 + 1.67 = 16.2$$

where DF is the degrees of freedom,  $r$  is the number of levels of gender,  $c$  is the number of levels of the voting preference,  $n_r$  is the number of observations from level  $r$  of gender,  $n_c$  is the number of observations from level  $c$  of voting preference,  $n$  is the number of observations in the sample,  $E_{r,c}$  is the expected frequency count when gender is level  $r$  and voting preference is level  $c$ , and  $O_{r,c}$  is the observed frequency count when gender is level  $r$  voting preference is level  $c$ .

The P-value is the probability that a chi-square statistic having 2 degrees of freedom is more extreme than 16.2.

We use the Chi-Square Distribution Calculator to find  $P(X^2 > 16.2) = 0.0003$ .

- **Interpret results.** Since the P-value (0.0003) is less than the significance level (0.05), we cannot accept the null hypothesis. Thus, we conclude that there is a relationship between gender and voting preference.

**Note:** If you use this approach on an exam, you may also want to mention why this approach is appropriate. Specifically, the approach is appropriate because the sampling method was simple random sampling, the variables under study were categorical, and the expected frequency count was at least 5 in each cell of the contingency table.

**Example: "Which holiday would you rather have?"**

	Beach	Cruise
Men	209	280
Women	225	248

### **Does Gender affect Preferred Holiday?**

If Gender (Man or Woman) **does** affect Preferred Holiday we say they are **dependent**.

By doing some special calculations (explained later), we come up with a "p" value:

p value is 0.13

Now,  $p < 0.05$  is the usual test for dependence. In this case **p is greater than 0.05**, so we believe the variables are **independent** (ie not linked together).

In other words Men and Women probably do **not** have a different preference for Beach Holidays or Cruises.

### **Understanding "p" Value**

"p" is the probability the variables are **independent**.

Imagine for the previous example you had tried to **fool the test** by choosing a random sample of

**Men** each time:





# M.I.E.T. ENGINEERING COLLEGE

(AUTONOMOUS)



(Approved by AICTE, New Delhi and Affiliated to Anna University, Chennai)

Accredited by NBA (CIVIL, CSE, ECE, EEE & MECH)

Accredited with 'A+' grade by NAAC  
(An ISO 9001:2015 Certified Institution)

(Recognized by UGC under section 2(f) & 12(B) of UGC Act, 1956)

TRICHY - PUDUKKOTTAI MAIN ROAD, TRICHY - 620 007



**Men:**

Beach 209, Cruise

280

**Men:**

Beach 225, Cruise

248

Is it **likely** you would get such different results surveying Men each time?

Well the "p" value of **0.132** says that it really could happen every so often.

Surveys are random after all. We expect slightly different results each time, right?

So most people want to see a **p** value less than **0.05** before they are happy to say the results show the groups have a different response.

Let's see another example:

**Example: "Which pet would you rather have?"**

	Cat	Dog
Men	207	282
Women	231	242

By doing the calculations (shown later), we come up with:

P value is 0.043

In this case **p < 0.05**, so this result is thought of as being "significant" meaning we think the variables are **not** independent.

In other words, because **0.043 < 0.05** we think that Gender is linked to Pet Preference (Men and Women have different preferences for Cats and Dogs).

*Just out of interest, notice that the numbers in our two examples are similar, but the resulting p-values are very different: **0.132** and **0.043**. This shows how sensitive the test is!*

**Why p<0.05 ?**

It is just a choice! Using **p<0.05 is common**, but we could have chosen p<0.01 to be even more sure that the groups behave differently, or any value really.

## Calculating P-Value

So how do we calculate this p-value? We use the Chi-Square Test!

### Chi-Square Test

Note: **Chi** Sounds like "Hi" but with a **K**, so say Chi-Square like "**Ki** square"

And Chi is the greek letter  $\chi$ , so we can also write it  $\chi^2$

Important points before we get started:

- This test only works for **categorical** data (data in categories), such as Gender {Men, Women} or color {Red, Yellow, Green, Blue} etc, but **not numerical** data such as height or weight.
- The numbers must be large enough. Each entry must be **5** or more. In our example we have values such as 209, 282, etc, so we are good to go.

### Our first step is to state our hypotheses:

**Hypothesis:** A statement that might be true, which can then be tested.

The two **hypotheses** are.

- Gender and preference for cats or dogs are **independent**.
- Gender and preference for cats or dogs are **not independent**.

**Subtract expected from actual, square it, then divide by expected:**

**Now add up those values:**

$$1.099 + 0.918 + 1.136 + 0.949 = 4.102$$

Chi-Square is 4.102

## From Chi-Square to p

To get from Chi-Square to p-value is a difficult calculation, so either look it up in a table, or use the [Chi-Square Calculator](#).

But first you will need a "Degree of Freedom" (DF)

### Calculate Degrees of Freedom

Multiply (rows - 1) by (columns - 1)

Example:  $DF = (2 - 1)(2 - 1) = 1 \times 1 = 1$

### Result

The result is:

$p = 0.04283$

Done!

### Chi-Square Formula

This is the formula for Chi-Square:

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

O = The observed (actual) value

E = The expected value

## Introduction to Correlation and Regression Analysis

In this section we will first discuss correlation analysis, which is used to quantify the association between two continuous variables (e.g., between an independent and a dependent variable or between two independent variables). Regression analysis is a related technique to assess the relationship between an outcome variable and one or more risk factors or confounding variables. The outcome variable is also called the **response** or **dependent variable** and the risk factors and confounders are called the **predictors**, or **explanatory** or **independent variables**. In regression analysis, the dependent variable is denoted "y" and the independent variables are denoted by "x".



[**NOTE:** The term "predictor" can be misleading if it is interpreted as the ability to predict even beyond the limits of the data. Also, the term "explanatory variable" might give an impression of a causal effect in a situation in which inferences should be limited to identifying associations. The terms "independent" and "dependent" variable are less subject to these interpretations as they do not strongly imply cause and effect.

### Correlation Analysis

In correlation analysis, we estimate a sample correlation coefficient, more specifically the Pearson Product Moment correlation coefficient. The sample correlation coefficient, denoted  $r$ , ranges between -1 and +1 and quantifies the direction and strength of the linear association between the two variables. The correlation between two variables can be positive (i.e., higher levels of one variable are associated with higher levels of the other) or negative (i.e., higher levels of one variable are associated with lower levels of the other).

The sign of the correlation coefficient indicates the direction of the association. The magnitude of the correlation coefficient indicates the strength of the association.

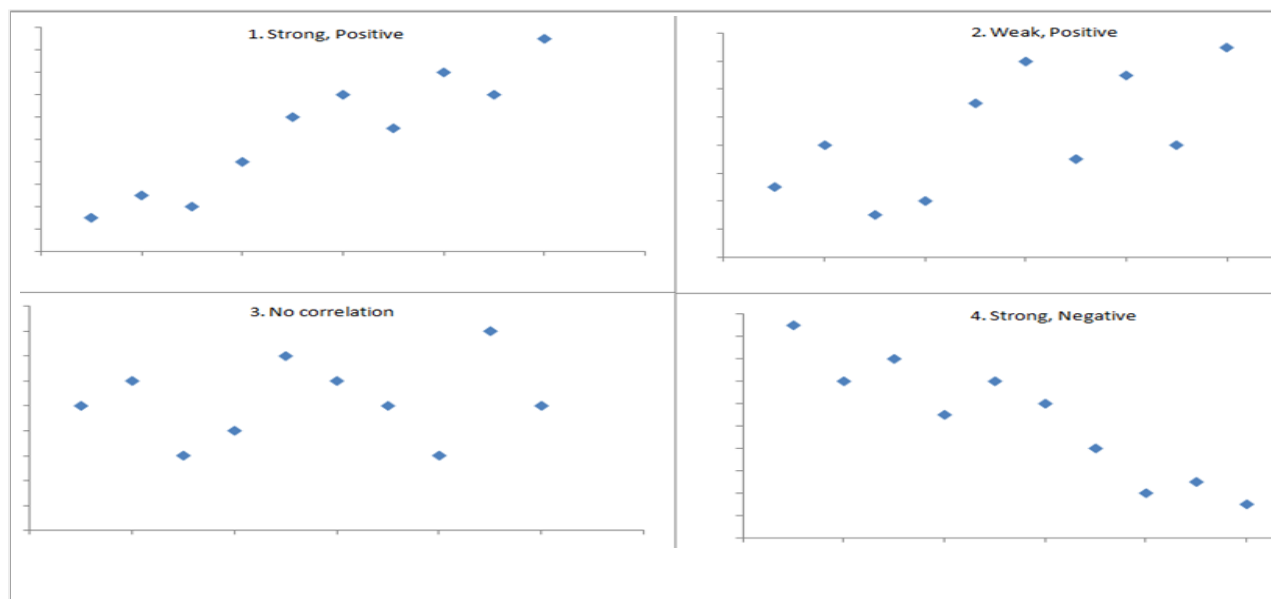
For example, a correlation of  $r = 0.9$  suggests a strong, positive association between two variables, whereas a correlation of  $r = -0.2$  suggest a weak, negative association. A correlation close to zero suggests no linear association between two continuous variables.

LISA: [I find this description confusing. You say that the correlation coefficient is a measure of the "strength of association", but if you think about it, isn't the slope a better measure of association? We use risk ratios and odds ratios to quantify the strength of association, i.e., when an exposure is present it has how many times more likely the outcome is. The analogous quantity in correlation is the slope, i.e., for a given increment in the independent variable, how many times is the dependent variable going to increase?

And " $r$ " (or perhaps better R-squared) is a measure of how much of the variability in the dependent variable can be accounted for by differences in the independent variable. The analogous measure for a dichotomous variable and a dichotomous outcome would be the attributable proportion, i.e., the proportion of  $Y$  that can be attributed to the presence of the exposure.]

It is important to note that there may be a non-linear association between two continuous variables, but computation of a correlation coefficient does not detect this. Therefore, it is always important to evaluate the data carefully before computing a correlation coefficient. Graphical displays are particularly useful to explore associations between variables.

The figure below shows four hypothetical scenarios in which one continuous variable is plotted along the X-axis and the other along the Y-axis.



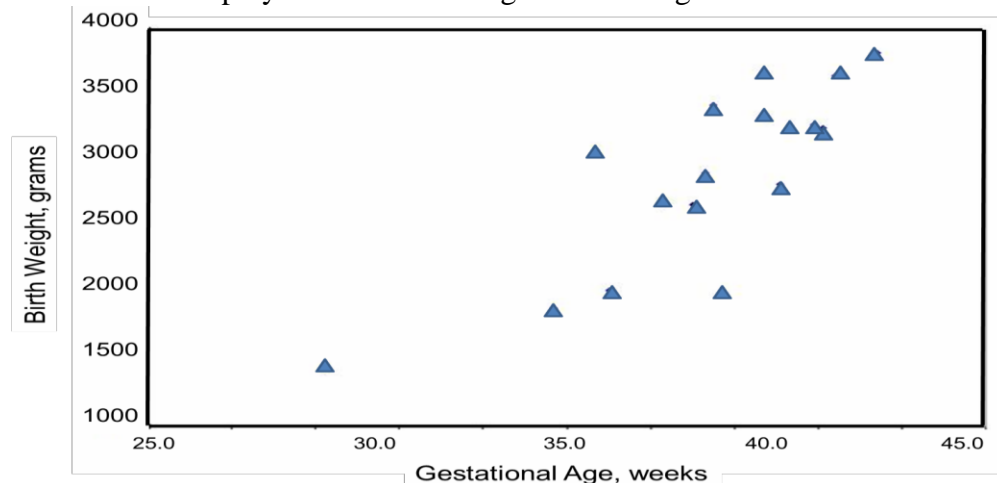
- Scenario 1 depicts a strong positive association ( $r=0.9$ ), similar to what we might see for the correlation between infant birth weight and birth length.
- Scenario 2 depicts a weaker association ( $r=0.2$ ) that we might expect to see between age and body mass index (which tends to increase with age).
- Scenario 3 might depict the lack of association ( $r$  approximately 0) between the extent of media exposure in adolescence and age at which adolescents initiate sexual activity.
- Scenario 4 might depict the strong negative association ( $r= -0.9$ ) generally observed between the number of hours of aerobic exercise per week and percent body fat.

### Example - Correlation of Gestational Age and Birth Weight

A small study is conducted involving 17 infants to investigate the association between gestational age at birth, measured in weeks, and birth weight, measured in grams.

Infant ID #	Gestational Age (wks)	Birth Weight (gm)
1	34.7	1895
2	36.0	2030
3	29.3	1440
4	40.1	2835
5	35.7	3090
6	42.4	3827
7	40.3	3260
8	37.3	2690
9	40.9	3285
10	38.3	2920
11	38.5	3430
12	41.4	3657
13	39.7	3685
14	39.7	3345
15	41.1	3260
16	38.0	2680
17	38.7	2005

We wish to estimate the association between gestational age and infant birth weight. In this example, birth weight is the dependent variable and gestational age is the independent variable. Thus  $y$ =birth weight and  $x$ =gestational age. The data are displayed in a scatter diagram in the figure below.



Each point represents an (x,y) pair (in this case the gestational age, measured in weeks, and the birth weight, measured in grams). Note that the independent variable is on the horizontal axis (or X-axis), and the dependent variable is on the vertical axis (or Y-axis). The scatter plot shows a positive or direct association between gestational age and birth weight. Infants with shorter gestational ages are more likely to be born with lower weights and infants with longer gestational ages are more likely to be born with higher weights.

The formula for the sample correlation coefficient is

$$r = \frac{\text{Cov}(x,y)}{\sqrt{s_x^2 * s_y^2}}$$

where Cov(x,y) is the covariance of x and y defined as

$$\text{Cov}(x,y) = \frac{\sum(X - \bar{X})(Y - \bar{Y})}{n - 1}$$

$s_x^2$  and  $s_y^2$  are the sample variances of x and y, defined as

$$s_x^2 = \frac{\sum(X - \bar{X})^2}{n - 1} \quad \text{and} \quad s_y^2 = \frac{\sum(Y - \bar{Y})^2}{n - 1}$$

The variances of x and y measure the variability of the x scores and y scores around their respective sample means (

$\bar{X}$  and  $\bar{Y}$ , considered separately). The covariance measures the variability of the (x,y) pairs around the mean of x and mean of y, considered simultaneously.

To compute the sample correlation coefficient, we need to compute the variance of gestational age, the variance of birth weight and also the covariance of gestational age and birth weight.

We first summarize the gestational age data. The mean gestational age is:

$$\bar{X} = \frac{\sum X}{n} = \frac{652.1}{17} = 38.4.$$

Infant ID #	Gestational Age	$(X - \bar{X})$	$(X - \bar{X})^2$
1	34.7	-3.7	13.69
2	36.0	-2.4	5.76
3	29.3	-9.1	82.81
4	40.1	1.7	2.89
5	35.7	-2.7	7.29
6	42.4	4.0	16.00
7	40.3	1.9	3.61
8	37.3	-1.1	1.21
9	40.9	2.5	6.25
10	38.3	-0.1	0.01
11	38.5	0.1	0.01
12	41.4	3.0	9.00
13	39.7	1.3	1.69
14	39.7	1.3	1.69
15	41.1	2.7	7.29
16	38.0	-0.4	0.16
17	38.7	0.3	0.09
	$\sum X = 652.1$	$\sum (X - \bar{X}) = 0$	$\sum (X - \bar{X})^2 = 159.45$

To compute the variance of gestational age, we need to sum the squared deviations (or differences) between each observed gestational age and the mean gestational age. The computations are summarized below.

The variance of gestational age is:

$$s_x^2 = \frac{\sum (X - \bar{X})^2}{n - 1} = \frac{159.45}{16} = 10.0.$$

Next, we summarize the birth weight data. The mean birth weight is:

The variance of birth weight is computed just as we did for gestational age as shown in the table below.

$$\bar{Y} = \frac{\sum Y}{n} = \frac{49,334}{17} = 2902.$$



Infant ID #	Birth Weight	$(Y - \bar{Y})$	$(Y - \bar{Y})^2$
1	1895	-1007	1,014,049
2	2030	-872	760,384
3	1440	-1462	2,137,444
4	2835	-67	4,489
5	3090	188	35,344
6	3827	925	855,625
7	3260	358	128,164
8	2690	-212	44,944
9	3285	383	146,689
10	2920	18	324
11	3430	528	278,784
12	3657	755	570,025
13	3685	783	613,089
14	3345	443	196,249
15	3260	358	128,164
16	2680	-222	49,284
17	2005	-897	804,609
	$\Sigma Y = 49,334$	$\Sigma (Y - \bar{Y}) = 0$	$\Sigma (Y - \bar{Y})^2 = 7,767,660$

The variance of birth weight is:

Next we compute the covariance,

$$\text{Cov}(x, y) = \frac{\Sigma(X - \bar{X})(Y - \bar{Y})}{n - 1}$$

To compute the covariance of gestational age and birth weight, we need to multiply the deviation from the mean gestational age by the deviation from the mean birth weight for each participant (i.e.,

$$(X - \bar{X})(Y - \bar{Y}))$$

The computations are summarized below. Notice that we simply copy the deviations from the mean gestational age and birth weight from the two tables above into the table below and multiply.

$$s_y^2 = \frac{\Sigma(Y - \bar{Y})^2}{n - 1} = \frac{7,767,660}{16} = 485,478.8.$$

Infant Identification Number	$(X - \bar{X})$	$(Y - \bar{Y})$	$(X - \bar{X})(Y - \bar{Y})$
1	-3.7	-1007	3725.9
2	-2.4	-872	2092.8
3	-9.1	-1462	13,304.2
4	1.7	-67	-113.9
5	-2.7	188	-507.6
6	4.0	925	3700.0
7	1.9	358	680.2
8	-1.1	-212	233.2
9	2.5	383	957.5
10	-0.1	18	-1.8
11	0.1	528	52.8
12	3.0	755	2265.0
13	1.3	783	1017.9
14	1.3	443	575.9
15	2.7	358	966.6
16	-0.4	-222	88.8
17	0.3	-897	-269.1
			$\Sigma (X - \bar{X})(Y - \bar{Y}) = 28,768.4$

The covariance of gestational age and birth weight is:

$$s_y^2 = \frac{\Sigma(Y - \bar{Y})^2}{n - 1} = \frac{7,767,660}{16} = 485,578.8.$$

We now compute the sample correlation coefficient:

$$r = \frac{\text{Cov}(x,y)}{\sqrt{s_x^2 * s_y^2}} = \frac{1798.0}{\sqrt{10.0 * 485,578.8}} = \frac{1798.0}{2199.4} = 0.82$$

Not surprisingly, the sample correlation coefficient indicates a strong positive correlation.

As we noted, sample correlation coefficients range from -1 to +1. In practice, meaningful correlations (i.e., correlations that are clinically or practically important) can be as small as 0.4 (or -0.4) for positive (or negative) associations. There are also statistical tests to determine whether an observed correlation is statistically significant or not (i.e., statistically significantly different from zero). Procedures to test whether an observed sample correlation is suggestive of a statistically significant correlation are described in detail in Kleinbaum, Kupper and Muller.<sup>1</sup>

### Supearman Rank Correlation :

The Pearson correlation coefficient between the ranked variables has been termed as the Spearman correlation coefficient. It is also referred as 'grade correlation'.

#### Formula :

$$R = 1 - ((6 \sum d^2) / (n^3 - n))$$

**Partial correlation analysis involves studying the linear relationship between two variables after excluding the effect of one or more independent factors.**

Simple [correlation](#) does not prove to be an all-encompassing technique especially under the above circumstances. In order to get a correct picture of the relationship between two variables, we should first eliminate the influence of other variables.

For example, study of partial correlation between price and demand would involve studying the relationship between price and demand excluding the effect of money supply, exports, etc.

**Partial correlation analysis involves studying the linear relationship between two variables after excluding the effect of one or more independent factors.**

Simple [correlation](#) does not prove to be an all-encompassing technique especially under the above circumstances. In order to get a correct picture of the relationship between two variables, we should first eliminate the influence of other variables.

For example, study of partial correlation between price and demand would involve studying the

relationship between price and demand excluding the effect of money supply, exports, etc.

**Partial correlation analysis involves studying the linear relationship between two variables after excluding the effect of one or more independent factors.**

Simple [correlation](#) does not prove to be an all-encompassing technique especially under the above circumstances. In order to get a correct picture of the relationship between two variables, we should first eliminate the influence of other variables.

For example, study of partial correlation between price and demand would involve studying the relationship between price and demand excluding the effect of money supply, exports, etc.

### **Multiple Correlation**

Another technique used to overcome the drawbacks of simple correlation is [multiple regression analysis](#).

Here, we study the effects of all the independent variables simultaneously on a dependent variable. For example, the correlation co-efficient between the yield of paddy ( $X_1$ ) and the other variables, viz. type of seedlings ( $X_2$ ), manure ( $X_3$ ), rainfall ( $X_4$ ), humidity ( $X_5$ ) is the multiple correlation co-efficient  $R_{1.2345}$ . This co-efficient takes value between 0 and +1.

The limitations of multiple correlation are similar to those of partial correlation. If multiple and partial correlation are studied together, a very useful analysis of the relationship between the different variables is possible.



# **M.I.E.T. ENGINEERING COLLEGE**

**(AUTONOMOUS)**

(Approved by AICTE, New Delhi and Affiliated to Anna University, Chennai)

**Accredited by NBA (CIVIL, CSE, ECE, EEE & MECH)**

**Accredited with 'A+' grade by NAAC**

**(An ISO 9001:2015 Certified Institution)**

**(Recognized by UGC under section 2(f) & 12(B) of UGC Act, 1956)**

**TRICHY - PUDUKKOTTAI MAIN ROAD, TRICHY - 620 007**



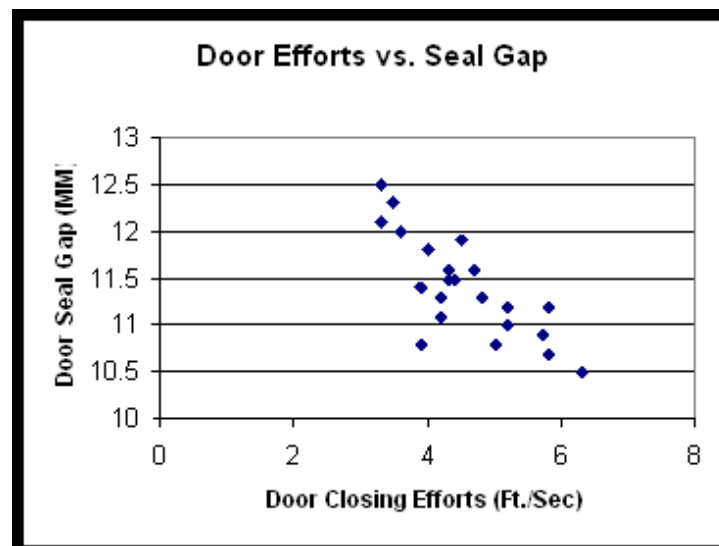
## **UNIT-V REGRESSION ANALYSIS**

### **Regression Analysis**

#### **Introduction**

As you develop Cause & Effect diagrams based on data, you may wish to examine the degree of correlation between variables. A statistical measurement of correlation can be calculated using the least squares method to quantify the strength of the relationship between two variables. The output of that calculation is the **Correlation Coefficient, or (r)**, which ranges between -1 and 1. A value of 1 indicates perfect positive correlation - as one variable increases, the second increases in a linear fashion. Likewise, a value of -1 indicates perfect negative correlation - as one variable increases, the second decreases. A value of zero indicates zero correlation.

Before calculating the Correlation Coefficient, the first step is to construct a scatter diagram. Most spreadsheets, including Excel, can handle this task. Looking at the scatter diagram will give you a broad understanding of the correlation. Following is a scatter plot chart example based on an automobile manufacturer.



In this case, the process improvement team is analyzing door closing efforts to understand what the causes could be. The Y-axis represents the width of the gap between the sealing flange of a car door and the sealing flange on the body - a measure of how tight the door is set to the body. The fishbone diagram indicated that variability in the seal gap could be a cause of variability in door closing efforts.

In this case, you can see a pattern in the data indicating a negative correlation (negative slope) between the two variables. In fact, the Correlation Coefficient is -0.78, indicating a strong inverse or negative relationship.



**More Steam Note:** *It is important to note that Correlation is not Causation* - two variables can be very strongly correlated, but both can be caused by a third variable. For example, consider two variables: A) how much my grass grows per week, and B) the average depth of the local reservoir. Both variables could be highly correlated because both are dependent upon a third variable - how much it rains.

In our car door example, it makes sense that the tighter the gap between the sheet metal sealing surfaces (before adding weatherstrips and trim), the harder it is to close the door. So a rudimentary understanding of mechanics would support the hypothesis that there is a causal relationship. Other industrial processes are not always as obvious as these simple examples, and determination of causal relationships may require more extensive experimentation (Design of Experiments).

### Simple Regression Analysis

While Correlation Analysis assumes no causal relationship between variables, Regression Analysis assumes that one variable is dependent upon: A) another single independent variable (Simple Regression) , or B) multiple independent variables (Multiple Regression).

Regression plots a line of best fit to the data using the least-squares method. You can see an example below of linear regression using the same car door scatter plot:

You can see that the data is clustered closely around the line, and that the line has a downward slope. There is strong negative correlation expressed by two related statistics: the  $r$  value, as stated before is,  $-0.78$  the  $r^2$  value is therefore  $0.61$ .  $R^2$ , called the **Coefficient of Determination**, expresses how much of the variability in the dependent variable is explained by variability in the independent variable. You may find that a non-linear equation such as an exponential or power function may provide a better fit and yield a higher  $r^2$  than a linear equation.

These statistical calculations can be made using Excel, or by using any of several statistical analysis software packages. More Steam provides links to statistical software downloads, including free software.

### Multiple Regression Analysis

Multiple Regression Analysis uses a similar methodology as Simple Regression, but includes more than one independent variable. Econometric models are a good example, where the dependent variable of GNP may be analyzed in terms of multiple independent variables, such as interest rates, productivity growth, government spending, savings rates, consumer confidence, etc.

Many times historical data is used in multiple regression in an attempt to identify the most significant inputs to a process. The benefit of this type of analysis is that it can be done very quickly and relatively simply. However, there are **several potential pitfalls**:

- The **data may be inconsistent** due to different measurement systems, calibration drift, different operators, or recording errors.
- The **range of the variables may be very limited**, and can give a false indication of low correlation. For example, a process may have temperature controls because temperature has been found in the past to have an impact on the output. Using historical temperature data may therefore indicate low significance because the range of temperature is already controlled in tight tolerance.
- There may be a **time lag that influences the relationship** - for example, temperature may be much more critical at an early point in the process than at a later point, or vice-versa. There also may be inventory effects that must be taken into account to make sure that all measurements are taken at a consistent point in the process.

Once again, it is critical to remember that correlation is not causality. As stated by Box, Hunter and Hunter: "Broadly speaking, **to find out what happens when you change something, it is necessary to change it**. To safely infer causality the experimenter cannot rely on natural happenings to choose the design for him; he must choose the design for himself and, in particular, must introduce randomization to break the links with possible lurking variables".<sup>1</sup>

Returning to our example of door closing efforts, you will recall that the door seal gap had an  $r^2$  of 0.61. Using multiple regression, and adding the additional variable "door weatherstrip durometer" (softness), the  $r^2$  rises to 0.66. So the durometer of the door weatherstrip added some

explaining power, but minimal. Analyzed individually, durometer had much lower correlation with door closing efforts - only 0.41.

This analysis was based on historical data, so as previously noted, the regression analysis only tells us what did have an impact on door efforts, not what could have an impact. If the range of durometer measurements was greater, we might have seen a stronger relationship with door closing efforts, and more variability in the output.

### Trend Analysis

There are no proven "automatic" techniques to identify trend components in the time series data; however, as long as the trend is monotonous (consistently increasing or decreasing) that part of data analysis is typically not very difficult. If the time series data contain considerable error, then the first step in the process of trend identification is smoothing.

**Smoothing.** Smoothing always involves some form of local averaging of data such that the nonsystematic components of individual observations cancel each other out. The most common technique is *moving average* smoothing which replaces each element of the series by either the simple or weighted average of  $n$  surrounding elements, where  $n$  is the width of the smoothing "window" (see Box & Jenkins, 1976; Velleman & Hoaglin, 1981). Medians can be used instead of means. The main advantage of median as compared to moving average smoothing is that its results are less biased by outliers (within the smoothing window). Thus, if there are outliers in

the data (e.g., due to measurement errors), median smoothing typically produces smoother or at least more "reliable" curves than moving average based on the same window width. The main disadvantage of median smoothing is that in the absence of clear outliers it may produce more "jagged" curves than moving average and it does not allow for weighting.

In the relatively less common cases (in time series data), when the measurement error is very large, the *distance weighted least squares smoothing* or *negative exponentially weighted smoothing* techniques can be used. All those methods will filter out the noise and convert the data into a smooth curve that is relatively unbiased by outliers (see the respective sections on each of those methods for more details). Series with relatively few and systematically distributed points can be smoothed with *bicubic splines*.

**Fitting a function.** Many monotonous time series data can be adequately approximated by a linear function; if there is a clear monotonous nonlinear component, the data first need to be transformed to remove the nonlinearity. Usually a logarithmic, exponential, or (less often) polynomial function can be used.

### Additive models

The models that we have considered in earlier sections have been **additive models**, and there has been an implicit assumption that the different components affected the time series additively.

$$\text{Data} = \text{Seasonal effect} + \text{Trend} + \text{Cyclical} + \text{Residual}$$

For monthly data, an additive model assumes that the difference between the January and July values is approximately the same each year. In other words, the **amplitude** of the seasonal effect is the same each year.

The model similarly assumes that the residuals are roughly the same size throughout the series -- they are a random component that adds on to the other components in the same way at all parts of the series.

### Multiplicative models

In many time series involving **quantities** (e.g. money, wheat production, ...), the absolute differences in the values are of less interest and importance than the percentage changes.

For example, in seasonal data, it might be more useful to model that the July value is the same **proportion** higher than the January value in each year, rather than assuming that their difference is constant. Assuming that the seasonal and other effects act proportionally on the series is equivalent to a **multiplicative model**,

$$\text{Data} = (\text{Seasonal effect}) \times \text{Trend} \times \text{Cyclical} \times \text{Residual}$$

Fortunately, multiplicative models are equally easy to fit to data as additive models! The trick to fitting a multiplicative model is to take logarithms of both sides of the model,

$$\begin{aligned}\log(\text{Data}) &= \log(\text{Seasonal effect} \times \text{Trend} \times \text{Cyclical} \times \text{Residual}) \\ &= \log(\text{Seasonal effect}) + \log(\text{Trend}) \\ &\quad + \log(\text{Cyclical}) + \log(\text{Residual})\end{aligned}$$

After taking logarithms (either natural logarithms or to base 10), the four components of the time series again act additively.

### **What is an additive model?**

A data model in which the effects of individual factors are differentiated and added together to model the data. They occur in several Minitab commands:

- An additive model is optional for Decomposition procedures and for Winters' method.
- An additive model is optional for two-way ANOVA procedures. Choose this option to omit the interaction term from the model.

### **What is a multiplicative model?**

This model assumes that as the data increase, so does the seasonal pattern. Most time series plots exhibit such a pattern. In this model, the trend and seasonal components are multiplied and then added to the error component.

### **Should I use an additive model or a multiplicative model?**

Choose the multiplicative model when the magnitude of the seasonal pattern in the data depends on the magnitude of the data. In other words, the magnitude of the seasonal pattern increases as the data values increase, and decreases as the data values decrease.

Choose the additive model when the magnitude of the seasonal pattern in the data does not depend on the magnitude of the data. In other words, the magnitude of the seasonal pattern does not change as the series goes up or down.

If the pattern in the data is not very obvious, and you have trouble choosing between the additive and multiplicative procedures, you can try both and choose the one with smaller accuracy measures.



## INDEX NUMBERS

### Introduction:

Index numbers are meant to study the change in the effects of such factors which cannot be measured directly. According to Bowley, -Index numbers are used to measure the changes in some quantity which we cannot observe directly. For example, changes in business activity in a country are not capable of direct measurement but it is possible to study relative changes in business activity by studying the variations in the values of some such factors which affect business activity, and which are capable of direct measurement.

Index numbers are commonly used statistical device for measuring the combined fluctuations in a group related variables. If we wish to compare the price level of consumer items today with that prevalent ten years ago, we are not interested in comparing the prices of only one item, but in comparing some sort of average price levels. We may wish to compare the present agricultural production or industrial production with that at the time of independence. Here again, we have to consider all items of production and each item may have undergone a different fractional increase (or even a decrease). How do we obtain a composite measure? This composite measure is provided by index numbers which may be defined as a device for combining the variations that have come in group of related variables over a period of time, with a view to obtain a figure that represents the net result of the change in the constitute variables.

Index numbers may be classified in terms of the variables that they are intended to measure. In business, different groups of variables in the measurement of which index number techniques are commonly used are (i) price, (ii) quantity, (iii) value and (iv) business activity. Thus, we have index of wholesale prices, index of consumer prices, index of industrial output, index of value of exports and index of business activity, etc. Here we shall be mainly interested in index numbers of prices showing changes with respect to time, although methods described can be applied to other cases. In general, the present level of prices is compared with the level of prices in the past. The present period is called the current period and some period in the past is called the base period.

### Index Numbers:



Index numbers are statistical measures designed to show changes in a variable or group of related variables with respect to time, geographic location or other characteristics such as income, profession, etc. A collection of index numbers for different years, locations, etc., is sometimes called an index series.

### **Simple Index Number:**

A simple index number is a number that measures a relative change in a single variable with respect to a base.

### **Composite Index Number:**

A composite index number is a number that measures an average relative changes in a group of relative variables with respect to a base.

### **Types of Index Numbers:**

Following types of index numbers are usually used:

#### **Price index Numbers**

Price index numbers measure the relative changes in prices of a commodities between two periods. Prices can be either retail or wholesale.

#### **Quantity Index Numbers:**

These index numbers are considered to measure changes in the physical quantity of goods produced, consumed or sold of an item or a group of items.

#### **Uses**

This index number is a useful number that helps us quantify changes in our field. It is easier to see one value than a thousand different values for each item in our field.

Take the stock market, for example. It is comprised of thousands of different public companies. We could, of course, look at the stock value of each of these companies to see how the companies are doing as a whole, or we can look at just one number, the stock index, to get a

general feel for how the companies are doing.

The same goes for the cost of goods. We could look at the cost of each item and compare it to its cost from last year. But that would mean looking at the cost of millions of items. Or we could look at the cost of goods index, just one number, to see whether prices have increased or decreased over the past year.

We can say that the index number is one simple number that we can look at to give us a general overview of what is happening in our field. Let's take a look at two real world index numbers.

### Line of Best Fit (Least Square Method)

A **line of best fit** is a straight line that is the best approximation of the given set of data.

It is used to study the nature of the relation between two variables. (We're only considering the two-dimensional case, here.)

A line of best fit can be roughly determined using an eyeball method by drawing a straight line on a scatter plot so that the number of points above the line and below the line is about equal (and the line passes through as many points as possible).

A more accurate way of finding the line of best fit is the **least square method**.

Use the following steps to find the equation of line of best fit for a set of ordered pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ .

Step 1: Calculate the mean of the  $xx$ -values and the mean of the  $yy$ -values.

$$\bar{X} = \frac{\sum_{i=1}^n x_i}{n}, \bar{Y} = \frac{\sum_{i=1}^n y_i}{n}$$

Step 2: The following formula gives the slope of the line of best fit:

$$m = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2}$$

Step 3: Compute the yy-intercept of the line by using the formula:

$$b = \bar{Y} - m\bar{X}$$

Step 4: Use the slope  $m$  and the  $yy$ -intercept  $b$  to form the equation of the line.

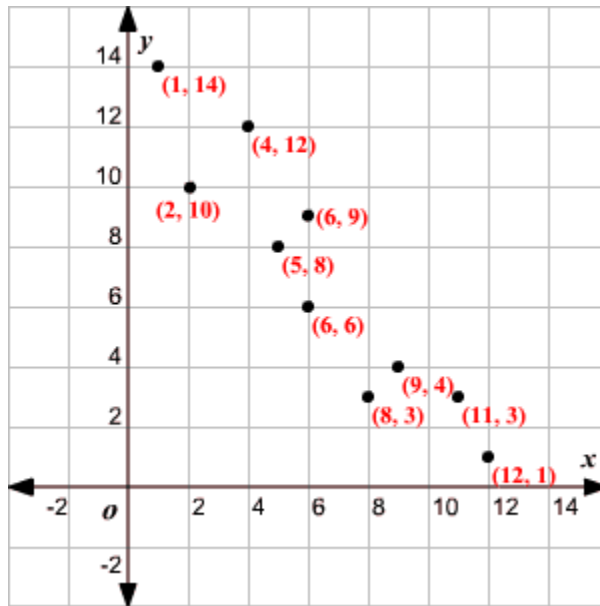
**Example:**

Use the least square method to determine the equation of line of best fit for the data. Then plot the line.

$x$	8	22	111	6	5	44	121	9	6	11
$x$	8		1	6	5		2	9	6	
$y$	3	101	33	6	8	121	11	4	9	141
$y$	3	0		6	8	2		4	9	4

**Solution:**

Plot the points on a [coordinate plane](#).



Calculate the means of the  $x$ -values and the  $y$ -values.

$$\bar{x} = \frac{1+2+4+5+6+6+8+9+11+12+12+12}{12} = 6.4$$

$$\bar{y} = \frac{14+10+12+8+9+6+3+4+3+1+1+1}{12} = 7$$

Now calculate  $x_i - \bar{x}$ ,  $y_i - \bar{y}$ ,  $(x_i - \bar{x})(y_i - \bar{y})$ , and  $(x_i - \bar{x})^2$  for each  $i$ .

$i$	$x_i$	$y_i$	$x_i - \bar{X}$	$y_i - \bar{Y}$	$(x_i - \bar{X})(y_i - \bar{Y})$	$(x_i - \bar{X})^2$
11	88	33	1.6	-4	-6.4	2.56
22	22	10	-4.4	3	-13.2	19.36
33	11	33	-4.6	-4	18.4	21.16
44	66	66	-0.4	-1	0.4	0.16
55	55	88	-1.4	1	-1.4	1.96
66	44	12	-2.4	5	-12	5.76
77	12	11	-5.6	-6	33.6	31.36
88	99	44	-2.6	-3	7.8	6.76
99	66	99	-0.4	2	-0.8	0.16
10	11	14	-5.4	7	-37.8	29.16

Calculate the slope.

$$m = \frac{\sum_{i=1}^n (x_i - \bar{X})(y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{X})^2} = \frac{-131.1}{118.4} \approx -1.1$$

Calculate the yy-intercept.

Use the formula to compute the yy-intercept.

$$b = \bar{Y} - m\bar{X} = 7 - (-1.1 \times 6.4) = 7 + 7.04 \approx 14.0$$

Use the slope and yy-intercept to form the equation of the line of best fit.

The slope of the line is -1.1 and the yy-intercept is 14.0.

Therefore, the equation is  $y = -1.1x + 14.0$ .

Draw the line on the scatter plot.