

**DIRECTORATE OF DISTANCE EDUCATION  
UNIVERSITY OF NORTH BENGAL**

**MASTER OF ARTS-POLITICAL SCIENCES  
SEMESTER -III**

**RESEARCH METHODOLOGY  
SOFT CORE 303  
BLOCK-2**

---

## UNIVERSITY OF NORTH BENGAL

Postal Address:

The Registrar,  
University of North Bengal,  
Raja Rammohunpur,  
P.O.-N.B.U., Dist-Darjeeling,  
West Bengal, Pin-734013,  
India.

Phone: (O) +91 0353-2776331/2699008

Fax: (0353) 2776313, 2699001

Email: [regnbu@sancharnet.in](mailto:regnbu@sancharnet.in) ; [regnbu@nbu.ac.in](mailto:regnbu@nbu.ac.in)

Website: [www.nbu.ac.in](http://www.nbu.ac.in)

First Published in 2019



All rights reserved. No Part of this book may be reproduced or transmitted, in any form or by any means, without permission in writing from University of North Bengal. Any person who does any unauthorised act in relation to this book may be liable to criminal prosecution and civil claims for damages.

This book is meant for educational and learning purpose. The authors of the book has/have taken all reasonable care to ensure that the contents of the book do not violate any existing copyright or other intellectual property rights of any person in any manner whatsoever. In the even the Authors has/ have been unable to track any source and if any copyright has been inadvertently infringed, please notify the publisher in writing for corrective action

## **FOREWORD**

The Self Learning Material (SLM) is written with the aim of providing simple and organized study content to all the learners. The SLMs are prepared on the framework of being mutually cohesive, internally consistent and structured as per the university's syllabi. It is a humble attempt to give glimpses of the various approaches and dimensions to the topic of study and to kindle the learner's interest to the subject

We have tried to put together information from various sources into this book that has been written in an engaging style with interesting and relevant examples. It introduces you to the insights of subject concepts and theories and presents them in a way that is easy to understand and comprehend.

We always believe in continuous improvement and would periodically update the content in the very interest of the learners. It may be added that despite enormous efforts and coordination, there is every possibility for some omission or inadequacy in few areas or topics, which would definitely be rectified in future.

We hope you enjoy learning from this book and the experience truly enrich your learning and help you to advance in your career and future endeavours.

---

---

# RESEARCH METHODOLOGY

---

## **BLOCK-1**

Unit 1: Contending perspective of social science research (Brief Outline)

Unit 2: Positivism

Unit 3: Positivism and its Critique

Unit 4: Marxism as a method

Unit 5: Post-structuralism

Unit 6: Two strands of research: Quantitative and Qualitative debates

Unit 7: Quantitative Research methods: Sampling

## **BLOCK -2**

**Unit 8: descriptive and inferential statistics (uni-variate and bivariate analysis) .....6**

**Unit 9: correlation and regression .....33**

**Unit 10: hypothesis testing, t-test, z-test, chisquare) .....62**

**Unit 11: QUALITATIVE RESEARCH: theoretical sampling, case studies .....100**

**Unit 12: Research process: review of literature, identifying research problems .....129**

**Unit 13: Hypothesis and variables .....152**

**Unit 14: research method, primary and secondary data, style and REFERENCE, RESEARCH report .....185**

---

# **BLOCK 2 : RESEARCH METHODOLOGY**

---

## **Introduction to the Block**

Unit 8 deals with the Descriptive and inferential statistics (uni-variate and bivariate analysis)

Unit 9 deals with the Correlation and Regression.

Unit 10 deals with the Hypothesis testing, t-test, z-test, chisquare

Unit 11 deal with Qualitative: Theoretical sampling and Case studies

Unit 12 deal with Research process: review of literature, identifying research problems

Unit 13 deals with Hypothesis and variables,

Unit 14 deal with Research method and research report.

---

# **UNIT 8: DESCRIPTIVE AND INFERENCEAL STATISTICS (UNI- VARIATE AND BIVARIATE ANALYSIS**

---

## **STRUCTURE**

- 8.0 Objectives
- 8.1 Introduction
- 8.2 Meaning of Descriptive Statistics
- 8.3 Organization of Data
  - 8.3.1 Classification
    - 8.3.1.1 Frequency Distribution can be with Ungrouped Data and Grouped Data
    - 8.3.1.2 Types of Frequency Distribution
  - 8.3.2 Tabulation
  - 8.3.3 Graphical Presentation of Data
    - 8.3.3.1 Cumulative Frequency Curve or Ogive
  - 8.3.4 Diagrammatic Presentation of Data
- 8.4 Summarization of Data
  - 8.4.1 Measures of Central Tendency
  - 8.4.2 Measures of Dispersion
  - 8.4.3 Skewness and Kurtosis
  - 8.4.4 Advantages and Disadvantages of Descriptive Statistics
- 8.5 Meaning of Inferential Statistics
  - 8.5.1 Estimation
  - 8.5.2 Point Estimation
  - 8.5.3 Interval Estimation
- 8.6 Hypothesis Testing
  - 8.6.1 Statement of Hypothesis
  - 8.6.2 Level of Significance
  - 8.6.3 One Tail and Two Tail Test
- 8.7 Errors in Hypothesis Testing
  - 8.7.1 Type I Error
  - 8.7.2 Type II Error
  - 8.7.3 Power of a Test
- 8.8 General Procedure for Testing A Hypothesis

- 8.9 Let us sum up
- 8.10 Key Words
- 8.11 Questions for Review
- 8.12 Suggested readings and references
- 8.13 Answers to Check Your Progress

---

## 8.0 OBJECTIVES

---

After going through this unit, you will be able to:

- To define the nature and meaning of descriptive statistics;
- To describe the methods of organising and condensing raw data;
- To explain concept and meaning of different measures of central tendency;
- To analyse the meaning of different measures of dispersion;
- To define inferential statistics;
- To explain the concept of estimation;
- To distinguish between point estimation and interval estimation;
- and
- To explain the different concepts involved in hypothesis testing.

---

## 8.1 INTRODUCTION

---

In this unit we will be dealing with descriptive and inferential statistics. First we start with defining descriptive statistics and indicate how to organise the data, classify, tabulate etc. This unit also presents as to how the data should be presented graphically. Once the data is collected the same has to be made meaningful which can be done through averaging the data or working out the variances in the data etc. Then we deal with the advantages and disadvantages of descriptive statistics. This is followed by defining what is inferential statistics and delineating its meaning. In this unit the student will also gain knowledge regarding point and interval estimation so as to validate the results. We also learn in this unit about hypothesis testing, how it is done and the methods

thereof. We also deal with different types of errors in hypothesis testing including sampling error etc.

---

## 8.2 MEANING OF DESCRIPTIVE STATISTICS

---

The word statistics has different meaning to different persons. For some, it is a onenumber description of a set of data. Some consider statistics in terms of numbers used as measurements or counts. Mathematicians use statistics to describe data in one word. It is a summary of an event for them. Number , n, is the statistic describing how big the set of numbers is, how many pieces of data are in the set. Also, knowledge of statistics is applicable in day to day life in different ways. Statistics is used by people to take decision about the problems on the basis of different types of information available to them. However, in behavioural sciences the word ‘statistics’ means something different, that is its prime function is to draw statistical inference about population on the basis of available quantitative and qualitative information. The word statistics can be defined in two different ways. In singular sense ‘Statistics’ refers to what is called statistical methods. When ‘Statistics’ is used in plural sense it refers to ‘data’. In this unit we will use the term ‘statistics’ in singular sense. In this context, it is described as a branch of science which deals with the collection of data, their classification, analysis and interpretations of statistical data. The science of statistics may be broadly studied under two headings:

- i) Descriptive Statistics, and
- ii) Inferential Statistics

i) Descriptive Statistics: Most of the observations in this universe are subject to variability, especially observations related to human behaviour. It is a well known fact that attitude, intelligence and personality differ from individual to individual. In order to make a sensible definition of the group or to identify the group with reference to their observations/ scores, it is necessary to express them in a precise manner. For this purpose observations need to be expressed as a single



estimate which summarises the observations. Descriptive statistics is a branch of statistics, which deals with descriptions of obtained data. On the basis of these descriptions a particular group of population is defined for corresponding characteristics. The descriptive statistics include classification, tabulation, diagrammatic and graphical presentation of data, measures of central tendency and variability. These measures enable the researchers to know about the tendency of data or the scores, which further enhance the ease in description of the phenomena. Such single estimate of the series of data which summarises the distribution are known as parameters of the distribution. These parameters define the distribution completely. Basically descriptive statistics involves two operations:

- (i) Organisation of data, and
- (ii) Summarination of data

---

## **8.3 ORGANIZATION OF DATA**

---

There are four major statistical techniques for organising the data. These are:

- i) Classification
- ii) Tabulation
- iii) Graphical Presentation, and
- iv) Diagrammatical Presentation

### **8.3.1 Classification**

The arrangement of data in groups according to similarities is known as classification. A classification is a summary of the frequency of individual scores or ranges of scores for a variable. In the simplest form of a distribution, we will have such value of variable as well as the number of persons who have had each value. Once data are collected, it should be arranged in a format from which they would be able to draw some conclusions. Thus by classifying data, the investigators move a step ahead in regard to making a decision. A much clear picture of the

information of score emerges when the raw data are organised as a frequency distribution. Frequency distribution shows the number of cases following within a given class interval or range of scores. A frequency distribution is a table that shows each score as obtained by a group of individuals and how frequently each score occurred.

### **8.3.1.1 Frequency Distribution can be with Ungrouped Data and Grouped Data**

i) An ungrouped frequency distribution may be constructed by listing all score values either from highest to lowest or lowest to highest and placing a tally mark (/) besides each scores every times it occurs. The frequency of occurrence of each score is denoted by 'f'.

ii) Grouped frequency distribution: If there is a wide range of score value in the data, then it is difficult to get a clear picture of such series of data. In this case grouped frequency distribution should be constructed to have a clear picture of the data. A group frequency distribution is a table that organises data into classes.

It shows the number of observations from the data set that fall into each of the class.

#### **Construction of frequency distribution**

To prepare a frequency distribution it is essential to determine the following:

- 1) The range of the given data =, the difference between the highest and lowest scores.
- 2) The number of class intervals = There is no hard and fast rules regarding the number of classes into which data should be grouped. If there are very few scores it is useless to have a large number of class-intervals. Ordinarily, the number of classes should be between 5 to 30.

3) Limits of each class interval = Another factor used in determining the number of classes is the size/ width or range of the class which is known as 'class interval' and is denoted by 'i'. Class interval should be of uniform width resulting in the same-size classes of frequency distribution. The width of the class should be a whole number and conveniently divisible by 2, 3, 5, 10, or 20. There are three methods for describing the class limits for distribution:

(i) Exclusive method, (ii) Inclusive method and (iii) True or actual class method.

- i) Exclusive method In this method of class formation, the classes are so formed that the upper limit of one class become the lower limit of the next class. In this classification, it is presumed that score equal to the upper limit of the class is exclusive, i.e., a score of 40 will be included in the class of 40 to 50 and not in a class of 30 to 40 (30-40, 40-50, 50-60)
- ii) Inclusive method In this method the classes are so formed that the upper limit of one class does not become the lower limit of the next class. This classification includes scores, which are equal to the upper limit of the class. Inclusive method is preferred when measurements are given in whole numbers. (30-39, 40-49, 50-59)
- iii) True or Actual class method Mathematically, a score is an internal when it extends from 0.5 units below to 0.5 units above the face value of the score on a continuum. These class limits are known as true or actual class limits. (29.5 to 39.5, 39.5 to 49.5) etc.

### 8.3.1.2 Types of Frequency Distribution

There are various ways to arrange frequencies of a data array based on the requirement of the statistical analysis or the study. A couple of them are discussed below:

- i) Relative frequency distribution: A relative frequency distribution is a distribution that indicates the proportion of the total number of cases observed at each score value or interval of score values.
- ii) Cumulative frequency distribution: Sometimes investigator may be interested to know the number of observations less than a particular value. This is possible by computing the cumulative frequency. A cumulative frequency corresponding to a class-interval is the sum of frequencies for that class and of all classes prior to that class.
- iii) Cumulative relative frequency distribution: A cumulative relative frequency distribution is one in which the entry of any score of class interval expresses that score's cumulative frequency as a proportion of the total number of cases.

### 8.3.2 Tabulation

Frequency distribution can be either in the form of a table or it can be in the form of graph. Tabulation is the process of presenting the classified data in the form of a table. A tabular presentation of data becomes more intelligible and fit for further statistical analysis. A table is a systematic arrangement of classified data in row and columns with appropriate headings and sub-headings. The main components of a table are: i) Table number: When there is more than one table in a particular analysis a table should be marked with a number for their reference and identification. The number should be written in the center at the top of the table. 23 ii) Title of the table: Every table should have an appropriate title, which describes the content of the table. The title should be clear, brief, and self-explanatory. Title of the table should be placed either centrally on the top of the table or just below or after the table number. iii) Caption: Captions are brief and self-explanatory headings for columns. Captions may involve headings and sub-headings. The captions should be placed in the middle of the columns. For example, we can divide students of a class into males and females, rural and urban, high SES and Low SES

etc. iv) Stub: Stubs stand for brief and self-explanatory headings for rows. v) Body of the table: This is the real table and contains numerical information or data in different cells. This arrangement of data remains according to the description of captions and stubs. vi) Head note: This is written at the extreme right hand below the title and explains the unit of the measurements used in the body of the tables. vii) Footnote: This is a qualifying statement which is to be written below the table explaining certain points related to the data which have not been covered in title, caption, and stubs. viii) Source of data: The source from which data have been taken is to be mentioned at the end of the table.

## TITLE

<b>Stub Head</b>	<b>Caption</b>			
<b>Stub Entries</b>	<b>Column Head I</b>		<b>Column Head II</b>	
	<b>Sub Head</b>	<b>Sub Head</b>	<b>Sub Head</b>	<b>Sub Head</b>
	<b>MAIN BODY</b>	<b>OF</b>	<b>THE TABLE</b>	
<b>Total</b>				

Footnote(s):

Source :

### 8.3.3 Graphical Presentation of Data

The purpose of preparing a frequency distribution is to provide a systematic way of “ looking at” and understanding data. To extend this understanding, the information contained in a frequency distribution often is displayed in graphic and/or diagrammatic forms. In graphical presentation of frequency distribution, frequencies are plotted on a pictorial platform formed of horizontal and vertical lines known as graph. A graph is created on two mutually perpendicular lines called the X and Y–axes on which appropriate scales are indicated. The horizontal

## Notes

line is called the abscissa and vertical the ordinate. Like different kinds of frequency distributions there are many kinds of graph too, which enhance the scientific understanding of the reader. The commonly used graphs are Histogram, Frequency polygon, Frequency curve, Cumulative frequency curve. Here we will discuss some of the important types of graphical patterns used in statistics. i) Histogram: It is one of the most popular methods for presenting continuous frequency distribution in a form of graph. In this type of distribution the upper limit of a class is the lower limit of the following class. The histogram consists of series of rectangles, with its width equal to the class interval of the variable on horizontal axis and the corresponding frequency on the vertical axis as its heights. ii) Frequency polygon: Prepare an abscissa originating from 'O' and ending to 'X'. Again construct the ordinate starting from 'O' and ending at 'Y'. Now label the class-intervals on abscissa stating the exact limits or midpoints of the class intervals. You can also add one extra limit keeping zero frequency on both side of the class-interval range. The size of measurement of small squares on graph paper depends upon the number of classes to be plotted. Next step is to plot the frequencies on ordinate using the most comfortable measurement of small squares depending on the range of whole distribution. To plot a frequency polygon you have to mark each frequency against its concerned class on the height of its respective ordinate. After putting all frequency marks a draw a line joining the points. This is the polygon. iii) Frequency curve: A frequency curve is a smooth free hand curve drawn through frequency polygon. The objective of smoothing of the frequency polygon is to eliminate as far as possible the random or erratic fluctuations that are present in the data.

### 8.3.3.1 Cumulative Frequency Curve or Ogive

The graph of a cumulative frequency distribution is known as cumulative frequency curve or ogive. Since there are two types of cumulative frequency distribution e.g., "less than" and "more than" cumulative frequencies. We can have two types of ogives. i) 'Less than' Ogive: In 'less than' ogive, the less than cumulative frequencies are plotted against

the upper class boundaries of the respective classes. It is an increasing curve having slopes upwards from left to right. ii) 'More than' Ogive : In more than ogive , the more than cumulative frequencies are plotted against the lower class boundaries of the respective classes. It is decreasing curve and slopes downwards from left to right.

### 8.3.4 Diagrammatic Presentation of Data

A diagram is a visual form for the presentation of statistical data. They present the data in simple , readily comprehensible form. Diagrammatic presentation is used only for presentation of the data in visual form, whereas graphic presentation of the data can be used for further analysis. There are different forms of diagram e.g., Bar diagram, Sub-divided bar diagram, Multiple bar diagram, Pie diagram and Pictogram.

- i) Bar diagram: Bar diagram is most useful for categorical data. A bar is defined as a thick line. Bar diagram is drawn from the frequency distribution table representing the variable on the horizontal axis and the frequency on the vertical axis. The height of each bar will be corresponding to the frequency or value of the variable.
- ii) Sub- divided bar diagram: Study of sub classification of a phenomenon can be done by using sub-divided bar diagram. Corresponding to each sub-category of the data the bar is divided and shaded. There will be as many shades as there will sub portion in a group of data. The portion of the bar occupied by each sub-class reflects its proportion in the total.
- iii) Multiple Bar diagram: This diagram is used when comparisons are to be shown between two or more sets of interrelated phenomena or variables. A set of bars for person, place or related phenomena are drawn side by side without any gap. To distinguish between the different bars in a set, different colours, shades are used.
- iv) Pie diagram: It is also known as angular diagram. A pie chart or diagram is a circle divided into component sectors

corresponding to the frequencies of the variables in the distribution. Each sector will be proportional to the frequency of the variable in the group. A circle represents 360°. So 360° angles is divided in proportion to percentages. The degrees represented by the various component parts of given magnitude can be obtained by using this formula. After the calculation of the angles for each component, segments are drawn in the circle in succession, corresponding to the angles at the center for each segment. Different segments are shaded with different colours, shades or numbers.

---

## 8.4 SUMMARIZATION OF DATA

---

In the previous section we have discussed about tabulation of the data and its representation in the form of graphical presentation. In research, comparison between two or more series of the same type is needed to find out the trends of variables. For such comparison, tabulation of data is not sufficient and it is further required to investigate the characteristics of data. The frequency distribution of obtained data may differ in two ways, first in measures of central tendency and second, in the extent to which scores are spread over the central value. Both types of differences are the components of summary statistics

### 8.4.1 Measures of Central Tendency

It is the middle point of a distribution. Tabulated data provides the data in a systematic order and enhances their understanding. Generally, in any distribution values of the variables tend to cluster around a central value of the distribution. This tendency of the distribution is known as central tendency and measures devised to consider this tendency is known as measures of central tendency. A measure of central tendency is useful if it represents accurately the distribution of scores on which it is based. A good measure of central tendency must possess the following characteristics: It should be clearly defined- The definition of a measure of central tendency should be clear and unambiguous so that it leads to



one and only one information. It should be readily comprehensible and easy to compute. It should be based on all observations- A good measure of central tendency should be based on all the values of the distribution of scores. It should be amenable for further mathematical treatment. It should be least affected by the fluctuation of sampling. In Statistics there are three most commonly used measures of central tendency. These are:

1) Arithmetic Mean

2) Median, and

3) Mode

1) Arithmetic Mean: The arithmetic mean is most popular and widely used measure of central tendency. Whenever we refer to the average of data, it means we are talking about its arithmetic mean. This is obtained by dividing the sum of the values of the variable by the number of values. It is also a useful measure for further statistics and comparisons among different data sets. One of the major limitations of arithmetic mean is that it cannot be computed for open-ended class-intervals.

2) Median: Median is the middle most value in a data distribution. It divides the distribution into two equal parts so that exactly one half of the observations is below and one half is above that point. Since median clearly denotes the position of an observation in an array, it is also called a position average. Thus more technically, median of an array of numbers arranged in order of their magnitude is either the middle value or the arithmetic mean of the two middle values. It is not affected by extreme values in the distribution.

3) Mode: Mode is the value in a distribution that corresponds to the maximum concentration of frequencies. It may be regarded as the most typical of a series value. In more simple words, mode is the point in the distribution comprising maximum frequencies therein.

## 8.4.2 Measures of Dispersion

In the previous section we have discussed about measures of central tendency. By knowing only the mean, median or mode, it is not possible to have a complete picture of a set of data. Average does not tell us about how the score or measurements are arranged in relation to the center. It is possible that two sets of data with equal mean or median may differ in terms of their variability. Therefore, it is essential to know how far these observations are scattered from each other or from the mean. Measures of these variations are known as the 'measures of dispersion'. The most commonly used measures of dispersion are range, average deviation, quartile deviation, variance and standard deviation.

- i) **Range** Range is one of the simplest measures of dispersion. It is designated by 'R'. The range is defined as the difference between the largest score and the smallest score in the distribution. It gives the two extreme values of the variable. A large value of range indicates greater dispersion while a smaller value indicates lesser dispersion among the scores. Range can be a good measure if the distribution is not much skewed.
- ii) **Average deviation** Average deviation refers to the arithmetic mean of the differences between each score and the mean. It is always better to find the deviation of the individual observations with reference to a certain value in the series of observation and then take an average of these deviations. This deviation is usually measured from mean or median. Mean, however, is more commonly used for this measurement. Merits: It is less affected by extreme values as compared to standard deviation. It provides better measure for comparison about the formation of different distributions.
- iii) **Standard deviation** Standard deviation is the most stable index of variability. In standard deviation, instead of the actual values of the deviations we consider the squares of deviations and the outcome is known as variance. Further, the square

root of this variance is known as standard deviation and designated as SD. Thus, standard deviation is the square root of the mean of the squared deviations of the individual observations from the mean. The standard deviation of the sample ( $\hat{\sigma}$ ) and population denoted by ( $\hat{\sigma}$ ) respectively. If all the score have an identical value in a sample, the SD will be 0 (zero). Merits: It is based on all observations. It is amenable to further mathematical treatments. Of all measures of dispersion, standard deviation is least affected by fluctuation of sampling.

### 8.4.3 Skewness and Kurtosis

There are two other important characteristics of frequency distribution that provide useful information about its nature. They are known as skewness and kurtosis. i) Skewness Skewness is the degree of asymmetry of the distribution. In some frequency distributions scores are more concentrated at one end of the scale. Such a distribution is called a skewed distribution. Thus, skewness refers to the extent to which a distribution of data points is concentrated at one end or the other. Skewness and variability are usually related, the more the skewness the greater the variability. ii) Kurtosis The term 'kurtosis' refers to the 'peakedness' or flatness of a frequency distribution curve when compared with normal distribution curve. The kurtosis of a distribution is the curvedness or peakedness of the graph. If a distribution is more peaked than normal it is said to be leptokurtic. This kind of peakedness implies a thin distribution. On the other hand, if a distribution is more flat than the normal distribution it is known as Platykurtic distribution. A normal curve is known as mesokurtic.

### 8.4.4 Advantages and Disadvantages of Descriptive Statistics

The Advantages of Descriptive statistics are given below:

## Notes

- It is essential for arranging and displaying data.
- It forms the basis of rigorous data analysis.
- It is easier to work with, interpret, and discuss than raw data.
- It helps in examining the tendencies, variability, and normality of a data set.
- It can be rendered both graphically and numerically.
- It forms the basis for more advanced statistical methods. The disadvantages of descriptive statistics can be listed as given below:
- It can be misused, misinterpreted, and incomplete.
- It can be of limited use when samples and populations are small.
- It offers little information about causes and effects.
- It can be dangerous if not analysed completely.
- There is a risk of distorting the original data or losing important detail.

### Self Assessment Questions

1) Which one of the alternative is appropriate for descriptive statistics?

i) In a sample of school children, the investigator found an average IQ was

110.

ii) A class teacher calculates the class average on their final exam. Was

64%.

2) State whether the following statements are True (T) or False (F).

- i) Mean is affected by extreme values ( )
- ii) Mode is affected by extreme values ( )
- iii) Mode is useful in studying qualitative facts such as intelligence ( )
- iv) Median is not affected by extreme values ( )
- v) Range is most stable measures of variability ( )
- vi) Standard deviation is most suitable measures of dispersion ( )
- vii) Skewness is always positive ( )

**Check Your Progress 1**

Note: a) Use the space provided for your answer.

b) Check your answers with those provided at the end of the unit.

1. What is the Meaning of Descriptive Statistics?

.....  
.....  
.....  
.....  
.....

2. Discuss about Organization of Data

.....  
.....  
.....  
.....  
.....

3. What are the Measures of Central Tendency

.....  
.....  
.....

.....  
.....

4. Discuss the Measures of Dispersion

.....  
.....  
.....  
.....  
.....

---

## 8.5 MEANING OF INFERENTIAL STATISTICS

---

In the previous section we discussed about descriptive statistics, which basically describes some characteristics of data. But the description or definition of the distribution or observations is not the prime objective of any scientific investigation. Organising and summarising data is only one step in the process of analysing the data. In any scientific investigation either the entire population or a sample is considered for the study. In most of the scientific investigations a sample, a small portion of the population under investigation, is used for the study. On the basis of the information contained in the sample we try to draw conclusions about the population. This process is known as statistical inference. Statistical inference is widely applicable in behavioural sciences, especially in psychology. For example, before the Lok Sabha or Vidhan Sabha election process starts or just before the declaration of election results print media and electronic media conduct exit poll to predict the election result. In this process all voters are not included in the survey, only a portion of voters i.e. sample is included to infer about the population. This is called inferential statistics. Inferential statistics deals with drawing of conclusions about large group of individuals ( population) on the basis of observation of a few participants from among them or about the events which are yet to occur on the basis of past events. It provides tools to compute the probabilities of future behaviour of the subjects. Inferential statistics is the mathematics and logic of how this

generalisation from sample to population can be made. There are two types of inferential procedures: (1) Estimation, (2) Hypothesis testing

### 8.5.1 Estimation

In estimation, inference is made about the population characteristics on the basis of what is discovered about the sample. There may be sampling variations because of chance fluctuations, variations in sampling techniques, and other sampling errors. Estimation about population characteristics may be influenced by such factors. Therefore, in estimation the important point is that to what extent our estimate is close to the true value. Characteristics of Good Estimator: A good statistical estimator should have the following characteristics, (i) Unbiased (ii) Consistent (iii) Accuracy . These are being dealt with in detail below. i) Unbiased An unbiased estimator is one in which, if we were to obtain an infinite number of random samples of a certain size, the mean of the statistic would be equal to the parameter. The sample mean, ( $\bar{x}$ ) is an unbiased estimate of population mean ( $\mu$ ) because if we look at possible random samples of size  $N$  from a population, then mean of the sample would be equal to  $\mu$ . ii) Consistent A consistent estimator is one that as the sample size increased, the probability that estimate has a value close to the parameter also increased. Because it is a consistent estimator, a sample mean based on 20 scores has a greater probability of being closer to ( $\mu$ ) than does a sample mean based upon only 5 scores ii) Accuracy The sample mean is an unbiased and consistent estimator of population mean ( $\mu$ ). But we should not over look the fact that an estimate is just a rough or approximate calculation. It is unlikely in any estimate that ( $\bar{x}$ ) will be exactly equal to population mean ( $\mu$ ). Whether or not  $\bar{x}$  is a good estimate of ( $\mu$ ) depends upon the representative ness of sample, the sample size, and the variability of scores in the population.

### 8.5.2 Point Estimation

We have indicated that  $\bar{x}$  obtained from a sample is an unbiased and consistent estimator of the population mean ( $\mu$ ). Thus, if an investigator

obtains Adjustment Score from 100 students and wanted to estimate the value of  $(\mu)$  for the population from which these scores were selected, researcher would use the value of  $x$  as an estimate of population mean  $(\mu)$ . If the obtained value of  $x$  were 45.0 then this value would be used as estimate of population mean  $(\mu)$ . This form of estimate of population parameters from sample statistic is called point estimation. Point estimation is estimating the value of a parameter as a single point, for example, population mean  $(\mu) = 45.0$  from the value of the statistic  $x = 45.0$

### 8.5.3 Interval Estimation

A point estimate of the population mean  $(\mu)$  almost is assured of being in error, the estimate from the sample will not equal to the exact value of the parameter. To gain confidence about the accuracy of this estimate we may also construct an interval of scores that is expected to include the value of the population mean. Such intervals are called confidence interval. A confidence interval is a range of scores that is expected to contain the value of  $(\mu)$ . The lower and upper scores that determine the interval are called confidence limits. A level of confidence can be attached to this estimate so that, the researcher can be 95% or 99% confidence level that encompasses the population mean.

---

## 8.6 HYPOTHESIS TESTING

---

Inferential statistics is closely tied to the logic of hypothesis testing. In hypothesis testing we have a particular value in mind. We hypothesize that this value characterise the population of observations. The question is whether that hypothesis is reasonable in the light of the evidence from the sample. In estimation no particular population value need be stated. Rather, the question is , what is the population value. For example, Hypothesis testing is one of the important areas of statistical analyses. Sometimes hypothesis testing is referred to as statistical decision-making process. In day-to-day situations we are required to take decisions about the population on the basis of sample information.



### 8.6.1 Statement of Hypothesis

A statistical hypothesis is defined as a statement, which may or may not be true about the population parameter or about the probability distribution of the parameter that we wish to validate on the basis of sample information. Most times, experiments are performed with random samples instead of the entire population and inferences drawn from the observed results are then generalised over to the entire population. But before drawing inferences about the population it should be always kept in mind that the observed results might have come due to chance factor. In order to have an accurate or more precise inference, the chance factor should be ruled out. The probability of chance occurrence of the observed results is examined by the null hypothesis ( $H_0$ ). Null hypothesis is a statement of no differences. The other way to state null hypothesis is that the two samples came from the same population. Here, we assume that population is normally distributed and both the groups have equal means and standard deviations. Since the null hypothesis is a testable proposition, there is counter proposition to it known as alternative hypothesis and denoted by  $H_1$ . In contrast to null hypothesis, the alternative hypothesis ( $H_1$ ) proposes that

- i) the two samples belong to two different populations,
- ii) their means are estimates of two different parametric means of the respective population, and
- iii) there is a significant difference between their sample means.

The alternative hypothesis ( $H_1$ ) is not directly tested statistically; rather its acceptance or rejection is determined by the rejection or retention of the null hypothesis. The probability 'p' of the null hypothesis being correct is assessed by a statistical test. If probability 'p' is too low,  $H_0$  is rejected and  $H_1$  is accepted. It is inferred that the observed difference is significant. If probability 'p' is high,  $H_0$  is accepted and it is inferred that the difference is due to the chance factor and not due to the variable factor.

### **8.6.2 Level of Significance**

The level of significance ( $p$ ). The selection of level of significance depends on the choice of the researcher. Generally level of significance is taken to be 5% or 1%, i.e., = .05 or = .01). If null hypothesis is rejected at .05 level, it means that the results are considered significant so long as the probability 'p' of getting it by mere chance of random sampling works out to be 0.05 or less ( $p < .05$ ). In other words, the results are considered significant if out of 100 such trials only 5 or less number of the times the observed results may arise from the accidental choice in the particular sample by random sampling.

### **8.6.3 One Tail and Two Tail Test**

Depending upon the statement in alternative hypothesis ( $H_1$ ), either a one-tail or two-tail test is chosen for knowing the statistical significance. A one-tail test is a directional test. It is formulated to find the significance of both the magnitude and the direction (algebraic sign) of the observed difference between two statistics. Thus, in two-tailed tests researcher is interested in testing whether one sample mean is significantly higher (alternatively lower) than the other sample mean.

---

## **8.7 ERRORS IN HYPOTHESIS TESTING**

---

In hypothesis testing, there would be no errors in decision making as long as a null hypothesis is rejected when it is false and also a null hypothesis is accepted when it is true. But the decision to accept or reject the null hypothesis is based on sample data. There is no testing procedure that will ensure absolutely correct decision on the basis of sampled data. There are two types of errors regarding decision to accept or to reject a null hypothesis.

### **8.7.1 Type I Error**

When the null hypothesis is true, a decision to reject it is an error and this kind of error is known as type I error in statistics. The probability of making a type I error is denoted as ' $\alpha$ ' (read as alpha). The null hypothesis is rejected if the probability ' $p$ ' of its being correct does not exceed the  $p$ . The higher the chosen level of  $p$  for considering the null hypothesis, the greater is the probability of type I error.

### 8.7.2 Type II Error

When null hypothesis is false, a decision to accept it is known as type II error. The probability of making a type II error is denoted as ' $\beta$ ' (read as beta). The lower the chosen level of significance  $p$  for rejecting the null hypothesis, the higher is the probability of the type II error. With a lowering of  $p$ , the rejection region as well as the probability of the type I error declines and the acceptance region  $(1-p)$  widens correspondingly. The goodness of a statistical test is measured by the probability of making a type I or type II error. For a fixed sample size  $n$ ,  $\alpha$  and  $\beta$  are so related that reduction in one causes increase in the other. Therefore, simultaneous reductions in  $\alpha$  and  $\beta$  are not possible. If  $n$  is increased, it is possible to decrease both  $\alpha$  and  $\beta$ .

### 8.7.3 Power of a Test

The probability of committing type II error is designated by  $\beta$ . Therefore,  $1-\beta$  is the probability of rejecting null hypothesis when it is false. This probability is known as the power of a statistical test. It measures how well the test is working. The probability of type II error depends upon the true value of the population parameter and sample size  $n$ .

---

## 8.8 GENERAL PROCEDURE FOR TESTING A HYPOTHESIS

---

Step 1. Set up a null hypothesis suitable to the problem.

Step 2. Define the alternative hypothesis.

## Notes

Step 3. Calculate the suitable test statistics.

Step 4. Define the degrees of freedom for the test situation.

Step 5. Find the probability level 'p' corresponding to the calculated value of the test statistics and its degree of freedom. This can be obtained from the relevant tables.

Step 6. Reject or accept null hypothesis on the basis of tabulated value and calculated value at practical probability level. There are some situations in which inferential statistics is carried out to test the hypothesis and draw conclusion about the population , for example (i) Test of hypothesis about a population mean (Z test), (ii) Testing hypothesis about a population mean (small sample ' t' test).

### Check Your Progress 2

Note: a) Use the space provided for your answer.

b) Check your answers with those provided at the end of the unit.

1. What is Skewness and Kurtosis?

.....  
.....  
.....

2. Statement of Hypothesis

.....  
.....  
.....

3. Level of Significance

.....  
.....  
.....

4. One Tail and Two Tail Test

.....  
 .....  
 .....

---

## 8.9 LET US SUM UP

---

Descriptive statistics are used to describe the basic features of the data in investigation. Such statistics provide summaries about the sample and measures. Data description comprises two operations: organising data and describing data. Organising data includes: classification, tabulation, graphical and diagrammatic presentation of raw scores. Whereas, measures of central tendency and measures of dispersion are used in describing the raw scores. In the above section, the basic concepts and general procedure involved in inferential statistics are also discussed. Inferential statistics is about inferring or drawing conclusions from the sample to population. This process is known as statistical inference. There are two types of inferential procedures: estimation and hypothesis testing. An estimate of unknown parameter could be either point or interval. Hypothesis is a statement about a parameter. There are two types of hypotheses: null and alternative hypotheses. Important concepts involved in the process of hypothesis testing example, level of significance, one tail test, two tail test, type I error, type II error, power of a test are explained. General procedure for hypothesis testing is also given.

---

## 8.10 KEY WORDS

---

Classification : A systematic grouping of data

Cumulative frequency : A classification, which shows the cumulative distribution frequency below, the upper real limit of the corresponding class interval.

Data : Any sort of information that can be analysed.

Discrete data : When data are counted in a classification.

Exclusive classification : The classification system in which the upper limit of the class becomes the lower limit of next class

## Notes

Frequency distribution : Arrangement of data values according to their magnitude.

Inclusive classification : When the lower limit of a class differs the upper limit of its successive class.

Mean : The ratio between total and numbers of scores.

Median : The mid point of a score distribution.

Mode : The maximum occurring score in a score distribution.

Central Tendency : The tendency of scores to bend towards center of distribution.

Dispersion : The extent to which scores tend to scatter from their mean and from each other.

Standard Deviation : The square root of the sum of squared deviations of scores from their mean.

Skewness : Tendency of scores to polarize on either side of abscissa.

Kurtosis : Curvedness of a frequency distribution graph.

Range : Difference between the two extremes of a score distribution.

Confidence level : It gives the percentage (probability) of samples where the population mean would remain within the confidence interval around the sample mean.

Estimation : It is a method of prediction about parameter value on the basis Statistic.

Hypothesis testing: The statistical procedures for testing hypotheses.

Level of significance: The probability value that forms the boundary between rejecting and not rejecting the null hypothesis.

Null hypothesis: The hypothesis that is tentatively held to be true (symbolised by  $H_0$ )

One-tail test : A statistical test in which the alternative hypothesis specifies direction of the departure from what is expected under the null hypothesis.

Parameter : It is a measure of some characteristic of the population.

Population : The entire number of units of research interest.

Power of a test : An index that reflects the probability that a statistical test will correctly reject the null hypothesis relative to the size of the sample involved.

Sample : A sub set of the population under study.

Descriptive and Inferential Statistics Introduction to Statistics

Statistical inference : It is the process of concluding about an unknown population from known sample drawn from it.

Statistical hypothesis : The hypothesis which may or may not be true about the population parameter.

t-test : It is a parametric test for the significance of differences between means.

Type I error : A decision error in which the statistical decision is to reject the null hypothesis when it is actually true.

Type II error : A decision error in which the statistical decision is not to reject the null hypothesis when it is actually false.

Two-tail test : A statistical test in which the alternative hypothesis does not specify the direction of departure from what is expected under the null hypothesis.

---

## 8.11 QUESTIONS FOR REVIEW

---

1. What is descriptive statistics? Discuss its advantages and disadvantages.
2. What do you mean by organisation of data? State different methods of organising raw data.
3. Define measures of dispersion. Why it is that standard deviation is considered as the best measures of variability?
4. Explain the importance of inferential statistics.
5. Describe the important properties of good estimators.
6. Discuss the different types of hypothesis formulated in hypothesis testing.
7. Discuss the errors involved in hypothesis testing.
8. Explain the various steps involved in hypothesis testing.
9. What is statistical inference?
10. What are the procedures involved in statistical inference?

---

## 8.12 SUGGESTED READINGS AND REFERENCES

---

## Notes

- Asthana, H. S. and Bhushan, B. (2007). Statistics for Social Sciences ( with SPSS Application). Prentice Hall of India, New Delhi.
- Garret, H. E. (2005). Statistics in Psychology and Education. Jain publishing, India.
- Elhance, D. N., and Elhance, V. (1988). Fundamentals of Statistics. Kitab Mahal, Allahabad
- Nagar, A. L., and Das, R. K. (1983). Basic Statistics. Oxford University Press, Delhi.
- Sani, F., and Todman, J. (2006). Experimental Design and Statistics for Psychology. Blackwell Publishing U.K.
- Yale, G. U., and M.G. Kendall (1991). An Introduction to the Theory of Statistics. Universal Books, Delhi.

---

## 8.13 ANSWERS TO CHECK YOUR PROGRESS

---

### Check Your Progress 1

1. See Section 8.2
2. See Section 8.3
3. See Sub Section 8.4.1
4. See Sub Section 8.4.2

### Check Your Progress 2

1. See Sub Section 8.4.3
2. See Sub Section 8.6.1
3. See Sub Section 8.6.2
4. See Sub Section 8.6.3



---

# UNIT 9: CORRELATION AND REGRESSION

---

## STRUCTURE

- 9.0 Objectives
- 9.1 Introduction
- 9.2 Correlation
- 9.3 Method of Calculating Correlation of Ungrouped Data
- 9.4 Method of Calculating Correlation of Grouped Data
- 9.5 Regression
- 9.6 Let us sum up
- 9.7 Key Words
- 9.8 Questions for Review
- 9.9 Suggested readings and references
- 9.10 Answers to Check Your Progress

---

## 9.0 OBJECTIVES

---

After studying Unit 9, you should be able to:

- Discuss the relevance of the analysis of co-variation between two or more variables;
- Describe different types of correlation;
- Elaborate methods of calculating correlation of both ungrouped and grouped data; and
- Understand the method of regression analysis that helps in estimating the values of a variable from the knowledge of one or more variables.

---

## 9.1 INTRODUCTION

---

Unit 9 is about correlation that is an analysis of co-variation between two or more variables. You would notice that the statistical tool of correlation helps to measure and express the quantitative relationship between two variables. Unit 9 elaborates the ways of applying the tool. It shows the relevance of coefficient of correlation, coefficient of determination and

## Notes

regression analysis in the social sciences. Further, it explains regression analysis, which is the method of estimating the values of a variable from the knowledge of one or more variables. The Unit tells you to use the statistical tool of correlation without fear or apprehension that its application is difficult and complex.

So far we have been dealing with the distributions of the data involving only one variable. Such a distribution is called a univariate distribution. Very often, we have to deal with the situations where more than one variable are involved. For example, we may like to study the relationship between the heights and weights of adult males, quantum of rainfall and the yield of wheat in India over a number of years, doses of drug and a response vi2.a dose of insulin and blood sugar levels in a person, the age of individuals and their blood pressure, etc. In such situations, our main purpose is to determine whether or not a relationship exists between the two variables. If such a relationship can be expressed by a mathematical formula, then we shall be able to use it for an analysis and hence make certain predictions. Correlation and regression are methods that deal with the analysis of such relationships between various variables and possible predictions. In this unit, we shall confine ourselves to analysing the linear relationship between two variables. However, we can extend the methods for two variables to the situations where more than two variables are studied simultaneously.

In statistics, dependence or association is any statistical relationship, whether causal or not, between two random variables or bivariate data. In the broadest sense correlation is any statistical association, though it commonly refers to the degree to which a pair of variables is linearly related. Familiar examples of dependent phenomena include the correlation between the physical statures of parents and their offspring, and the correlation between the demand for a limited supply product and its price.

Correlations are useful because they can indicate a predictive relationship that can be exploited in practice. For example, an electrical utility may

produce less power on a mild day based on the correlation between electricity demand and weather. In this example, there is a causal relationship, because extreme weather causes people to use more electricity for heating or cooling. However, in general, the presence of a correlation is not sufficient to infer the presence of a causal relationship (i.e., correlation does not imply causation).

Formally, random variables are dependent if they do not satisfy a mathematical property of probabilistic independence. In informal parlance, correlation is synonymous with dependence. However, when used in a technical sense, correlation refers to any of several specific types of mathematical operations between the tested variables and their respective expected values. Essentially, correlation is the measure of how two or more variables are related to one another. There are several correlation coefficients, often denoted  $\rho$  or  $r$ , measuring the degree of correlation. The most common of these is the Pearson correlation coefficient, which is sensitive only to a linear relationship between two variables (which may be present even when one variable is a nonlinear function of the other). Other correlation coefficients have been developed to be more robust than the Pearson correlation – such as the Spearman's rank correlation that is, more sensitive to nonlinear relationships. Mutual information can also be applied to measure dependence between two variables.

---

## 9.2 CORRELATION

---

In studying the linear relationship between two variables, we try to examine the question "Are the two variables mutually related to each other?" In other words, we may ask whether the changes in one variable are accompanied by some corresponding changes in the other variable. For example, to find the relationship between the heights and weights of 100 persons, we can arrange them in increasing order of their heights and see whether or not the weight increases as the height increases. In other words, we are asking, "Do taller people tend to weigh more than shorter people?" Note carefully that we are not saying that if an individual is

## Notes

taller than another, he has to necessarily weigh heavier. Very seldom, a taller person may weigh less than a shorter person, but quite often taller persons have higher weights than shorter persons. That is, in general, we may expect to see that as the heights of 100 individuals are arranged in increasing order and the corresponding weights written down the weights will show a tendency to increase.

In such a situation, two variables, then, are said to be mutually related or correlated. This process of mutual relationship is called correlation between two variables. Note that correlation need not be only in one direction. As one variable shows an increase, the second variable may show an increase or a decrease. We know, for example, as the altitudes of places increase, the atmospheric pressure decreases. Hence, whenever two variables are related to each other in such a way that change in the one creates a corresponding change in the other, then the variables are said to be correlated. An easy way of studying the correlation of two quantitative variables is to plot them on a graph sheet taking one of the variables on the X-axis and the other on the Y-axis. The resulting diagram is called a scatter diagram because it shows how the pairs of observations are scattered on the graph sheet. Note that the points representing the values of  $x$  and  $y$  may lie very close to a straight line. This means that we can approximate the relationship between the values of  $x$  and  $y$  by a straight line or by some other geometrical curve. If this is a straight line, then we say that the relationship between  $x$  and  $y$  is linear. The relationship between  $x$  and  $y$  may be a curve other than a straight line. The study of such relationships is beyond the scope of the present syllabus. Hence, we shall confine our discussion to the linear relationship between  $x$  and  $y$ .

Correlation is an analysis of the co-variation between two or more variables. When the relationship between the two variables is quantitative, the statistical tool for measuring the relationship and expressing it in a brief formula is known as correlation. If a change in one variable results in a corresponding change in the other, the two variables are correlated. Let us look at types of correlation.

## Types of Correlation

Probing into the types of correlation, we contemplate two types of correlation:

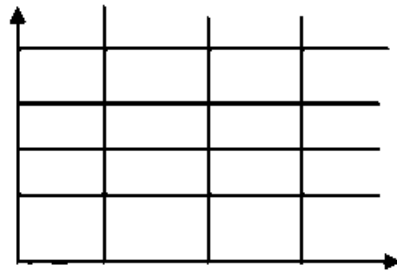
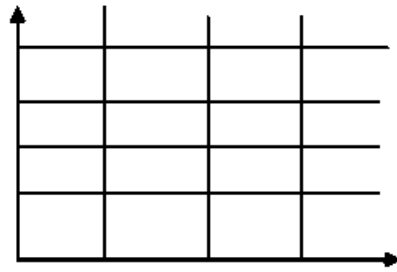
A) Positive and Negative correlation;

B) Linear and Non-linear correlation

A) Positive and negative correlation If the values of the two variables deviate in the same direction, i.e., if an increase in the value of one results on an average in a corresponding increase in the value of the other, or if decrease in the value of one variable results in a decrease in the value of the other, then correlation is said to be positive or direct. Some examples of a series of positive correlation are (i) height and weight (ii) land owned and household income. On the other hand, if the variables deviate in the opposite directions, i.e. if an increase (decrease) in the value of one variable, on an average, results in a decrease (increase) in the value of the other variable, then the correlation is negative or indirect. Some examples of negative correlation are (i) physical assets and the level of poverty, (ii) muscle strength and age. Figure 9.1 shows the positive and negative types of correlation

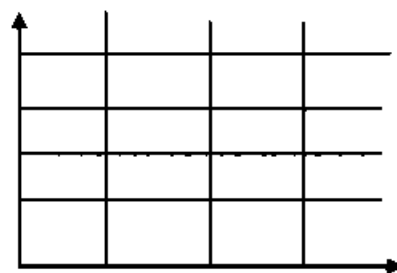
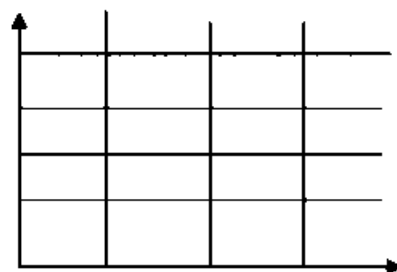
Figure 9.1 (a) Positive Correlation and (b) Negative Correlation

## Notes



The values of correlation range from -1 to +1. When  $r = +1$ , it means there is perfect positive correlation between the variables. When  $r = -1$ , there is perfect negative correlation. When  $r = 0$ , it means there is no correlation between the two variables (see Figure 9.2).

Figure 23.2 (a) Perfect Positive Correlation ( $r = +1$ ) and (b) Perfect Negative Correlation ( $r = -1$ )



## B) Linear and non-linear correlation

The correlation between two variables is said to be linear if corresponding to a unit change in one variable, there is a constant change in the other variable over the entire range of the values. Consider the following data in Figure 9.3.

Figure 9.3 Constant Change Figuring in the Entire Range of Values

X	1	2	3	4	5	6
Y	3	5	7	9	11	13

In this case, the data in Figure 9.3 can be represented by the relation  $Y = 1 + 2X$ . In general, two variables are said to be linearly related if there exists a relationship of the form  $Y = a + bX$ . On the other hand, the relationship between the two variables is said to be non-linear or curvilinear if corresponding to a unit change in one variable, the other variable does not change at a constant but a fluctuating rate. Example of a non-linear correlation is given by the following data set in Figure 9.4.

In the example in Figure 9.4,- there is fluctuating (not constant) change in the value of Y corresponding to a unit change in the value of X, and thus it represents a non-linear correlation. You would like to know how to study correlation. Let us briefly discuss the methods of studying correlation. But before going on to methods of studying correlation, let us complete Reflection and Action 9.1. Relating to your hypothesis, draw the figure of its positive and negative correlations. Next draw another figure of perfect positive and perfect negative correlations. In addition, draw two more figures of constant change reflected in the entire range of values and non-linear correlation. You may take help of Figures 18.1 to 18.4 in the text above for drawing your figures. Methods of studying correlation The various methods to determine whether there is a correlation between two variables are (i)

## Notes

Scatter diagram; (ii) Graphic method; (iii) Karl Pearson's coefficient of correlation; (iv) Rank method; (v) Concurrent deviation method; and (vi) Method of least squares. Of these, the first two are based on the knowledge of diagrams and graphs and the rest on mathematical tools. Of the several mathematical tools used, the most popular is the Karl Pearson coefficient of correlation ( $r$ ) and thus we will focus on this method. The procedure is different for calculating correlation from 6 ungrouped and grouped data.

In the previous section, we saw how a scatter diagram helps us visually to judge whether pairs of values of two variables are correlated with each other. If related, whether positively or negatively. It would, however, be essential to measure the strength of this correlation in order that we may; compare two sets of variables both indicating positive or negative correlation. As a first step in this direction, we define a measure called 'Co-Variance'. You proved that the variance of a variable  $x$  is defined as the average sum of the squares of deviations of the values of the variable from its mean value. Algebraically, it is given by the formula,

---

### **9.3 METHOD OF CALCULATING CORRELATION OF UNGROUPED DATA**

---

There are various methods for the calculation of the coefficient of correlation from ungrouped data. i) Using actual mean ii) Using assumed mean iii) Direct method The use of all these methods is illustrated with the help of the following example. Example: Find out the correlation coefficient (Karl Pearson's) between the age at marriage of husbands and wives using the following data in Figure 9.5

Figure 9.5 Correlation Coefficient between the Age at Marriage of Husbands and Wives



Age at Marriage	Case1	Case2	Case3	Case4	Case5	Case6	Case7	Case8	Case 9	Case 10
Husbands	28	25	24	29	31	22	21	25	26	28
Wives	22	23	21	25	26	20	19	21	21	24

### Method of calculating correlation coefficient using the actual mean

You would first learn the method of calculating correlation coefficient using the actual mean and then you would actually carry out the calculation itself. The formula used for calculating  $r$  is:  $r = \frac{\sum xy}{N \cdot \sigma_X \cdot \sigma_Y}$  Where,  $x = (X - M_X)$  in which  $M_X$  is the mean of series of  $X$  values;  $y = (Y - M_Y)$  in which  $M_Y$  is the mean of series of  $Y$  values;  $\sigma_X =$  Standard deviation of series  $X$

$\sigma_Y =$  Standard deviation of series  $Y$

$N =$  Number of pairs of observations

This formula can also be expressed as:

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

2

\*  $\sigma_Y$

2

]

The following steps elucidate the calculation of the coefficient of correlation.

I Take deviations of  $X$  series from the mean of  $X$  and denote them by  $x$ ;

II. Square these deviations and obtain the total, i.e.  $\sum x^2$

2

;

III. Take deviations of  $Y$  series from the mean of  $Y$  and denote them by

$y$ ;

IV. Square these deviations and obtain the total, i.e.  $\sum y^2$

2

V. Multiply the deviations of  $x$  and  $y$  and obtain the total  $\sum xy$ ; and

VI. Substitute the values of  $\sum x^2$ ,  $\sum y^2$  and  $\sum xy$  in the above formula. Calculation of correlation coefficient using actual mean

## Notes

After learning the method, let us now make the calculation as reflected in Figure 9.6

Figure 9.6 Calculation of Correlation Coefficient using Actual Mean

X	$x=X-M_x$	$X^2$	Y	$y=Y-M_y$	$Y^2$	Xy
28	2.1	04.41	22	-0.2	00.04	-00.42
25	-0.9	00.81	23	0.8	00.64	-00.72
24	-1.9	03.61	21	-1.2	01.44	02.28
29	3.1	09.61	25	2.8	07.84	08.68
31	5.1	26.01	26	3.8	14.44	19.38
22	-3.9	15.21	20	-2.2	04.84	08.58
21	-4.9	24.01	19	-3.2	10.24	15.68
25	-0.9	00.81	21	-1.2	01.44	01.08
26	0.1	00.01	21	-1.2	01.44	-00.12
28	2.1	04.41	24	1.8	03.24	03.78
259	0	88.90	222	0	45.60	58.20

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

$$M_x = 259/10 = 25.9 \quad M_y = 222 / 10 = 22.2 \quad (\sum x^2) = 88.9 \quad (\sum y^2) = 45.6, \\ \sum xy = 58.2 \quad r = 58.2 / \sqrt{[88.9 * 45.6]} = 0.914$$

### Method of calculating correlation coefficient using assumed mean

The only difference in this method as compared to the above method is that in the former, the deviations are taken from the actual mean, and in this case from the assumed mean (i.e. by looking at the series of X and Y, assume means for X and Y and proceeding in the same manner). Calculation of correlation coefficient using assumed mean You would now calculate as per Figure 9.7

Figure 9.7 Calculation of correlation coefficient using assumed mean

X	$D_x = X - A_x$	$d_x^2$	Y	$d_y = Y - A_y$	$d_y^2$	$d_x * d_y$
28	3	9	22	0	0	0
25	0	0	23	1	1	0
24	-1	1	21	-1	1	1
29	4	16	25	3	9	12
31	6	36	26	4	16	24
22	-3	9	20	-2	4	6
21	-4	16	19	-3	9	12
25	0	0	21	-1	1	0
26	1	1	21	-1	1	-1
28	3	9	24	2	4	6
259	9	97	222	2	46	60

$$r = \frac{N \sum d_x * d_y - (\sum d_x * \sum d_y)}{\sqrt{\{N \sum d_x^2 - (\sum d_x)^2\} * \{N \sum d_y^2 - (\sum d_y)^2\}}}$$

$$r = \frac{10 * 60 - (9 * 2)}{\sqrt{\{10 * 97 - (9)^2\} * \{10 * 46 - (2)^2\}}}$$

$$r = \frac{582}{636.697}$$

$$r = 0.914$$

Direct method of calculating correlation coefficient The coefficient can also be calculated by taking actual X and Y values, without taking deviations either from the actual or assumed mean. The formula for its calculation is as follows.

$$r = \frac{N * \sum XY - \sum X * \sum Y}{\sqrt{[N * \sum X^2 - (\sum X)^2] * [N * \sum Y^2 - (\sum Y)^2]}}$$

## Notes

The direct method gives the same answer as one gets when deviations are taken from the assumed or actual means. The example demonstrates this point in Figure 9.8

Figure 9.8 Calculation of correlation coefficient using direct method

X	Y	$X_2$	$Y_2$	XY
28	22	784	484	616
25	23	625	529	575
24	21	576	441	504
29	25	841	625	725
31	26	961	676	806
22	20	484	400	440
21	19	441	361	399
25	21	625	441	525
26	21	676	441	546
28	24	784	576	672
259	222	6797	4974	5808

Let us now complete Reflection and Action 9.2 and then learn in Section 9.4 the methods of calculating correlation of grouped data. Reflection and Action 9.2 Select one of the following two calculations and carry it out in relation to your hypothesis. You need not worry about making mistakes in your calculations. At the moment the idea is to learn the procedure. This is not to be a part of your report. i) Calculation of correlation coefficient using assumed mean ii) Calculation of Correlation Coefficient using Direct Method

---

## 9.4 METHOD OF CALCULATING CORRELATION OF GROUPED DATA

---

With a large number of observations, the data is concealed into a twoway frequency distribution called correlation table. The class intervals of Y series are written as column headings and that of the X series are written

as row headings. The frequency distribution for the two variables is written in the respective cells. The formula for calculating the coefficient of correlation is:

$$r = \frac{\sum f \cdot dX \cdot dY - (\sum fX \cdot dX \cdot \sum fY \cdot dY) / N}{\sqrt{\{\sum fX \cdot dX^2 - (\sum fX \cdot dX)^2 / N\}} \cdot \sqrt{\{\sum fY \cdot dY^2 - (\sum fY \cdot dY)^2 / N\}}}$$

Steps: i) Take the step deviations of variable X and denote these deviations by  $dX$  ii) Take the step deviations of variable Y and denote these deviations by  $dY$  iii) Multiply  $dX \cdot dY$  and the respective frequencies for each cell and write the figure obtained in the right hand upper corn of the cell. iv) Add together all values to obtain  $\sum f \cdot dX \cdot dY$  v) Multiply all the frequencies of the variable X by the deviations of X and obtain the total  $\sum fX \cdot dX$  vi) Take the squares of the deviations of the variable X and multiply by respective frequencies to obtain  $\sum fX \cdot dX^2$  vii) Multiply all the frequencies of the variable Y by the deviations of Y and obtain the total  $\sum fY \cdot dY$  viii) Take the squares of the deviations of the variable Y and multiply by respective frequencies to obtain  $\sum fY \cdot dY^2$  ix) Substitute the values for  $\sum fY \cdot dY^2$ ,  $\sum fY \cdot dY$ ,  $\sum fX \cdot dX^2$ ,  $\sum fX \cdot dX$ ,  $\sum f \cdot dX \cdot dY$  in the above formula to get the value of r. Let us now take an example to calculate the Karl Pearson's coefficient of correlation using the data in Figure 9.9.

Figure 9.9 Coefficient Correlation regarding Expenditure on Luxury Items

Expenditure on Luxury Items	(Income in Thousand Rs.)				
	2 - 25	25 - 30	30 - 35	35 - 40	40 - 45
0-4	28	12	05	-	-
4-8	41	22	09	03	-
8-12	09	33	28	14	16
12-16	-	18	22	29	37
16-20	-	-	03	09	12

## Notes

We can calculate correlation coefficient in grouped data using direct method as seen in Figure 9.10 (See figure 9.10).

Figure 9.10 Calculation of Correlation Coefficient in Grouped Data

Expenditure on Luxury Items Income in Thousand Rs.)									
	20- 25	25 - 30	30 - 35	35 - 40	40-45	fy	DY	fy*dy	fy*dy*d
0-4	28	12	5			45	-2	-90	180
4-8	41	22	9	3		75	-1	-75	75
8- 12	9	33	28	14	16	100	0	0	0
12-16		18	22	29	37	106	1	106	106
16-20			3	9	12	24	2	48	96
<b>Fx</b>	<b>78</b>	<b>85</b>	<b>67</b>	<b>55</b>	<b>65</b>	<b>350</b>		<b>-11</b>	<b>457</b>
<b>dx</b>	<b>-2</b>	<b>-1</b>	<b>0</b>	<b>1</b>	<b>2</b>				
<b>fx*dx</b>	<b>-156</b>	<b>-85</b>	<b>0</b>	<b>55</b>	<b>130</b>	<b>-56</b>			
<b>fx*dx*dx 312</b>		<b>85</b>	<b>0</b>	<b>55</b>	<b>260</b>	<b>712</b>			

Now we can proceed to calculate  $fx*dx*dy$  using direct method as given in Figure 9.11.

Figure 9.11 Calculation of Correlation Coefficient of Grouped Data

Expenditure on Luxury Items	Income in Thousand Rs.)					fx*dx*dy
	20- 25	25 - 30	30 - 35	35 - 40	40-45	
0-4	112	24	0	0	0	136
4-8	82	22	0	-3	0	101
8-12	0	0	0	0	0	0
12-16	0	-18		29	74	85
16 - 20	0	0	0	18	48	66
<b>fx*dx*dy 194</b>	<b>28</b>	<b>0</b>	<b>44</b>	<b>122</b>	<b>388</b>	

$$N = 350 \quad \sum f*dx*dY = 388 \quad \sum fx* dx = -56$$

$$\sum fy*dy = -11 \quad \sum fx*dx^2 = 712 \quad \sum fy*dy^2 = 457$$

$$\sum f*dx *dy - (\sum f x*dx*\sum fy*dy)/N$$

$$r = \frac{\sum f x dx^2 - (\sum f x dx)^2 / N}{\sqrt{\sum f y dy^2 - (\sum f y dy)^2 / N}}$$

$$388 - (-56 * -11/350)$$

$$r = \frac{\sqrt{712 - (-56)^2 / 350} * \sqrt{457 - (-11)^2 / 350}}{r = 386.24 / (26.515 * 21.369) = .682}$$

Most of the variables show some kind of relationship. With the help of correlation one can measure the degree of relationship between two or more variables. Correlation, however, does not tell us anything about the cause and effect relationship. Even a high degree of relationship does not necessarily imply that a cause and effect relationship exists. Conversely, however the cause and effect relationship (or functional relationship) would always result in the expression of correlation. We would now discuss regression analysis.

---

## 9.5 REGRESSION

---

In the previous section you have seen that the data giving the corresponding values of two variables can be graphically represented by a scatter diagram. Also, you were introduced to a method of finding the relationship between these two variables in terms of the correlation coefficient. Very often, in the study of relationship of two variables, we come across instances where one of the two variables depends on the other. In other words, what is the possible value of the dependent variable when the value of independent variable is known? For example, the bodyweight of a growing child depends on the nutrient intake of the child or the weight of an individual may be dependent on his height or the response L. a drug can be dependent on the dose of the drug or the agricultural yield may depend on the quantum of rainfall. In such situations, where one of the variables is dependent and the other independent, you may ask "can we find a method of estimating the numerical relationship between two variables so that given a value of the independent variable, we can predict the average value of the dependent

## Notes

variable?". Note that we are trying to predict or estimate the average value of the dependent variable for a given value of the independent variable. We cannot determine the exact value of the dependent variable when the value of the independent variable is known. What perhaps we can do is just to make an estimation of the value of the dependent variable, knowing fully well that there could be an error in our estimation. This is because of the reason that there is no certainty that the estimated value of the variable would be exactly the same as the value actually observed. This is also because for a given value of the independent variable, the dependent variable will usually show some variations in its values. For example, not all persons of a given height, say of 5' 6" have the same weight. Some will be heavier than others. This is why we talk of predicting the average value of the dependent variable for a given value of the independent variable.

Regression analysis is the method of estimating the values of a variable from the knowledge of one or more variables. The variable that the researcher tries to estimate is called dependent variable (denoted as Y), whereas the variable used for prediction is independent variable (denoted as X). In a regression equation, there may be one or more independent variables, but there is only one dependent variable. Depending on whether there are one or more independent variables, the regression equation is called simple or multiple. The term „linear“ is added if the relationship between the dependent and the independent variable is linear. Thus a simple linear regression equation is represented as  $Y = a + bX$  Where, Y is dependent variable X is independent variable

„a“ is regression constant

„b“ is regression coefficient. It measures the change in Y corresponding to a change in X.

Similarly a multilinear regression equation is represented as

$$Y = a + b_1 X_1 + b_2 X_2 + \dots + b_n X_n$$

Where, Y is dependent variable

$X_1, X_2, \dots, X_n$ , are independent variables

„a“ is regression constant



„ $b_1, b_2 \dots b_n$ “ are respective regression coefficients.

Like the calculation of coefficient of the correlation, there are various methods of calculating regression equation: 1. From actual mean values of X and Y. 2. From assumed mean values of X and Y.

Calculation of regression equation using actual mean Regression equation (of Y on X) can be calculated using the following formula:  $Y - M_Y = b_{YX} (X - M_X)$  or  $Y - M_Y = r(\sigma_Y / \sigma_X)(X - M_X)$  As,  $b_{YX} = r(\sigma_Y / \sigma_X) = (\sum xy / \sum x^2)$ , the regression equation may be calculated using the following formula.  $Y - M_Y = (\sum xy / \sum x^2) (X - M_X)$  Where, Y and X are dependent and independent variables respectively;  $M_Y$  and  $M_X$  are means of Y and X variable respectively; and  $y = Y - M_Y$  and  $x = X - M_X$  The following example illustrates the calculation of the regression equation.

Example: Calculate the regression equation using the following data, taking age at marriage of husbands as independent variable and that of wives as dependent variable (see Figure 9.11)

Age at Marriage	Case1	Case2	Case3	Case4	Case5	Case6	Case7	Case8	Case9	Case10
Husbands	28	25	24	29	31	22	21	25	26	28
Wives	22	23	21	25	26	20	19	21	21	24

Calculation of regression equation using actual mean (see Figure 9.12)

Age of Y =  $Y - m_y$     Age of X =  $X - M_x$      $\sum xy$     Wives    Husbands

Figure 9.12 Calculation of Regression Equation using Actual Mean

## Notes

Y			X			
22	-0.2	00.04	28	2.1	4.41	-0.42
23	0.8	00.64	25	-0.9	0.81	-0.72
21	-1.2	01.44	24	-1.9	3.61	02.28
25	2.8	07.84	29	3.1	9.61	08.68
26	3.8	14.44	31	5.1	26.01	19.38
20	-2.2	04.84	22	-3.9	15.21	8.58
19	-3.2	10.24	21	-4.9	24.01	15.68
21	-1.2	01.44	25	-0.9	0.81	01.08
21	-1.2	01.44	26	0.1	0.01	-0.12
24	1.8	03.24	28	2.1	4.41	03.78
222	0	45.60	259	0	88.9	58.20

$$M_y = 222/10 = 22.2 \quad M_x = 259/10 = 25.9$$

$$y - M_y = (\sum xy / \sum x^2) * (X - M_x)$$

$$Y - 22.2 = (58.2 / 88.9) * (X - 25.2)$$

$$Y - 22.2 = 0.655 * (X - 25.2)$$

$$Y - 22.2 = 0.655X - 16.96$$

$$Y = 5.24 + 0.655X$$

Calculation of regression equation using assumed mean (see Figure 9.13)

Regression equation (of Y on X) can be calculated using the following formula, taking the assumed mean:  $Y - M_y = b_{yx} * (X - M_x)$  Where,  $b_{yx} = [\sum dX * dY - (\sum dX * \sum dY) / N] / [\sum dX^2 - (\sum dX)^2 / N]$

Y and X are dependent and independent variables respectively;  $M_y$  and  $M_x$  are mean of Y and X variables respectively  $dY = Y - A_{M_y}$  and  $dX = X - A_{M_x}$   $A_{M_y}$  and  $A_{M_x}$  are the assumed mean of Y and X variable respectively; and

Calculation of regression equation using assumed mean

Figure 9.13 Calculation of Regression Equation using Assumed Mean

Age of Wives	$d_y = Y - AMy$	$Y \cdot d_y^2$	Age of Husbands	$d_x = X - AM_{xx}$	$Dx^2$	$dx * dy$
22	0	0	28			0
23	1	1	25	0	0	0
21	-1	1	24	-1	1	1
25	3	9	29	4	16	12
26	4	16	31	6	36	24
20	-2	4	22	-3	9	6
19	-3	9	21	-4	16	12
21	-1	1	25	0	0	0
21	-1	1	26	1	1	-1
24	2	4	28	3	9	6
222	2	46	259	9	97	60

$$My = 222 / 10 = 22.2 \quad MX = 259 / 10 = 25.9$$

$$r = \frac{[\sum dx * dy - (\sum dx * \sum dy) / N]}{[\sum dx^2 - (\sum dx)^2 / N]}$$

$$byx = \frac{60 - (9 * 2) / 10}{97 - 9 * 9 / 10}$$

$$byx = 58.2 / 88.9 = 0.655$$

$$Y - My = byx * (X - MX)$$

$$Y - 22.2 = 0.655 * (X - 25.2)$$

$$Y - 22.2 = 0.655X - 16.96$$

$$Y = 5.24 + 0.655X$$

Standard error of estimate: Perfect prediction, using a regression equation is not possible (except when correlation value is -1 or + 1). Thus the researcher is interested in finding the accuracy of estimation of a regression equation. Standard error of estimate measures the error involved in using a regression equation as a basis of estimation. It can be calculated using the following equation:

$$SEE_{y..x} = \sqrt{\frac{\sum (Y - Y_c)^2}{N - 2}}$$

Where,  $SEE_{y..x}$  is Standard error of estimate

$Y_c$  is dependent variable

$Y_c$  is predicted value of  $Y$

$N$  is the number of observations

$$\text{It can also be calculated from the following formula } SEE_{y..x} = \sqrt{\frac{\sum Y^2 - a^2 Y - b \sum XY}{N - 2}}$$

Where, SEE Y ..Xis

Standard error of estimate Y is dependent variable X is independent variable „a“ is regression constant „b“ is regression coefficient. N is the number of observations Coefficient of determination: Coefficient of determination ( $r^2$ ) is the square of correlation coefficient (r) and is often used in interpreting the value of the coefficient of correlation. If the value of r were 0.8 then the coefficient of determination or  $r^2$  would be 0.64. This would mean that 64% of variance of one variable (dependent) is explained in terms of the other variable (independent).

Reflection and Action 9.3 I tried to understand how to make the calculation of regression equation using assumed mean. I could not succeed. May be you can explain it to me with an example. Write out on a separate sheet of paper your explanation with one or two examples.

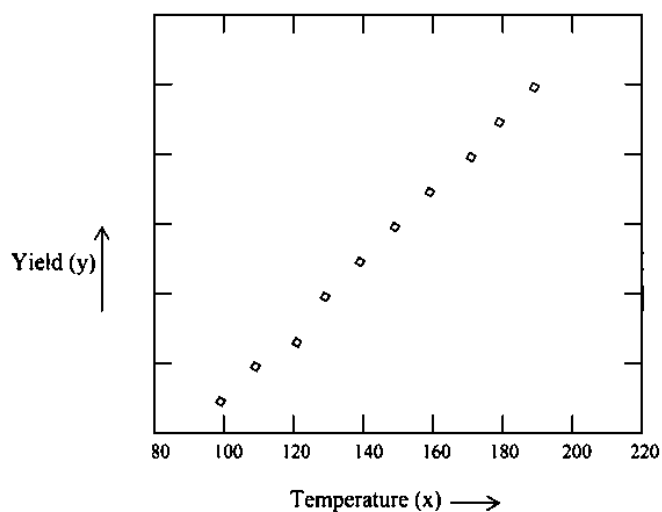
In many problems there are two or more variables that are inherently related and it may be necessary to explore the nature of their relationship. Regression analysis is a statistical technique for modeling and investigating the relationship between two or more variables. For example in a chemical process suppose that the yield of the product is related to the process operating temperature. Regression analysis can be used to build a model that expresses yield as a function of temperature. This model can be used to predict yield at a given temperature level. It can also be used for process optimization or process control purposes. In general, suppose that there is a single dependent variable or response variable  $y$  and that is related to  $k$  independent or regressor variables say  $x_1, \dots, x_k$ . The response variable  $y$  is a random variable and the regressor variables  $x_1, \dots, x_k$  are measured with negligible error. The relationship between  $y$  and  $x_1, \dots, x_k$  is characterized by a mathematical model and it is known as the regression model. It is also known as the regression of  $y$  on  $x_1, \dots, x_k$ . This regression model is fitted to a set of data. In many situations the experimenter knows the exact form of the functional relationship between  $y$  and  $x_1, \dots, x_k$ , say

$\varphi(x_1, \dots, x_k)$ , except for a set of unknown parameters. When the functional form is unknown, it has to be approximated on the basis of past experience or from the existing information. Because of its tractability, a polynomial function is popular in the literature. In this unit we will be mainly discussing the linear regression model and when  $k = 1$ , that is only one regressor variables. We will be discussing in details how to estimate the regression line and how it can be used for prediction purposes from a given set of data. We will also discuss briefly how we can estimate the function  $\varphi$ , if it is not linear.

We wish to determine the relationship between a single regressor variable  $x$  and a response variable  $y$  (note: The linear regression with one independent variable is referred to as simple linear regression). We will refer to  $y$  as the dependent variable or response and  $x$  as the independent variable or regressor. The regressor variable  $x$  is assumed to be a continuous variable controlled by the experimenter. You know that it is often easy to understand data through a graph. So, let us plot the data on Scatter diagram (a set of points in a 2-D graph where horizontal axis is regressor and vertical axis is response) Suppose that the true relationship between  $y$  and  $x$  is straight line. Therefore, each observation  $y$  can be described by the following mathematical relation (model)

Scatter diagram of yield versus temperature

$$y = \beta_0 + \beta_1 x + \epsilon \quad (1)$$



## Notes

where  $\epsilon$  is a random variable with mean 0 and variance  $\sigma^2$ . The  $\epsilon$  is known as the error component and it is assumed to be small. If the error  $\epsilon$  was absent then it was a perfect relation between the variables  $y$  and  $x$  which may not be very practical. Let us look at the following example. Example 1: A chemical engineer is investigating the effect of process operating temperature on product yield. The study results in the following data.

Temperature °C (x)	100	110	120	130	140	150	160	170	180	190
Yield, % (y)	45	51	54	61	66	70	74	78	85	89

The scatter diagram between the temperature and the yield is presented in the Figure 1 above. From the Figure 1 it is clear that there is a linear relationship between yield and temperature but clearly it is not perfect. For example we can not write the relationship between  $y$  and  $x$  as follows

$$y = \beta_0 + \beta_1 x$$

Clearly the presence of the error  $\epsilon$  is needed. Moreover the error  $\epsilon$  is a random variable because it is not fixed and it varies from one temperature to another. It may also vary when two observations are taken at the same temperature. If there was a perfect linear relationship between  $y$  and  $x$  we would have required just two points to find the relationship. Since the relationship is not perfectly linear it is usually required much more than two data points to find their relationship. Our main objective is to find the relationship between them from the existing information (data points). Since it is assumed that the relationship between  $x$  and  $y$  is linear therefore the relationship can be expressed by the equation (1) and finding the relationship basically boils down finding the unknown constants  $\beta_0$  and  $\beta_1$  from the observations. Let us discuss this concept of linear regression by one more illustration/collection of data described in the table 1 given below. This table encloses the data of 25 samples of cement, for each sample we have a pair of observation

(x,y) where x is percentage of SO<sub>3</sub>, a chemical and y is the setting time in minutes. These two components are strongly related; it is the percentage of SO<sub>3</sub> which influences the setting time of any cement sample, the recorded observations are given in table 1 below.

**Table 1: Data on SO<sub>3</sub> and Setting Time**

S.No. i	Percentage of SO <sub>3</sub> x	Setting Time Y (in minutes)
1	1.84	190
2	1.91	192
3	1.90	210
4	1.66	194
5	1.48	170
6	1.26	160
7	1.21	143
8	1.32	164
9	2.11	200
10	0.94	136
11	2.25	206
12	0.96	138
13	1.71	185
14	2.35	210
15	1.64	178
16	1.19	170
17	1.56	160
18	1.53	160
19	0.96	140
20	1.7	168
21	1.68	152
22	1.28	160
23	1.35	116
24	1.49	145
25	1.78	170
<b>Total</b>	<b>39.04</b>	<b>4217</b>
<b>Sum of Squares</b>	<b>64.446</b>	<b>726539</b>

From the table 1, you see that setting time y increases as percentage of SO<sub>3</sub> increases. Whenever you find this type of increasing (or decreasing) trend in a table, same will be reflected in the scatter diagram, and it indicates that there is a linear relationship between x and y. By drawing the scatter diagram you can observe that the relationship is not perfect in the sense that a straight line cannot be drawn through all the points in the scatter diagram. Nevertheless, we may approximate it with some linear

## Notes

equation. What formula shall we use? Suppose, we use the formula  $y = 90 + 50x$  to predict  $y$  based on  $x$ . To examine how good this formula is, we need to compare the actual values of  $y$  with the corresponding predicted values. When  $x = 0.96$ , the predicted  $y$  is equal to  $138 (= 90 + 50 \times 0.96)$ . Let  $(x_i, y_i)$  denote the values of  $(x, y)$  for the  $i$ th sample. From Table-1, notice that  $x_{12} = x_{19} = 0.96$ , whereas  $y_{12} = 138$  and  $y_{19} = 140$ . Let  $\hat{y}_i$  be the predicted value of  $y$  (then using  $y = 90 + 50x$  for the  $i$ th sample). Since,  $x_{12} = x_{19} = 0.96$ , both  $\hat{y}_{12}$  and  $\hat{y}_{19}$  are equal to 138. Thus the difference, the error in prediction, also called residual is observed to be  $e_{12} = 0$  and  $e_{19} = 2$ . The formula we have considered above,  $y = 90 + 50x$ , is called a simple linear regression equation, we will study these terms in detail in our successive sections.

### Least squares estimation

Suppose that we have  $n$  pairs of observations, say  $(x_1, y_1), \dots, (x_n, y_n)$ . It is assumed that the observed  $y_i$  and  $x_i$  satisfy a linear relation as given in the model (1). These data can be used to estimate the unknown parameters  $\beta_0$  and  $\beta_1$ . The method we are going to use is known as the method of least squares, that is, we will estimate  $\beta_0$  and  $\beta_1$  so that the sum of squares of the deviations from the observations to the regression line is minimum. We will try to explain it first using a graphical method in Figure 2. For illustrative purposes we are just taking 5 data points  $(x, y) = (0.5, 57), (0.75, 64), (1.00, 59), (1.25, 68), (1.50, 74)$ . The estimated regression line can be obtained as follows. For any line we have calculated the sum of the differences (vertical distances) squares between the  $y$  value and the value, which is obtained using that particular line. Now the estimated regression line is that line for which the sum of these differences squares is minimum.

### NON-LINEAR REGRESSION

Linear regression is a widely used method for analyzing data described by models which are linear in parameters. However, in many practical



situations, people come across with data where the relationship between the independent variable and the dependent variable is no more linear. In that case definitely one should not try to use a linear regression model to represent the relationship between the independent and dependent variable.

**Check Your Progress 2**

Note: a) Use the space provided for your answer.

b) Check your answers with those provided at the end of the unit.

1. Define Correlation.

.....  
.....  
.....  
.....  
.....

2. Discuss the Method of Calculating Correlation of Ungrouped Data.

.....  
.....  
.....  
.....  
.....

3. Discuss the Method of Calculating Correlation of Grouped Data.

.....  
.....  
.....  
.....  
.....

4. Define Regression.

.....  
.....  
.....  
.....  
.....

---

## 9.6 LET US SUM UP

---

Unit 9 is the last unit on Quantitative Methods. All five units of this block have emphasized that quantitative methods should be used in social research when they are necessary and relevant and can provide superior results. Sometimes you can use them in combination with the qualitative methods. You need not avoid the quantitative methods because of lack of information or apprehension that it is difficult to understand them.

In this unit you have seen :

- that regression analysis is an important technique, which can be used to verify the results of any experiment.
- How to determine the relationship between a dependent and an independent variable by using the Scatter diagram
- that by knowing the technique of regression you have an edge to analyse the results in an organized way. Further this analysis is smoothened by application of the concepts like least square estimation, goodness to fit and residual analysis.
- that many times the data obtained by conducting an experiment does not follow the linear relation. So, to handle such aspects we have also discussed the concept of non linear regression, under we have emphasized least square estimation technique.
- Formulas and applications of following topics: *f* Simple Linear Regression

- Least Squares Estimation
- Goodness to Fit
- Residual Analysis *f* Non-Linear Regression
- Least Squares Estimation.

In statistical modeling, regression analysis is a set of statistical processes for estimating the relationships between a dependent variable (often called the 'outcome variable') and one or more independent variables (often called 'predictors', 'covariates', or 'features'). The most common form of regression analysis is linear regression, in which a researcher finds the line (or a more complex linear function) that most closely fits the data according to a specific mathematical criterion. For example, the method of ordinary least squares computes the unique line (or hyperplane) that minimizes the sum of squared distances between the true data and that line (or hyperplane). For specific mathematical reasons (see linear regression), this allows the researcher to estimate the conditional expectation (or population average value) of the dependent variable when the independent variables take on a given set of values. Less common forms of regression use slightly different procedures to estimate alternative location parameters (e.g., quantile regression or Necessary Condition Analysis) or estimate the conditional expectation across a broader collection of non-linear models (e.g., nonparametric regression).

Regression analysis is primarily used for two conceptually distinct purposes. First, regression analysis is widely used for prediction and forecasting, where its use has substantial overlap with the field of machine learning. Second, in some situations regression analysis can be used to infer causal relationships between the independent and dependent variables. Importantly, regressions by themselves only reveal relationships between a dependent variable and a collection of

independent variables in a fixed dataset. To use regressions for prediction or to infer causal relationships, respectively, a researcher must carefully justify why existing relationships have predictive power for a new context or why a relationship between two variables has a causal interpretation. The latter is especially important when a researcher hopes to estimate causal relationships using observational data.

---

### 9.7 KEY WORDS

---

**Regression:** In statistical modeling, regression analysis is a set of statistical processes for estimating the relationships between a dependent variable and one or more independent variables.

**Correlation:** In statistics, dependence or association is any statistical relationship, whether causal or not, between two random variables or bivariate data. In the broadest sense correlation is any statistical association, though it commonly refers to the degree to which a pair of variables is linearly related.

---

### 9.8 QUESTIONS FOR REVIEW

---

5. Define Correlation.
6. Discuss the Method of Calculating Correlation of Ungrouped Data.
7. Discuss the Method of Calculating Correlation of Grouped Data.
8. Define Regression.

---

### 9.9 SUGGESTED READINGS AND REFERENCES

---

- William H. Kruskal and Judith M. Tanur, ed. (1978), "Linear Hypotheses," International Encyclopedia of Statistics. Free Press, v. 1,
- Evan J. Williams, "I. Regression," pp. 523–41.
- Julian C. Stanley, "II. Analysis of Variance," pp. 541–554.

- Lindley, D.V. (1987). "Regression and correlation analysis," New Palgrave: A Dictionary of Economics, v. 4, pp. 120–23.
- Birkes, David and Dodge, Y., *Alternative Methods of Regression*. ISBN 0-471-56881-3
- Chatfield, C. (1993) "Calculating Interval Forecasts," *Journal of Business and Economic Statistics*, 11. pp. 121–135.
- Draper, N.R.; Smith, H. (1998). *Applied Regression Analysis* (3rd ed.). John Wiley. ISBN 978-0-471-17082-2.
- Fox, J. (1997). *Applied Regression Analysis, Linear Models and Related Methods*. Sage
- Hardle, W., *Applied Nonparametric Regression* (1990), ISBN 0-521-42950-1
- Meade, Nigel; Islam, Towhidul (1995). "Prediction intervals for growth curve forecasts". *Journal of Forecasting*. 14 (5): 413–430. doi:10.1002/for.3980140502.
- Sen, M. Srivastava, *Regression Analysis — Theory, Methods, and Applications*, Springer-Verlag, Berlin, 2011 (4th printing).
- T. Strutz: *Data Fitting and Uncertainty (A practical introduction to weighted least squares and beyond)*. Vieweg+Teubner, ISBN 978-3-8348-1022-9.
- Malakooti, B. (2013). *Operations and Production Systems with Multiple Objectives*. John Wiley & Sons.

---

## 9.10 ANSWERS TO CHECK YOUR PROGRESS

---

### Check Your Progress 1

1. See Section 9.2
2. See Section 9.3

### Check Your Progress 2

1. See Section 9.4
2. See Section 9.5

---

# UNIT 10: HYPOTHESIS TESTING, T-TEST, Z-TEST, CHISQUARE)

---

## STRUCTURE

- 10.0 Objectives
- 10.1 Introduction
- 10.2 Classification of Statistical Tests
- 10.3 Parametric Tests
  - 10.3.1 Sampling Distribution of Means
    - A Large Samples
    - B Confidence Intervals and Levels of Significance
    - C Small Samples
    - D Degree of Freedom
  - 10.3.2 Application of Parametric Tests
    - A Application of Z-test
    - B Two-tailed and one-tailed tests
    - C Application of t-test
    - D Application of F-test
    - E Factor Analysis
- 10.4 Non-parametric Tests and Application of Chi-square Test
  - A Application of Chi-square test
  - B Application of Median test
- 10.5 Let us sum up
- 10.6 Key Words
- 10.7 Questions for Review
- 10.8 Suggested readings and references
- 10.9 Answers to Check Your Progress

---

## 10.0 OBJECTIVES

---

This Unit aims to provide you with detailed information about the nature and use of parametric and non-parametric tests in general and the application of some of these tests for drawing inferences and generalizations. On the completion of this

Unit, you should be able to:

- Classify various statistical tests,
- Describe the nature of parametric tests alongwith the assumptions on which they are based,
- Work out sampling distribution of means in the context of (i) large samples, and (ii) small samples,
- Define and illustrate the concept of confidence intervals and levels of significance
- Define and illustrate the concept of degrees of freedom,
- Use Z-test and t-test in testing the significance of the difference between means,
- Define and illustrate the concept of one-tailed and two-tailed tests of
- significance,
- Describe the nature and uses of analysis of variance,
- Describe the nature of the non-parametric tests alongwith their assumptions,
- Use of chi-square test, and
- Describe the use of median test and its application.

---

## 10.1 INTRODUCTION

---

In Unit 10, we focussed on descriptive statistics including the various measures of central tendency, variability, relative positions and relationships. These measures are used to describe the properties of particular samples. In this Unit, we shall introduce inferential or sampling statistics. The knowledge of these statistics is useful for testing the hypothesis(es) related to your research problems, and to make generalisations about populations on the basis of data analysis. This requires you to be familiar with certain statistical tests - parametric and nonparametric.

---

## 10.2 CLASSIFICATION OF STATISTICAL TESTS

---

The descriptive statistics already discussed in Unit 1 are used to explain the properties of samples drawn from a population. The researcher computes certain 'statistics' (sample values) as the basis for inferring the corresponding 'parameters' (population values). Ordinarily, a single sample is drawn from a given population so as to determine how well a researcher can infer or estimate the 'parameter' from a computed sample 'statistics'. For making the inferences about the various parameters, the researcher makes use of parametric and non-parametric tests.

---

### **10.3 PARAMETRIC TESTS**

---

Under this section we discuss two sub-themes: sampling distribution of means and application of parametric tests. Sampling distribution of means covers a) large samples, b) confidence intervals and levels of significance, c) small samples, and d) degree of freedom. Application of parametric tests covers three tests, namely Z-test, t-test and F-test.

Parametric tests are the most powerful statistical tests for testing the significance of the computed sampling statistics. These tests are based on the following assumptions:

1. the variables described are expressed in interval or ratio scales and not in nominal or ordinal scales of measurement,
2. the population values are normally distributed,
3. the samples have equal or nearly equal variances-this condition is known as 'equality or homogeneity of variances' and is particularly important to determine for small samples,
4. the selection of one case in the sample is not dependent upon the selection of any other. Z-test, t-test and F-test are the most commonly used parametric tests. Before discussing the application of these tests, it is necessary to describe certain concepts relating to 'sampling distribution of means', 'confidence intervals', 'levels of confidence of significance', and 'degrees of freedom'.



### 10.3.1 Sampling Distribution of Means

#### A Large Samples

An important principle, known as the 'central limit theorem', describes the characteristics of sample means. If a large number of equal-sized samples (greater than 30) are selected at random from an infinite population,

the distribution of 'sample means' is normal and it possesses all the characteristics of a normal distribution, the average value of 'sample means' will be the same as the mean of the population, the distribution of the sample means around the population mean will have its own standard deviation, known as 'standard error of mean, which is denoted as SEM or OM. It is computed by the formula:

$$SE_M = \sigma_M = \frac{\bar{\sigma}}{\sqrt{N}} \dots\dots\dots(1)$$

in which  $\bar{\sigma}$  = Standard deviation of the population and  
 $N$  = The number of cases in the sample.

Since the value of  $\bar{\sigma}$  (i.e. standard deviation of population) is usually not known, we make an estimate of this standard error of mean by the formula:

$$\sigma_M = \frac{\sigma}{\sqrt{N}} \dots\dots\dots(2)$$

in which  $\sigma$  = Standard deviation of the sample  
 $N$  = The number of cases in the sample.

To illustrate the use of formula (2), we assume that the mean of the attitude scores of a sample of 100 distance learners enrolled with IGNOU towards student support services is 25 and the standard deviation is 5. The standard error of mean can be calculated accordingly:

$$SE_M = \sigma_M = \frac{5.0}{\sqrt{100}} = 0.50$$

This standard 'error of mean' may be assumed as the standard deviation of a distribution of sample means, around the fixed population mean of all distance learners. In the case of large randomly selected samples, the sampling distribution of sample means is assumed to be normal.

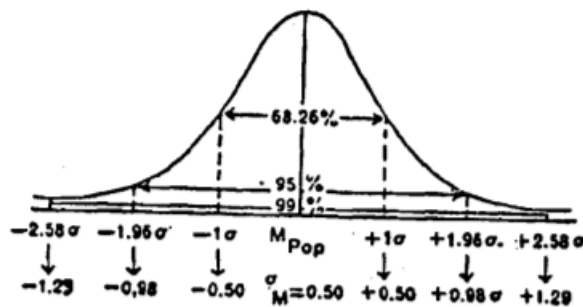


Fig. 10.1: Sampling Distribution of Means showing Variability of obtained Means around the Population Mean in terms of  $\sigma_M$

The normal curve in Figure 1 shows that this sampling distribution is centered around the unknown population mean with standard deviation 0.50. The sample means often fall between the positive and the negative side of the population mean. About 213 of our sample means (exactly 68.26 per cent) will lie within  $\pm 1.00$  cr. of the population mean, i.e., within a range of  $\pm 1 \times 0.50 = \pm 0.50$ . Furthermore, 95 of our 100

sample means will lie within  $\pm 2.00$  a, (more exactly  $\pm 1.96$  a, ) of the population mean, i.e. 95 of 100 sample means will lie within  $\pm 1.96 \times 0.50$  or  $\pm 0.98$  of the population mean. In other words, the probability that our sample mean of 25 does not miss the population mean ( $M_{pop.}$ ) by more than  $\pm 0.98$  is 0.95. Also, 99 of our sample means will be within  $\pm 3.00$  a, (more exactly  $\pm 2.58$  a, ) of the population mean. This indicates that 99 out of 100 sample means will fall within  $\pm 2.58 \times 0.50$  or  $\pm 1.29$  of the population mean. The probability (P) that our sample mean of 25 does not miss the  $M_{pop.}$  by more than  $\pm 1.29$  is .99. Thus, the value of a population mean, to be inferred from a randomly selected sample mean, can be estimated on a probability basis.

Thus, the value of a population mean, to be inferred from a randomly selected sample mean, can be estimated on a probability basis. In case of proportion (7c) the SE(p) can be estimated using the formula:

$$SE_{(p)} = \sqrt{\frac{p(1-p)}{n}}$$

**Check Your Progress 1**

- 1) Describe the assumptions on which the use of parametric tests are based.

.....  
 .....  
 .....

- 2) Given a sample of 100 children 1 - 3 year of age with mean (SD) intake of calcium = 1.75 (5.82). Compute the standard error of mean.

.....  
 .....  
 .....

**B Confidence Intervals and Levels of Significance**

## Notes

When we draw a large random sample from the population to obtain measures of a variable and compute the mean for the sample, we can use the 'central limit theorem' and 'normal probability curve' to have an estimate of the population mean. We can say that  $\bar{M}$  has a 95 percent chance of being within 1.96 standard error units of  $M_{pop}$ . In other words, a mean for a random sample has a chance of 95 percent of being within 1.96  $\sigma_{\bar{M}}$  units from  $M_{pop}$ . It may also be said that there is a 99 percent chance that the sample mean lies within 2.58  $\sigma_{\bar{M}}$  units of  $M_{pop}$ . To be more specific, it may be stated that there is a 95 percent probability that the limits  $\bar{M} \pm 1.96 \sigma_{\bar{M}}$  enclose the population mean, and the limits  $\bar{M} \pm 2.58 \sigma_{\bar{M}}$  enclose the population mean with 99 percent probability. Such limits enclosing the population mean are known as the 'confidence intervals'. These limits help us to adopt particularly two levels of confidence. One is known as 5 percent level or 0.05 level, and the other is known as 1 per cent level or .01 level. The .05 level of confidence indicates that the probability  $M_{pop}$  that lies within the interval  $\bar{M} \pm 1.96 \sigma_{\bar{M}}$  is 0.95 and that it falls outside of these limits is .05. By saying that probability is 0.99, it is meant that  $M_{pop}$  lies within the interval  $\bar{M} \pm 2.58 \sigma_{\bar{M}}$ , and that the probability of its falling outside of these limits is .01.

To illustrate, let us apply the concept to the previous problem. Taking as our limits  $\bar{M} \pm 1.96 \sigma_{\bar{M}}$ , we have  $25 \pm 1.96(0.50)$  or a confidence interval marked off by the limits 24.02 and 25.98. Our confidence that this interval contains  $M_{pop}$  is expressed by a probability of .95. If we want a higher degree of confidence, we can take the .99 level of confidence for which the limits are  $\bar{M} \pm 2.58 \sigma_{\bar{M}}$ , or a confidence interval given by the limits 23.71 and 26.29. We may be quite confident that  $M_{pop}$  is not lower than 23.71 nor higher than 26.29, i.e., the chances are 99 in 100 that the  $M_{pop}$  lies between 23.71 and 26.19.

C Small Samples

When the number of cases in the sample is less than 30, we may estimate the value of  $\sigma$ , by the formula:

$$SE_M = \frac{S}{\sqrt{N}} \dots\dots\dots(3)$$

In which

S = Standard deviation of the small sample.  
 N = The number of cases in the sample.

The formula for computing S is

$$S = \sqrt{\frac{\sum x^2}{N-1}} \dots\dots\dots(4)$$

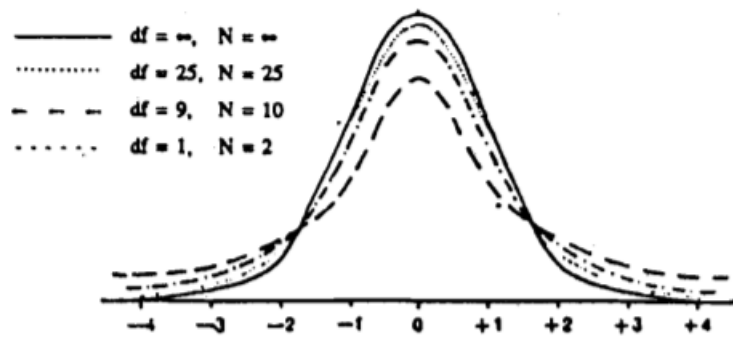
In which

$\sum x^2$  = Sum of the squares of deviations of individual scores from the sample mean.

N = The number of cases in the sample.

The concept of small size was developed by William Seely Gosset, a consulting statistician for Guinness Breweries of Dublin (Ireland) around 1915. The principle is that we should not assume that the sampling distribution of means of small samples is normally distributed. He found that the distribution curves of small sample means were somewhat different from the normal curve. When the size of the sample is small, the t-distribution lies under the normal curve, but the tails or ends of the curve are higher than the corresponding parts of the normal curve. Figure 2 shows that the t-distribution does not differ significantly from the normal distribution unless the sample size is quite small. Further, as the sample size increases in size, the t-distribution approaches more and more closely to the normal curve.

Fig. 2: Distribution of t-values for Degrees of Freedom from 1 to 2 (when df is very large, the distribution of t approaches the normal).



For small samples, it is necessary to make use of selected points in the table of Gosset's t-critical values or student's t-values, given in Appendix (Table 11). As the sample size increases, the student's t-values approach the z-values of Normal probability Table. When small samples are used, the use of t-values involves an important concept known as 'degrees of freedom', which we shall now discuss separately.

#### D Degree of Freedom

While finding the standard deviation of small samples we use  $N-1$  in the denominator instead of  $N$  in the basic formula for standard deviation. The difference in the two formulae may seem very little, if  $N$  is sufficiently large. But there is a very important difference in the 'meaning' in the case of small samples.  $N-1$  is known as the 'number of degrees of freedom', denoted by  $df$ . The 'number of degrees of freedom' in a distribution is the number of observations or values that are independent of each other and cannot be deduced from each other. In other words, we may say that the 'degrees of freedom' connote freedom to vary. To illustrate as to why the  $df$  used here is  $N-1$ , we take 5 scores, i.e., 5, 6, 7, 8 and 9, the mean of which is 7. This mean score is to be used as an estimate of the population mean. The deviations of the scores from the mean 7 are - 2, - 1, 0, +1 and +2. A mathematical requirement of the mean is that the sum of these deviations should be zero. Of the five deviations, only 4, i.e.,  $N-1$  can be chosen freely (independently) as the condition that the sum is equal to zero restricts the value of the 5th deviate. With this condition, we can arbitrarily change any four of the

five deviates and thereby fix the fifth. We could take the first four deviates as -2, -1, XI and +1, which would mean that for the sum of deviates to be zero, the fifth deviate has to be +2. Similarly, we can try other changes and if the sum is to remain zero, one of the five deviates is automatically determined. Hence, only 4, i.e., (5-1)'s are free to vary within the restrictions imposed. When a statistic is to be used to estimate a parameter, the number of degrees of freedom depends upon the restrictions imposed. One df is lost for each of the restrictions imposed. Therefore, the number of df varies from one statistics to another. For example, in estimating and computing the population mean ( $M_{pop}$ ) from the sample Mean ( $M$ ), we lose 1 df. So, the number of degrees of freedom is (N-1). Let us determine the .95 and .99 confidence intervals for the population mean ( $M_{pop}$ ) of the scores 10, 15, 10, 25, 30, 20, 25, 30, 20 and 15, obtained by 10 distance learners on an attitude scale. The mean of the scores is

$$= \frac{10 + 15 + 10 + 25 + 30 + 20 + 25 + 30 + 20 + 15}{10}$$

$$= \frac{200}{10} = 20.00$$

Using formula (4) we compute the standard deviation as:

## Notes

X	$x = X - M$	$x^2$
10	-10	100
15	-5	25
10	-10	100
25	5	25
30	10	100
20	0	0
25	5	25
30	10	100
20	0	0
15	-5	25
$\sum x = 0$		$\sum x^2 = 500$

$$\begin{aligned}
 s &= \sqrt{\frac{\sum x^2}{N-1}} \\
 &= \sqrt{\frac{500}{10-1}} \\
 &= 7.45
 \end{aligned}$$

From formula (3) we compute

$$\begin{aligned}
 SE_M &= \frac{7.45}{\sqrt{10}} \\
 &= 2.36
 \end{aligned}$$

For estimating the M pop from the sample mean of 20.00, we determine the value oft at the selected points using appropriate number of degrees of freedom. The available df for determining t is N-1 or 9. Entering Table II (See Appendix) with 9 df, we read that t = 2.26 at .05 level and 3.25 at .01 level. From the first t-value we know that 95 of our 100 sample means will lie within f 2.26 SE, or f 2.26 x 2.36 of the population mean and 5 out of 100 fall outside of these limits. The probability (P) that our sample mean 20.00 does not miss the M pop by more than f 2.26 x 2.36 or f 5.33 is .95 From the second t-value, we know that 99 percent of our sample mean will lie between M pop and 3.25 SE, or f 3.25 x 2.36, and that 1 percent fall will beyond these limits. So, the probability (P) that our sample mean of 20.00 does not miss the M pop by more than f 3.25 x 2.36 or f 7.67 is .99. Taking our limits as M f 2.26 SE, ,we have 20.00 f 2.26 x 2.36 or 14.67 and 25.33 as the limits of the .95 confidence



interval. The probability (P) that M pop is not less than 14.67 nor greater than 25.33 is .95. Taking the limits  $M \pm 3.25 \text{ SE}$ , we have  $20.00 \pm 3.25 \times 2.36$ , or 12.33 and 27.67 as the limits of the .99 confidence interval, and the probability (P) so that M pop is not less than 12.33 and not greater than 27.67 is .99. The use of small samples to build generalizations in educational research should be made cautiously as it is difficult to ensure that a small sample adequately represents the population from which the sample is drawn. Further more, conclusions drawn from small samples are usually unsatisfactory because of the great variability from sample to sample. In other words, large samples drawn randomly from the population will provide a more accurate basis than will small samples for inferring population parameters.

### 10.3.2 Application of Parametric Tests

In this subsection, we shall discuss the application of three parametric tests, namely Z-test, t-test and F-test.

#### A Application of Z-test

It has already been explained in earlier sections that the frequency distribution of large sample means drawn from the same population fall into a normal distribution around M pop as their measure of central tendency. It is also reasonable to expect that the frequency distribution of the difference between the means computed from the two samples will also tend to be normal with a mean of zero and standard deviation of 1. It is termed the 'standard error of the difference between two means' and is denoted by  $\sigma_{d_M}$ . It is computed by the formula:

$$\sigma_{d_M} = \sqrt{\sigma_{M_1}^2 + \sigma_{M_2}^2} \dots\dots\dots(5)$$

In which

## Notes

$\sigma_{M_1}$  = the SE of the mean of the first sample

$\sigma_{M_2}$  = the SE of the mean of the second sample

To illustrate, we apply formula (5) to a problem. Suppose a test of creativity was administered on two groups, one of 120 males and the other of 75 females enrolled in Post Graduate Diploma in Distance Education (PGDDE) course of IGNOU. The results are summarized in table 10.1 below.

Table 10.1 Means and Standard Deviations of Two Independent Large Samples

Statistics	Boys	Girls
N	120	75
Mean (M)	57.50	55.75
Standard Deviation		
$\sigma$	8.42	8.13

Assuming that our samples are random, it is to be ascertained whether the I difference between the means 57.50 and 55.75 is significant. 1 I Using formula (5) we compute the 'standard error of the difference between means'

$$\begin{aligned}\sigma_{d_M} &= \sqrt{\frac{(8.42)^2}{120} + \frac{(8.13)^2}{75}} \\ &= 1.21\end{aligned}$$

The obtained difference between the means of males and females is 1.75 (i.e., 57.50 - 55.75); and the SE of this difference ( $\sigma_{d_M}$ ) = 1.21. To determine whether two groups of males and females actually differ in creative thinking ability, we set-up a null hypothesis, i.e., the difference between the population means of males and females is zero and that, except for sampling errors, means of males and females is zero and that,

except for sampling errors, mean differences from sample to sample will also be zero. In accordance with a null hypothesis, we assume a sampling distribution of differences to be normal with the mean at zero, (or at  $M_{\text{pop (males)}} - M_{\text{pop (females)}} = 0$ ). The deviation of each sample difference,  $[M_{\text{males}} - M_{\text{females}}] - [M_{\text{pop (males)}} - M_{\text{POP (females)}}]$  or  $[M_{\text{males}} - M_{\text{females}}] - 0$ . The deviation of each sampled difference between two means, given in terms of standard measure, would be the deviation divided by the standard error, which gives us a Z - value in terms of a general formula:

$$\bar{z} = \frac{|M_1 - M_2|}{d_M} \dots\dots\dots(6)$$

Using formula (6)

$$\bar{z} = \frac{1.75}{1.21} = 1.45$$

For the sake of convenience, we use .05 and .01 levels of significance as two arbitrary standards for accepting or rejecting a hypothesis. From the normal distribution Table 1 (See Appendix) we read that + 1.960 mark off points along the base line of the normal curve to the left and right of which lie 5 percent of the cases (i.e., 2.5 percent at each end of the curve). When a Z - value is 1.96 or more, we reject a null hypothesis at .05 level of significance. The computed Z -value of 1.45 in our problem falls short of 1.96, i.e., it does not reach the .05 level. Accordingly, we retain the null hypothesis and conclude that two groups of males and females actually do not differ in their mean performance on creative-thinking-test. Furthermore, from Table I (See Appendix) we know that + 2.580 mark off points to the left and right of which lie 1 percent (0.5 percent at each end of the curve) of the cases in the normal distribution. If the Z -value is 2.58 or more, we reject the null hypothesis at .01 level and the probability (P) is that not more than once in 100 trials would a difference of this size arise if the true difference ( $M_{\text{pop}_1} - M_{\text{pop}_2}$ ) was zero.

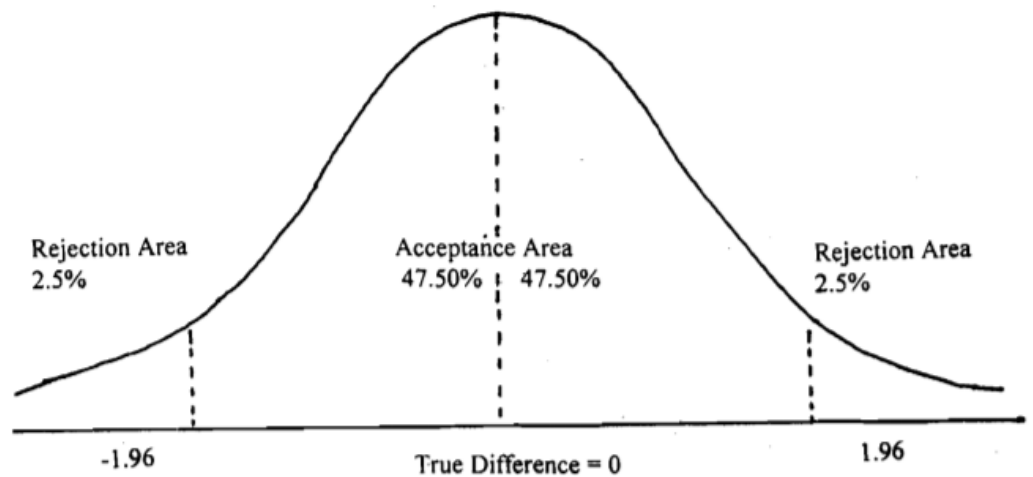
B Two-tailed and one-tailed tests

## Notes

Suppose a null-hypothesis were set up that there was no difference, other than a sampling error difference, between the mean height of two groups, A and B. We would be concerned only with the difference and not with the superiority or inferiority in height of either group. To test this hypothesis, we apply two-tailed test as the difference between the obtained means of height of two groups may be as often in one direction (plus) as in the other (minus) from the true difference of zero. Moreover, for determining probability, we take both tails of sampling distribution.

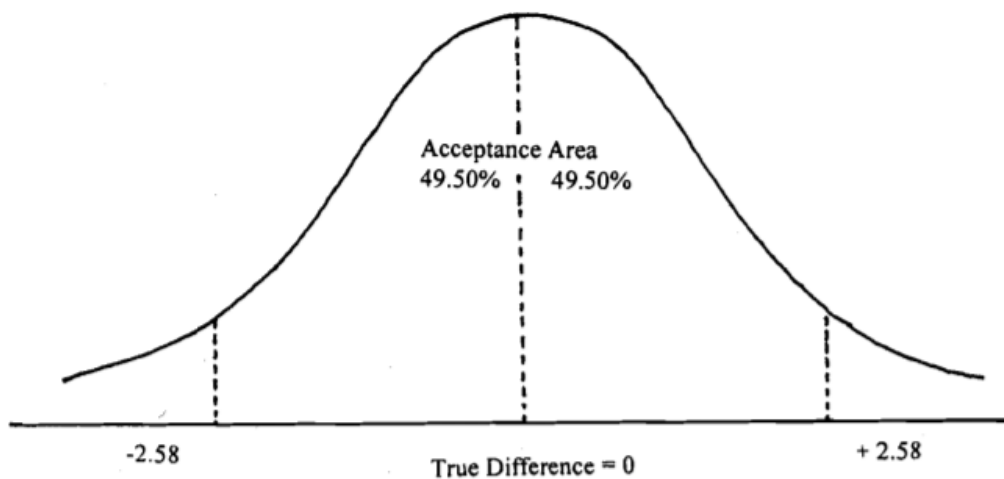
For a large sample two-tailed test, we make use of a normal distribution curve. The 5 percent area of rejection is divided equally between the upper and the lower tails of this curve and we have to go out to  $\pm 1.96$  on the base line of the curve to reach the area of rejection as shown in the Figure 10.3.

Fig. 10.3 A Two-tailed Test at .05 level (2.5 per cent at each level)



Similarly, if we have 0.5 percent area at each end of the normal curve where 1 percent area of rejection is to be divided equally between its upper and lower tails, it is necessary to go out to  $\pm 2.58$  on the base line to reach the area of rejection as shown in the Fig.10.4.

Fig. 10.4. A Two-tailed Test at .01 level (0.5 per cent at each level)



In the case of the above example a null hypothesis was set up that there was no difference other than a sampling error difference between the mean creative thinking score of males and females of PGDDE. Thus, we were concerned with a difference, and not in superiority or inferiority of either group in the creative thinking ability. To test this hypothesis, we applied 'two-tailed test' as the difference between the two means might have been in one direction (plus) or in the other (minus) from the true difference of zero; and we took both tails of sampling distribution in determining probabilities. As is evident from the above example, we make use of a normal distribution curve in the case of a large sample 'two-tailed test'. The 5 percent area of rejection is equally divided between the upper and lower tails of the curve and we have to go out to + 1.96 on the base line of the curve to reach the area of rejection. Similarly, we have 0.5 percent area at each end of the normal curve when 1 percent of rejection is to be divided equally between its upper and lower tails and it is necessary to go out to + 2.58 on the base line to reach the area of rejection. In the above problem, if we change the null hypothesis as: male group of PGDDE have significantly higher creative thinking than that of the female group; or male group have significantly lower creative thinking than the female group of PGDDE course, then each of these hypotheses indicates a direction of difference. In such situations, the use of 'one-tailed test' is made. For such a test, the 5 percent area or 1 per

**Notes**

cent area of rejection is either at the upper tail or at the lower tail of the curve, to be read from 0.10 column (instead of .05) and .02 column (instead of .01).

**C Application of t-test**

We have already discussed that the frequency distribution of small sample means drawn from the same population forms a t-distribution and it is reasonable to expect that the sampling distribution of the difference between the means computed from two different populations will also fall under the category of t-distribution. Fisher provided the formula for testing the difference between the means computed from independent small samples.

$$t = \frac{|M_1 - M_2|}{\sqrt{\left(\frac{\sum x_1^2 + \sum x_2^2}{N_1 + N_2 - 2}\right)\left(\frac{N_1 + N_2}{N_1 \times N_2}\right)}} \dots\dots\dots(7)$$

in which

- $M_1$  and  $M_2$  = means of two samples
- $\sum x_1^2$  and  $\sum x_2^2$  = sums of squares of the deviations from the means in the two samples
- $N_1$  and  $N_2$  = number of cases in the two samples.
- df = degrees of freedom =  $N_1 + N_2 - 2$

To illustrate the use of the formula, let us test the significance of the difference between mean scores of 7 boys and 10 girls on an intelligence test.

Table 10.2: Scores of 7 boys and 10 girls on an intelligence test

$(N_1 = 7)$			$(N_2 = 10)$		
Boys $X_1$	$x_1$	$x_1^2$	Girls $X_2$	$x_2$	$x_2^2$
13	0	0	10	-4	16
14	1	1	16	2	4
11	-2	4	12	-2	4
12	-1	1	13	-1	1
15	2	4	18	4	16
13	0	0	13	-1	1
13	0	0	19	5	25
			14	0	0
			13	-1	1
			12	-2	4
<hr/> $\Sigma X_1 = 91$ <hr/>			<hr/> $\Sigma X_2 = 140$ <hr/>		
<hr/> $\Sigma x_1^2 = 10$ <hr/>			<hr/> $\Sigma x_2^2 = 72$ <hr/>		
$M_1 = \frac{91}{7} = 13$			$M_2 = \frac{140}{10} = 14$		
<hr/> $df = N_1 + N_2 - 2 = 7 + 10 - 2 = 15$ <hr/>					

Using formula (7)

$$\begin{aligned}
 t &= \frac{|14 - 13|}{\sqrt{\left(\frac{10 + 72}{7 + 10 - 2}\right)\left(\frac{7 + 10}{7 \times 10}\right)}} \\
 &= \frac{1}{\sqrt{\left(\frac{82}{15}\right)\left(\frac{17}{70}\right)}} \\
 &= \frac{1}{\sqrt{9.29}} \\
 &= 0.33
 \end{aligned}$$

To test the significance of difference between the two means by making use of two-tailed test (null hypothesis, i.e., no differences between the two groups), we look for the t-critical values for rejection of null hypothesis in Table II (Appendix) for  $(7 + 10 - 2)$  or 15 df. These t-values are 2.13 at .05 and 2.95 at .01 levels of significance. Since the obtained t-value 0.33 is less than the table value necessary for the rejection of the null hypothesis at .05 level for df 15, the null hypothesis

## Notes

is accepted and it may be concluded that there is no significant difference in the mean intelligence scores of males and females.

If we change the null hypothesis as: boys will have higher intelligence scores than girls, or males will have lower intelligence scores than females, then each of these hypotheses indicates a direction of difference rather than simply the existence of the difference. So, we make use of one-tailed test. For given degrees of freedom, i.e., 15, the .05 level is read from the 0.10 column ( $\alpha = .05$ ) and the .01 level from 0.02 column ( $\alpha = .01$ ) of the t-table. In the one-tailed test, for 15 df t-critical values at .05 and .01 levels, as read from the 0.10 and the 0.02 columns of Table II are 1.75 and 2.60 respectively. Since the computed t-value of 0.33 does not reach the table value at .05 level (i.e., 1.75 for .10), we may conclude that the difference in two groups is there merely because of chance factors.

### D Application of F-test

The use of 2 and t-test is made by a researcher to determine whether there is any significant difference between the mean of two random samples. Suppose we have seven randomly drawn samples from a population and we want to determine whether there are any significant differences among their 7(7 - 1) means. This will require computation of 21 t-tests to determine the 2 significance of difference between the seven means by taking two means at a time. This procedure is time consuming as well as cumbersome. The technique of analysing of variance is applied to determine if any two of the seven means differ significantly from each other by a single test, known as F-test, rather than 21 t-tests. The F-test makes it possible to determine whether the sample means differ from one another (between group variance) to a greater extent than the test scores differ from their own sample means (within group variance). Using the ratio:

$$F = \frac{\text{Variance between the groups}}{\text{Variance within groups}} \dots\dots\dots(8)$$



The values of F-ratio are given in the Appendix (Table 111). This table indicates the F-critical values necessary for rejecting the null hypothesis at selected levels of significance, usually .05 and .01 levels.

### E Factor Analysis

Factor Analysis is a technique mainly used in research in psychology and education. We will not examine factor analysis in detail here, but very briefly describe about this technique. Factor Analysis is a procedure for determining the number and nature of constructs that underlie a set of measures (Wiersman 1986). Construct as you are already aware, is a trait or an attribute that explains some phenomena, e.g., anxiety, intelligence, motivation, attitude etc. In factor analysis artificial variables are generated and called factors. These factors represent the constructs. Factor analysis is initiated from the correlation matrix of the variables. The variables which are highly correlated are grouped together. Suppose we have scores of 100 students on 10 different tests. Question here is - How many different traits or constructs these 10 tests measure. The possibility can be that three or four tests measure the same trait or one test may measure two or more traits. The researcher can determine the correlation co-efficient among the 10 different tests. High correlation's between test scores indicate that common constructs are measured. Low or zero correlation's indicate the absence of common constructs. There are few terms used in factor analysis. If a test measures only one construct, it is labeled as factorially pure. A factorially complex test is one that measures two or more factors. Factor loading is the extent to which a test measures a factor. Factor loading is very important in factor analysis because factors which are artificial variables generated from the data must be described and integrated. It is a correlation coefficient between a test and a factor.

### Uses of Factor Analysis in Research

## Notes

In conducting research, the aim of using factor analysis is to identify the nature and number of constructs that underlie a set of variables. Factor analysis is associated with construct related evidence when establishing validity. Factor analysis is used in confirmatory analysis and exploratory analysis. Confirmatory factor analysis is used in studies where hypothesized constructs measured by a set of variables are either confirmed or refuted. It is also used to analyze a single test by factor analyzing the item scores. Suppose 50 item test measures three traits, a confirmatory analysis of the item scores would support or refute this proposition. On the other side, in the exploratory analysis the number of variables are reduced to a manageable number of explanatory purposes. A set of measures can be factor analyzed to enhance the explanation of what is measured in a more parsimonious manner. Eg. A group of teachers in an open university were observed and measured on 2 different competencies. A factor analysis of the competency scores undoubtedly would generate a smaller number of factors, say three or four, that represents the constructs underlying the performance of teacher. Thus factor analysis in any research analysis provide valuable insights into the nature of phenomena.

### Check Your Progress 3

Notes: (a) Space is given below for writing your answer.

(b) Compare your answer with the given at the end of the unit.

1. The following scores were obtained on an interest test for 5 males and 8 females of PGDDE enrolled with IGNOU. Male : 20,22, 30, 32, 26. Female: 34,25,16,30,22,27,20,26.

Is the difference between the Mean Interest Scores of the Males and the Females significant?

.....  
.....  
.....

Basic Assumptions for the Analysis of Variance Certain basic assumptions underlying the technique of analysing of variance are :

1. The population distribution should be normal. This assumption is, however, not so important. The study of Norton (Guilford, 1965; pp.300-301) also points out that F is rather insensitive to variations in the shape of population distribution.
2. All groups of a certain criterion or of the combination of more than one criterion should be randomly chosen from the sub-population having the same criterion or having the same combination of more than one criterion. For example, if we wish to select two groups from a study centre, one belonging to a rural area and the other to the urban area, we must, choose the groups randomly from the respective sub-populations.
3. The sub-groups under investigation must have the same variability. In other words, there should be homogeneity of variance. It is tested either by applying Bartlett's test of homogeneity or by applying Hartley's test. To illustrate the use of F-test, let us consider an example of twenty distance learners who have been randomly assigned to 4 groups of 5 each, to be taught by different methods, i.e., A, B, C and D. Their performance scores on an achievement test, administered after the completion of experiment are given in Table 10.3.

Table 10.3 Achievement Test Scores of the Four Groups Taught through Four Different Methods

	Methods or Groups				
	A ( $X_1$ )	B ( $X_2$ )	C ( $X_3$ )	D ( $X_4$ )	
	14	19	12	17	
	15	20	16	17	
	11	19	16	14	
	10	16	15	12	
	12	16	12	17	
$\sum X$	62	90	71	77	300
$\sum X^2$	786	1634	1025	1207	4652

## Notes

We may compute the analysis of variance using the following steps:

$$1. \quad \text{Correction} = \frac{(\sum X)^2}{N} = \frac{(300)^2}{20} = 4500$$

$$2. \quad \text{Total sum of squares (Total SS)}$$

$$= \sum X^2 - \text{Correction}$$

$$= 14^2 + 15^2 + \dots \dots \dots 12^2 + 17^2 - 4500$$

$$= 4652 - 4500$$

$$= 152$$

$$3. \quad \text{Sum of squares between means of treatments (Methods) A, B, C, and D (between means):}$$

$$= \frac{(\sum X_1)^2}{N_1} + \frac{(\sum X_2)^2}{N_2} + \frac{(\sum X_3)^2}{N_3} + \frac{(\sum X_4)^2}{N_4} - \text{Correction}$$

$$= \frac{(62)^2}{5} + \frac{(90)^2}{5} + \frac{(71)^2}{5} + \frac{(77)^2}{5} - 4500$$

$$= 4582.8 - 4500$$

$$= 82.8$$

$$4. \quad \text{Sum of squares within treatments (Methods) A, B, C, and D (SS within means):}$$

$$= \text{Total SS} - \text{SS between means}$$

$$= 152 - 82.8$$

$$= 69.2$$

$$5. \quad \text{Calculation of variances from each SS and analysis of the total variance into its components.}$$

Each SS becomes a variance when divided by the degrees of freedom (df) allotted to it. There are 20 scores in all in Table 3, and hence there are (N-1) or (20-1) = 19 df in all. These 19 df are allocated in the following ways: If N = number of scores in all and K = number of treatments or groups, we have df for total SS = N-1 = 20-1 = 19, df for within treatments = N-K = 20-4 = 16; and df for between the means of treatments = K- 1 = 4- 1 = 3. - The variance among means of treatments is 82.813 or 27.60; and the variance within means is 69.2116 or 4.33.

The summary of the analysis of variance may be presented in tabular form as follows:

Table 10.4 : Summary of Analysis of Variance

Source of Variance	df	Sum of Squares (SS)	Mean Square (Variance)
Between the means of treatment	3	82.8	27.60
Within treatment	16	69.2	4.33
Total	19	152.0	

Using formula (8)

$$F = \frac{27.60}{4.33} = 6.374$$

In the present problem, the null hypothesis asserts that four sets of scores are in reality the scores of four random samples drawn from the same normally distributed population, and that the means of the four groups A, B, C, and D will differ only through fluctuations of sampling. For testing this hypothesis, we divided the 'between means' variance by the 'within treatments' variance and compared the resulting variance ratio, called F, with the F-values in Table 111. The F value of 6.374 in the present case is to be checked for table value for df 3 and 16 (the degrees of freedom for numerator and denominator). The table values for .05 and .01 levels of significance are 3.24 and 5.29. Since the computed F-value of 6.374 is greater than the table values, we reject the null hypothesis and conclude that the means of the four groups differ significantly.

---

## 10.4 NON-PARAMETRIC TESTS AND APPLICATION OF CHI-SQUARE TEST

---

In the preceding section, we described/discussed some important parametric tests involving the assumptions based upon the nature of the population ' distribution. There are some test's which do not make

## Notes

numerous or stringent assumptions about the nature of the population distribution. These tests are known as distribution-free or non-parametric tests. The non-parametric tests are based upon the following assumptions:

1. The nature of the population, from which samples are drawn, is not known to be normal.
2. The variables are expressed in nominal form, that is, classified in categories and represented by frequency counts.
3. . The variables are expressed in ordinal form, that is, ranked in order or expressed in numerical scores which have the strength of ranks.
4. The sample sizes are small. The most frequently used non-parametric tests are: Chi-square test, the median test, the sign test, the Mann-Whitney U test, the KolmogorovSmirnov Two Sample Test, the Wilcoxon-Matched-Pairs Signed-Ranks Test, the McNemar Test for' Significance of changes, contingency coefficient, etc. For the present unit, we will discuss the applications of Chisquare 'test and median test only.

### A Application of Chi-square test

The Chi-square (pronounced as Ki-square) test is used with discrete data in the form of frequencies. It is a test of independence and is used to estimate the likelihood that some factor other than chance accounts for the observed relationship. Since the null hypothesis states that there is no relationship between the variables under study, the Chi-square test merely evaluates the probability that the observed relationship results from chance. The formula for Chi-square ( $\chi^2$ ) is:

$$\chi^2 = \sum \left[ \frac{(f_o - f_e)^2}{f_e} \right] \dots\dots\dots (9)$$

in which

$f_o$  = frequency of the occurrence of observed or experimentally determined facts

$f_e$  = expected frequency of occurrence

For example consider the following data (in Table 10.5) of 500 distance learners who have been categorised into three groups, elder, middle-aged and younger on the basis of age and their preference for four colours, red, blue, yellow and green.

Table 10.5 The Chi-square Test of Independence in Contingency Table

<i>Colour</i> → <i>Age Group</i> ↓	<i>Red</i>	<i>Blue</i>	<i>Yellow</i>	<i>Green</i>	<i>Total</i>
Elder	40 (38.42)	50 (45.22)	35 (39.10)	45 (47.26)	170
Middle	35 (36.16)	42 (42.56)	44 (36.80)	39 (44.48)	160
Younger	38 (38.42)	41 (45.22)	36 (39.10)	55 (47.26)	170
Total	113	133	115	139	500

Across the first row of the table, we find that out of 170 distance learners in the older age-group, 40 have given their preference for red colour, 50 for blue, 35 for yellow, and 45 for green. Reading down the first column, we find that out of 113 distance learners giving preference for red colour, 40 belong to the older age-group, 35 to middle and 38 to younger age-group. The other columns and rows are interpreted in the same way. The hypothesis to be tested is the null hypothesis, that is, age and colour preferences are essentially unrelated or independent. To compute Chisquare we must calculate an independent value, i.e., expected frequency for each cell in the contingency table. Independent values are represented by the figures in parentheses within the different cells. They give the number of students whom we should expect to fall in a particular age-group, showing their preference for a particular colour in

## Notes

the absence of any real association. The calculation of expected frequencies ( $f_e$ ) and Chi-square ( $\chi^2$ ) are shown as under:

### 1. Calculation of expected frequencies ( $f_e$ )

$$\text{Row I: } \frac{113 \times 170}{500} = 38.42; \quad \frac{133 \times 170}{500} = 45.22;$$

$$\frac{115 \times 170}{500} = 39.10; \quad \frac{139 \times 170}{500} = 47.26;$$

$$\text{Row II: } \frac{113 \times 160}{500} = 36.16; \quad \frac{133 \times 160}{500} = 42.56;$$

$$\frac{115 \times 160}{500} = 36.80; \quad \frac{139 \times 160}{500} = 44.48;$$

$$\text{Row III: } \frac{113 \times 170}{500} = 38.42; \quad \frac{133 \times 170}{500} = 45.22;$$

$$\frac{115 \times 170}{500} = 39.10; \quad \frac{139 \times 170}{500} = 47.26.$$



## 2. Computation of the Chi-square value

Using formula (9)

$$\begin{aligned}
 X^2 &= \sum \left[ \frac{(f_o - fe)^2}{fe} \right] \\
 &= \frac{(40 - 38.42)^2}{38.42} + \frac{(50 - 45.22)^2}{45.22} + \frac{(35 - 39.10)^2}{39.10} + \frac{(45 - 47.26)^2}{47.26} + \\
 &\quad \frac{(35 - 36.16)^2}{36.16} + \frac{(42 - 42.56)^2}{42.56} + \frac{(44 - 36.80)^2}{36.80} + \frac{(39 - 44.48)^2}{44.48} + \\
 &\quad \frac{(38 - 38.42)^2}{38.42} + \frac{(41 - 45.22)^2}{45.22} + \frac{(36 - 39.10)^2}{39.10} + \frac{(55 - 47.26)^2}{47.26}
 \end{aligned}$$

$$X^2 = 5.182$$

$$\begin{aligned}
 3. \quad df &= (r-1)(c-1) \\
 &= (3-1)(4-1) \\
 &= 8
 \end{aligned}$$

The  $\chi^2$  critical values for 8 df as given in Table IV (See Appendix) are 15.507 and 20.090 respectively for .05 and .01 levels of significance and the obtained value, 5.182, is less than the table value even at .05 level. This indicates that there is no relationship between the age and the colour preference and thus the hypothesis that age and colour preference are essential! independent may be accepted at .05 level of significance. In the case of 2 x 2 contingency table, with  $(r-1)(c-1) = 1$  df, there is no need of computing the expected frequencies (independence values) for each cell. The formula is:

$$X^2 = \frac{N (|AD - BC|)^2}{(A+B)(C+D)(A+C)(B+D)} \quad \dots\dots\dots (10)$$

In the above formula A, B, C and D are the frequencies in the first, second, third and fourth cells respectively and the vertical lines in  $|AD - BC|$  mean that the difference is to be taken as positive. To illustrate the use of formula (10), let us determine whether item 5 of an achievement

## Notes

test differentiates between high and low achievers. The responses to items are given in the following 2 x 2 contingency Table.

Table 10.6 The Chi-square Test in 2 x 2 Fold Contingency Table

	Passed item 5	Failed item 5	Total
	(A)	(B)	(A + B)
High Achiever	115	35	150
	(C)	(D)	(C + D)
Low Achiever	40	90	130
	(A+C)	(B+D)	
Total	155	125	280

Using formula (10)

$$X^2 = \frac{280(|115 \times 90 - 35 \times 40|)^2}{(115 + 35)(40 + 90)(115 + 40)(35 + 90)}$$

$$= \frac{280(|10350 - 1400|)^2}{(150)(130)(155)(125)} = 59.36$$

Since the computed  $\chi^2$  value of 59.36 exceeds the critical  $\chi^2$  value of 6.635 to be significant at .01 level, we reject the hypothesis that item 5 of the test, does not discriminate significantly between high and low achievers. In other words, it may be concluded that item 5 of the achievement test discriminates significantly between the two groups, namely high and low achieving students. Further, when entries in 2x2 table are less than 10, Yate's correction for continuity is applied to formula (10). The corrected formula reads:

$$X_c^2 = \frac{N(|AD - BC| - N/2)}{(A + B)(C + D)(A + C)(B + D)} \dots\dots\dots (11)$$

The following example illustrates the use of formula (11).

Fifteen male and twelve female counsellors of a study centre were asked to express their attitude towards population education. Both the groups of counsellors were administered the attitude scale and were classified as

having either positive or negative attitude towards population education. The distribution of the sample is shown in table 10.7.

Table 10.7 Distribution of Male and Female Counsellors in Terms of their Positive or Negative Attitude towards Population Education.

	Positive Attitude	Negative Attitude	Total
	(A)	(B)	(A + B)
Female Counsellors	7	5	12
	(C)	(D)	(C + D)
Male Counsellors	9	6	15
	(A+C)	(B+D)	
Total	16	11	27

$$\begin{aligned}
 X_c^2 &= \frac{27(|7 \times 6 - 5 \times 9| - 27/2)^2}{(7+5)(9+6)(7+9)(5+6)} \\
 &= \frac{27(|42 - 45| - 13.5)^2}{12 \times 15 \times 16 \times 11} \\
 &= 0.23
 \end{aligned}$$

Since the obtained value of  $x^*$ , 0.23, is less than the table value of 3.842 to be significant at .05 level of significance, it may be inferred that there is no true difference in the attitude of male and female counsellors towards population education.

### B Application of Median test

The median test is used for testing whether two independent samples differ in central tendencies. It gives information as to whether it is likely that two independent samples have been drawn from populations with the same median. It is particularly useful whenever the measurements for two samples are expressed in an ordinal scale. In the median test, we first compute the combined median for all rank measures in both samples.

## Notes

Then both sets of rank measures :it the combined median are dichotomized and the data are set in a 2x2 Table as show in Table 10.8

Table 10.8: 2x2 Table for use of the Median test

	Group I	Group II	Total
No. of measures above the combined median	(A)	(B)	(A+B)
No. of measures below the combined median	(C)	(D)	(C+D)
<b>Total</b>	(A+C)	(B+D)	

Under the null-hypothesis, we would expect about half of each group's measures +o be above the combined median and about half to be below it, that is, we would expect frequencies A and C to be equal, and frequencies B and D to be nearly equal. In order to test this null-hypothesis, we calculate, f using the formula (12):

$$f^2 = \frac{N(|AD - BC| - N/2)^2}{(A+B)(C+D)(A+C)(B+D)} \dots\dots\dots(12)$$

Let us illustrate the use of the formula (12) in the following example. Eighteen male and female distance learners of a study centre of an open university were asked to express their attitude towards the functioning of a study centre. Both the groups were administered an attitude scale and a common median attitude measure was worked out. The number of cases from both the groups falling above and below the median point is shown in Table 10.9

Table 10.9 Distribution of Distance Learners (Male and Female) Below and Above the Common Median

	Below Median	Above Median	Total
Female Distance Learners	10	4	14
Male Distance Learners	6	12	18
Total	16	16	32

Using formula (12),

$$f^2 = \frac{32 \left( |(10)(12) - (6)(4)| - \frac{32}{4} \right)^2}{16 \times 16 \times 18 \times 14}$$

$$= 3.17$$

Since the obtained value 3.17 for 1 df does not exceed the f 'critical value of 3.84 for a two-tailed test at .05 level, the null hypothesis is retained and we may conclude that there is no difference in the attitude of male and female distance learners towards the functioning of the centre.

---

## 10.5 LET US SUM UP

---

In this unit we described the nature of parametric and non-parametric tests along with the assumptions they are based on. The applications of some parametric tests (Z, t and F) and non-parametric tests (Chi-square) in the analysis of data have also been discussed. 1) Parametric tests are used in the analysis of data available in interval or ratio scales of measurement. . 2) Parametric tests assume that the data are normally or nearly normally distributed. 3) Z-test, t-test and F-test are the most commonly used parametric tests. 4) If large-sized samples, i.e., those greater than 30 in size, are selected at random from an infinite population, the distribution of sample means is called the 'sampling distribution of means'. This distribution is normal and it possesses all the following characteristics of a normal distribution: i) The average value of sample means will be the same as the mean of population. ii) The distribution of the sample means around the population mean has its own standard deviation which is known as the 'standard error of the mean' (SE, or  $\sigma_{\bar{x}}$ ). iii) The sampling distribution is centered at the unknown population mean with its standard deviation 0.50. iv) The sample means often fall equally on the positive and negative sides of the population mean. v) About 213 of the sample means (exactly 68.26 per

## Notes

cent) will lie within  $\pm 1 \cdot 00 \sigma$ , of the population mean, i.e., within a range of  $\pm 1 \times 0.50$  or  $0.50$  of the population mean.

vi) 95 of the 100 sample means will lie within  $\pm 1.96 \sigma_M$  of the population mean, i.e., 95 of 100 sample means will lie within  $\pm 1.96 \times 0.50$  or  $\pm 0.98$  of the population mean. vii) 99 of the 100 sample means

will be within  $\pm 2.58 \sigma_M$  of the population mean, i.e., 99 of our sample means will lie within  $\pm 2.58 \times 0.50$  or  $\pm 1.29$  of the population mean. 5)

If we draw a large sample randomly from a population and compute its mean, the mean has 95 per cent chance of being within  $1.96 \sigma_M$  units from the population mean. Also, there is a 99 per cent chance that the sample mean lies within  $2.58 \sigma_M$  units from the population mean. To be

more specific, there is 95 per cent probability that the limits  $M \pm 1.96 \sigma_M$  enclose the population mean and 99 per cent probability that the limits  $M \pm 2.58 \sigma_M$  enclose the population mean. 6) The limits  $M \pm 1.96 \sigma$ , and  $M \pm 2.58 \sigma_M$  are called 'confidence intervals' for .05 and .01 levels of

confidence respectively. 7) When a 'statistic' is used to estimate, a 'parameter', the number of 'degrees of freedom' depends upon the

restrictions placed. Therefore, the number of degrees of freedom (df) will vary from one statistic to another. In estimating the population mean from the sample mean, for example, 1 df is lost and so the number of

degrees of freedom is  $N - 1$ . 8) The 2-test is used for testing the significance of the difference between the means of two large samples. 9)

The t-test is used for testing the significance of the difference between the means of two small samples. 10) Under the null hypothesis, the difference between the sample means may be either plus or minus and as often in one direction as in the other from the true (population) difference

of zero, so that in determining probabilities we take both tails of the normal sampling distribution and make use of two-tailed test. But in

many situations our primary concern is with the direction of the difference rather than with its existence in absolute terms. In such

situations, we make use of onetailed test. 11) We use Z-test and t-test to determine whether there is any significant difference between means of

two random samples. But when the number of samples is more than two, F-test, based on the technique of analysis of variance, is used for testing

the significance of the sample means. 12) Non-parametric tests are used

in the analysis of non-parametric data, i.e., when the data are available in nominal or ordinal scales of measurement. 13) Non-parametric tests are distribution-free tests and do not rest upon the more stringent assumptions of normally distributed population. 14) Chi-square test, median test, Sign-test, Mann-Whitney U test, Wilcoxon Matched Pairs Signed Ranks test and Kolmogorov Smirnov two Sample test are examples of some non-parametric tests. 15) Chi-square is used with discrete data in the form of frequencies. It is a test of independence, and is used to estimate the likelihood that some factor other than 'chance' accounts for the observed relationship between the variables. 16) Median test is used for testing whether two independent samples differ in central tendencies. It is particularly useful whenever the measurements for the two samples are expressed in an ordinal scale. Now, use the following check list and see whether you have learnt to:

- classify various statistical tests,
- describe the nature of parametric tests along with the assumptions on which
- they are based,
- define sampling distribution of means,
- define the standard error of mean,
- define the confidence intervals and levels of confidence,
- compute 0.95 and 0.99 confidence intervals for the true mean from a large
- sample mean,
- define and illustrate the concept of degrees of freedom,
- compute 0.95 and 0.99 confidence intervals for the true mean from the
- sample mean,
- apply Z-test for testing the significance of the difference between means of
- two independent large samples involving: (i) one-tailed and (ii) two-tailed
- tests,

## Notes

- apply t-test for testing the significance of the difference between means of
- two independent small samples involving: (i) one-tailed and (ii) two-tailed
- tests,
- describe the nature and uses of the analysis of variance,
- state the basic assumptions of the technique of analysis of variance,
- apply F-test for testing the significance of the difference between means,
- describe the nature of the non-parametric tests along with their assumptions,
- name various non-parametric tests,
- describe the use of Chi-square test,
- illustrate the application of Chi-square test, and
- describe the use of median test.

---

## 10.6 KEY WORDS

---

**Parametric Tests** : these are statistical tests which are used for analyzing parametric data and making inferences about the parameters from the statistics. These tests are based upon certain assumptions about the nature of data distributions and the types of measure used.

**Non-parametric Tests** : these are statistical tests which are used for analyzing non-parametric data and making possible useful inferences without any assumptions about the nature of data distributions.

**Standard Error of Mean** : it is the standard deviation of a distribution of sample means.

**Degrees of Freedom** : the number of degrees of freedom in a distribution is the number of observations or values that are independent of each other, and cannot be deduced from each other.

---

## 10.7 QUESTIONS FOR REVIEW

---

- 1) Describe the assumptions on which the use of parametric tests is based.



---

## 10.8 SUGGESTED READINGS AND REFERENCES

---

- Fisher, R. A., *Statistical Methods for Research Workers*, New York: Hagner Publishing, 1950
- Guilford, J.P., *Psychometric Methods*, New York: McGraw Hill, 1954.
- Wiersmar, William, *Research Methods in Education: An Introduction*, Allyn and Bacon, Inc. Massachusetts, 1986

---

## 10.9 ANSWERS TO CHECK YOUR PROGRESS

---

### Check Your Progress 1

- 1) Parametric data should be used if the following basic assumptions are met. These assumptions are based on the nature of the population distribution and on the way the scale is used to quantify the data observations.
  - i) The observations are independent. The selection of one case is not dependent Biostatistics upon the selection of any other case.
  - ii) The population values are normally distributed.
  - iii) The samples have equal or near variances. This condition is known as equality or homogeneity of variances and is particularly important to determine when the samples are small.
  - iv) The variables described are expressed in interval or ratio scale and not in nominal or ordinal scales of measurement.
- 2) The standard error of Mean is

## Notes

$$\begin{aligned}
 SE_M &= \sigma_M = \frac{\sigma}{\sqrt{N}} \\
 &= \frac{5.82}{\sqrt{100}} \\
 &= \frac{5.82}{10} \\
 &= 0.582.
 \end{aligned}$$

### Check Your Progress 2

1) Find it

Boys ( $N_1 = 5$ )			Girls ( $N_2 = 10$ )		
$X_1$	$x_1$	$x_1^2$	$X_2$	$x_2$	$x_2^2$
20	-6	36	34	9	81
22	-4	16	25	0	0
30	4	16	16	-9	81
32	6	36	30	5	25
26	0	0	22	-3	9
$\Sigma X_1 = 130$		$\Sigma X_1^2 = 104$	27	2	4
			20	2	25
			26	-5	1
			$\Sigma X_2 = 200$	1	$\Sigma X_2^2 = 226$
Mean = $M_1 = \frac{\Sigma X_1}{N_1} = \frac{130}{5}$			Mean = $M_2 = \frac{\Sigma X_2}{N_2} = \frac{200}{8}$		
= 26.00			= 25.00		

ii)  $df = N_1 + N_2 - 2 = 5 + 8 - 2 = 11$

iii) Using formula (7)

$$\begin{aligned}
 t &= \frac{|M_1 - M_2|}{\sqrt{\left( \frac{\Sigma x_1^2 + \Sigma x_2^2}{N_1 + N_2 - 2} \right) \left( \frac{N_1 + N_2}{N_1 \times N_2} \right)}} \\
 &= \frac{|26 - 25|}{\sqrt{\left( \frac{104 + 226}{5 + 8 - 2} \right) \left( \frac{5 + 8}{5 \times 8} \right)}} \\
 &= \frac{1}{\sqrt{9.75}} = 0.32
 \end{aligned}$$

iv) We used a two-tailed test as we are not hypothesizing a direction of the difference between the means. The t-values as given in Table I1 in the Appendix for 11 df for .05 and .01 columns are 2.20 and 3.1 respectively. Since the obtained value of 0.32 is less than these table values, the difference

between the mean interest scores of boys and girls is not significant.

### Check Your Progress 3

1) Non-parametric tests are distribution-free tests and are based on the following assumptions:

- i) The nature of the population distribution, from which samples are drawn is not known to be normal.
- ii) The variables are expressed in nominal form (classified in categories and represented by frequency counts).
- iii) The variables are expressed in ordinal form (ranked in order).

2) i) The number of the males and the females who have passed or failed the test item is given in the following 2 x 2 table.

	Number Passed	Number Failed	Total
Female	30 (A)	20 (B)	50 (A+B)
Male	25 (C)	15 (D)	40 (C+D)
<b>Total</b>	55 (A+C)	35 (B+D)	90

ii) Using formula (10)

$$\begin{aligned}\chi^2 &= \frac{N(AD-BC)^2}{(A+B)(C+D)(A+C)(B+D)} \\ &= \frac{90(1(30)(15) - (20)(25))^2}{(30+20)(25+15)(30+25)(20+15)} \\ &= 0.58\end{aligned}$$

iii) Since the obtained value 0.058 of  $\chi^2$  does not exceed the Table value 3.841 of  $\chi^2$  at .05 level of significance, we may conclude that the test item does not differentiate between the two groups of males and females.

---

# **UNIT 11: QUALITATIVE RESEARCH: THEORETICAL SAMPLING, CASE STUDIES**

---

## **STRUCTURE**

- 11.0 Objectives
- 11.1 Introduction
- 11.2 Definition of Qualitative Research
- 11.3 Characteristics of Qualitative Research
- 11.4 Reasons for Conducting Qualitative Research
- 11.5 Types of Qualitative Methods
  - 11.5.1 Biographical Studies
  - 11.5.2 Phenomenological Study
  - 11.5.3 Grounded Theory Study
  - 11.5.4 Ethnography
  - 11.5.5 Case Study
- 11.6 Context of Qualitative Methods
- 11.7 Common Steps of Conducting Qualitative Studies
- 11.8 Verification of Trustworthiness of Qualitative Research
- 11.9 Let us sum up
- 11.10 Key Words
- 11.11 Questions for Review
- 11.12 Suggested readings and references
- 11.13 Answers to Check Your Progress

---

## **11.0 OBJECTIVES**

---

After going through this unit, you will be able to:

- To define qualitative research
- To discuss needs and relevance of qualitative research in education

- To explain main features of different types of qualitative inquiry viz, biography,
- To phenomenology, grounded theory approach, ethnography and case study
- To describe basic procedural details of such methods
- To design qualitative research in the context of any educational problem.

---

## **11.1 INTRODUCTION**

---

As you know, events concerning education phenomena are very complex in nature. They involve mostly sociological and behavioral manifestations of human relations in specific contexts. The methods adopted to explore the meaning and processes of educational phenomena are basically qualitative in nature. Most of the information in education is in the form of verbal and other symbolic behavior. The verbal data gathered through questionnaires, observation or interviews are mostly qualitative in nature. These data provide depth and detail which emerge through direct quotation and careful description of phenomena. Unlike scientific explanations to educational phenomena qualitative research concentrates on understanding the process of dynamic situations and intervenes with the reality in its natural setting. The procedures adopted are of humanistic nature. Of late, major attention has been paid by educational researchers to qualitative methods. Having their origin in the typical nature of inquiries in sociology, anthropology, political science and psychology, qualitative methods have been given refined shape through current literature. In this unit, you will be exposed to the nature of qualitative research, its types, context and steps involved in conducting qualitative research.

---

## **11.2 DEFINITION OF QUALITATIVE RESEARCH**

---

As you know, the emergence of qualitative methods of research is of recent origin. Social scientists, anthropologists and sociologists have

## Notes

given a shape to the concept of qualitative research mostly during the end of twentieth century. The procedural clarity is still in the process of refinement. Clarity of the definitions of qualitative research is as follows:

1. In the initial stage of conceptual analysis of qualitative research there was a trend of defining qualitative research as an opposite pole of scientific (quantitative) inquiry. Quantitative research deals with a few variables and many cases whereas qualitative research is came out with a few cases and many variables.

2. The above definition focuses on the scope of qualitative methods of inquiry. However, from the point of view of methodological and process dimensions the following two definitions may clarify your understanding of qualitative research. Cresswell (1998) defined qualitative research in a similar tone. To him "qualitative research is an inquiry process of understanding based on distinct methodological traditions of inquiry that explore a social or human problem. The researcher builds a complex holistic picture, analyses words, reports detailed views of informants and conducts the study in a natural setting".

3. Denzin and Lincoln (1994) say "Qualitative research is multi-method focus, involving an interpretive, naturalistic approach to its subject matter Qualitative research involves the studied use and collection of a variety of empirical materials - case study, personal experience, introspection, life story, interview, observational, historical, interactional and visual texts-that describe routine and problematic and meaning in individuals' lives". You may trace at least four clear meanings from the above definitions. One, qualitative research focuses on the study of complex human and social problems in totality, unlike scientific method of concentrating on study of fragmented variables or situations or events. Two, qualitative research places the main emphasis on the researcher who narrates and interprets phenomena in terms of meanings derived from people's experiences, events etc. Hence, human and subjective approach is highlighted. Three, the studies are conducted in a natural setting i.e., to observe the events without making any manipulations or controls on variables studied. Four, it involves a variety of data gathering

techniques and approaches of qualitative nature viz., case study, interviews, dialogues, observations, personal experience, life story, visual data like photography etc. These data are gathered from varieties of Qualitative Research sources.

---

## 11.3 CHARACTERISTICS OF QUALITATIVE RESEARCH

---

The following are the characteristics of qualitative research:

- i) **Multiple realities** First, qualitative research assumes that there exist multiple realities in social and educational situations. These realities exist in concrete forms. They are perceived by people differently and thus become different mental constructs for different people. In other words, realities are taken to be what people perceive them to be at a particular point of time. Since social and educational situations keep on changing from time to time, the realities too keep on changing. Furthermore, since the realities are context specific, they cannot be tangible in a generalized form.
- ii) **Meanings and interpretations** Qualitative research emphasises on study of meanings given to or interpretations made about objects, events and processes concerning educational situations. To them changes in terms of social and behavioral phenomena cannot be identified with the concept of physical movements which can be identified by external observation alone. Rather the understanding of human behaviour or a social phenomenon involves understanding of how men are, what they are doing or participating in an activity.
- iii) **Generation of knowledge** 1 Qualitative inquiry insists on generation of knowledge resulting from the interaction between the inquirer and the respondents. The respondents answer the questions put by the inquirer in terms of their perception or meanings they attach to their actions. Moreover,

interactions take place between the inquirer and his/her respondents to achieve maximum levels of responsiveness and insights concerning the problem under investigation.

- iv) **Generalization** As stated above, the researchers do not believe in the process of generalisation as propounded by scientists. They argue that in the process of making a generalisation a lot of meaningful information existing in individual units is undermined; hence generalized knowledge does not represent real knowledge. For them, the process of knowledge generation must take into account the differences or real evidences existing in different specific situations.
- v) **Human relations** In the case of human relations, several intrinsic factors, events and processes keep on influencing each other constantly. Therefore, it is not possible to identify one to one cause and effect relationships at this level of naturalistic studies. The causality in social sciences cannot be demonstrated in the 'hard' sense as it is done in the case of physical sciences. Rather, only patterns of plausible influences can be inferred from social and behavioural studies.
- vi) **Value systems** Qualitative Researchers do not believe in value-free inquiry. The influence of value systems is recognized in the identification of problems, selection of samples, use of tools, data collection, the conditions in which data are gathered, and the possible interaction that takes place between the inquirer and the respondents.

---

## **11.4 REASONS FOR CONDUCTING QUALITATIVE RESEARCH**

---

As a research student you may be curious to explore in what circumstances one opts for qualitative methods of inquiry in comparison



to quantitative methods of inquiry. Some of the situations which prompt a researcher to opt for the qualitative method are:

- i) As you know, there are different kinds of research questions about the phenomena under investigation, such as What happens? How does it take place and Why did it take place? or What are the causes of certain effect's? The first two types of questions involving What (The present context and status) and How (the process) are dealt with very well through qualitative methods. Unlike quantitative inquiry, where you may intend to explain causal relationships (Why question) through comparison of groups or to establish associations between different variables, the qualitative method focuses on exploring the intricacies involved in process dimensions of an event in totality. For example, if you are interested in understanding the curricular practices of best quality universities, the qualitative method will suit you best.
- ii) Qualitative method is also preferred in the context of initial level of theoretical research. In social sciences, particularly in education, theories are not available to explain many completed events. The variables concerning the events are not easily identifiable. In such cases the qualitative method is used to identify significant variables associated with the event. For instance, formal education practices taking place in a tribal setup vis-8-vis human and social development need to be tackled through qualitative method since different variables involved in the situation need to be explored. Identification of different variables and trends lead towards refinement and theoretical explanation.
- iii) There are certain research problems which need holistic treatment. In other words, many variables need to be studied in depth in the context,of one case or unit. For instance, designing a grass root level plan for educational development

## Notes

of a village would require application of a qualitative method of inquiry where multifaceted data need to be gathered through various sources for intervention and development which treat the uniqueness of the concerned village. iv) Unlike experimental studies- where a researcher intends to explain a cause and effect relationship in a controlled laboratory setting, the qualitative method is preferred to conduct the study without disturbing neutrality in setting. For instance, curricular practices of best quality schools need to be explored in natural settings through adopting qualitative methods like participant observations, visual documents, verification of project works, interviewing students, teachers and parents, etc.

- iv) One of the major limitations of quantitative inquiry is associated with the study of limited variables. Moreover, those variables need to be amenable to measurement. However, if your research problem is of such type where variable are not easily amenable to measurement and extensive data situations, you will have to opt for qualitative inquiry.
  
- v) The nature of certain problems is such that your role as a researcher plays a vital subjective role than objective and impersonal role. For instance, in the process of exploring meaning of certain behavioural exposures with its motives your role is to be valued with a lot of significance since you will have to describe the event from the participant's point of view. Moreover, in such a process you can explore the reality, being a part of it as well as interacting with situations. You will have to depict the situations as an active learner and not by passing judgment as an expert. Such kind of studies demand application of qualitative methods in comparison to quantitative methods.

**Check Your Progress 1**

Notes: a) Space is given below for writing your answers.

b) Compare your answers with those given at the end of the unit.

- 1. Differentiate between quantitative studies and qualitative studies.

.....  
.....  
.....  
.....  
.....  
.....

- 2. What are the basic factors of qualitative studies?

.....  
.....  
.....  
.....  
.....

- 3. Why should you choose the qualitative method in educational research?

.....  
.....  
.....  
.....  
.....

---

**11.5 TYPES OF QUALITATIVE METHODS**

---

Many varieties of traditions of qualitative studies exist in social sciences. They have been categorized by Creswell(1998) in the context of their forms, terminologies and focus as under:

## Notes

- Biographical study
- Phenomenological study
- Grounded Theory study .
- Ethnography
- Case study

### 11.5.1 Biographical Studies

#### Main features

Biographical study is the study of an individual and his or her experiences as narrated to the researcher or found in different sources. You can come across biographical writings in different fields like literature, history, anthropology, sociology, education, psychology etc. Biographies are presented with different perspectives like literary, historical, anthropological, sociological, educational, psychological as well as of interdisciplinary nature. The focus of biography remains on telling and inscribing the stories of others. It explores history of life e.g. accounts of major achievements of life. There are different connotations linked with biographical study viz., individual biographies, autobiography, life history, oral history. In all these cases the researchers must take care of objectivity in expression with little research interpretation. It must be written in a scholarly way with a strong historical background of the subject and chronological organisation. The account must be presented artistically from the perspective of presenting details in a lively and interesting manner.

#### Steps of Biographical Studies Qualitative Research ;

Different steps involved in writing biography are:

- (i) The first step is to identify the experiences in an individual's life to be arranged chronologically in different stages of life.
- (ii) The researcher gathers concrete contextual materials through interviews. He gathers stories.

- (iii) The stories are organized around different themes. The themes indicate major events of individual's life.
- (iv) The researcher explores the meaning of these stories.
- (v) He also looks for larger structures to explain cultural issues etc.

## 11.5.2 Phenomenological Study

### Main features

Phenomenological study's focus remains on describing the meaning of live experiences for several individuals about a concept or the phenomenon. It is said that through phenomenological approach the researcher explores the structures of consciousness in human experiences. Here experiences contain both the outward appearance and inward consciousness based on memory image and meaning. Steps of Phenomenological Studies The procedural details of phenomenological studies are listed as under:

- i) At first the researcher must recognize his or her own pre-conceived ideas about the phenomena to understand it through the voices of several experienced persons.
- ii) Second, he writes research questions that explore the meaning of that experience for individuals.
- iii) Third, the researcher collects data from individuals who have experienced the phenomenon under investigation. Usually the data are gathered through long interviews of 5 to 25 experienced persons.
- iv) Data are analysed in the form of statements and units. Then the units are transformed into clusters of meanings.
- v) Finally such analysis is linked with general description of experiences incorporating what was experienced and how it was I experienced.

### 11.5.3 Grounded Theory Study

#### Main Features

This kind of study aims at discovering or generating a theory. Here theory means an abstract analytical scheme of phenomenon. In other words, a theory is understood as a plausible relationship, as any concept or sets of concepts. In this case, theory is discovered in the context of a particular situation. This situation is one in which individuals interact, take actions, or engage in a process in response to a phenomenon. The researcher intends to explore how people act and react to a phenomenon. The process involved in data collection can be through continuous visits to the field, interviews with participants, in-depth observations of activities etc. The researcher develops and interrelates categories of information and writes theoretical propositions or hypotheses. Through the grounded theory method, a theory is generated in the context of a phenomenon being studied. Hence, the researcher goes beyond one step ahead of understanding the complexities of processes involved in a situation. He tries to help others to comprehend such complexities with the help of a theoretical framework developed by him. Such a theory is evolved towards the end of the study. As stated in Different Types of Studies in above, the theoretical framework may be depicted in the form of a narrative statement, Educational Research a series of hypotheses or propositions or in the form of a visual picture. For example, the researcher studies the complexities involved in the teaching learning processes where the sole responsibilities of such processes are shared by learners under the guidance of a teacher. The researcher would like to build a theory in this context through application of the grounded theory method, and shares different steps involved in this method

#### Steps of Grounded Theory Study

- i) In the context of a particular situation, the researcher makes several visits to the field, makes in-depth observations and conducts interviews. The data collection process continues till

the researcher comprehends the totality in a situation. However, the major focus remains on in-depth interviews.

- ii) With the help of data available, the researcher makes categories of information. A unit of information composed of events, happenings, and instances is presented in the form of one category.
- iii) Often it happens that during the data collection process the researcher analyses the data and gradually categories are formed on specific instances or happenings. Hence the data collection process and data analysis are integrated on many occasions.
- iv) After data analysis, the researcher may like to be back to the field to gather more information, analyse the data and so forth. Hence, this zigzag process continues till a saturated state is identified for arriving at a theoretical framework.
- v) Towards the end of the study, the researcher presents an elaborate theoretical framework to understand the complexity of an event.

#### **11.5.4 Ethnography**

Main features Ethnography can be understood as a description and interpretation of a cultural or social group or system. Here the focus of the study remains on examining the patterns of behaviour of a group, its customs and ways of life. This method involves prolonged observation of events where the researcher becomes a part and parcel of the day-to-day lives of the people. One to one interviews with the members of group corroborated with participant observation can form the base of such a method. The researcher makes use of ethnography to study the meanings of behaviour, language and interactions of the culture sharing group. For instance, educational processes of particular tribe or a rural village can be studied applying ethnography where cultural and behavioural interactions involved in education processes can be studied in totality. As an outcome, the researcher comes out with a report almost in a book form.

### 11.5.5 Case Study

Main features Case study as a method of research focuses on indepth study of a unit or case in totality. The case may be an individual, programme, an event, an institution, an activity, etc. The case study method was originally used in medicine to examine the patient's previous development, his health and physical state from the beginning and many other factors in the past, besides making a careful study of the patient's present condition and symptoms. Freud used the case study method to assist his subjects in solving their personality problems. The published detailed accounts of his interviews with patients and his interpretations of their thoughts, dreams and actions provide excellent examples of case studies. The investigation of a case is of exploratory nature. It involves detailed, in-depth data collection employing multiple sources of information concerning all pertinent aspects of a case. It is also interpreted that a case may be a unique and bounded system. This means the case under investigation is bound by time and place. Qualitative Research The uniqueness of a case refers to the typical characteristics of a case, such as a quality institution or an ashram school, a particular programme or course of study, or university or an innovation etc. The researcher tries to explore in detail about what events occur, and how they occur. Multiple sources include observations, interviews, audio-visual materials, documents, records, etc. A case is to be studied in a given context, i.e., study of a case is conducted within its setting. The setting may be a physical setting, social, historical or economic setting, etc. As a whole, the case study helps a researcher to understand the complexities of an event or events with contextuality and develops insight about with the nature and process dimensions of the events studied. In the Indian context, such a method has been used in the case of doctoral and research projects like: 1. "A study of distance education in an Indian university". This study's focus is an identification and description of different underlying factors. Contributing towards success of distance education programme of all university as a whole. 2. "A case study of management processes of National Adult Education Programme



in Orissa". This study aims at description of the management processes interlinked with national level, state level, district level, project level and village level adult education programmes in terms of the different dimensions of management like policy making, planning, communication, staffing, direction, co-ordination, budgeting and evaluation. Besides these kind of case studies we have come across some other studies like 'A case study of school products in Delhi', 'A study of school involvement in a village of Nagaland', 'Case studies of innovative institutions at secondary level in Tamil Nadu', 'Management of Medical Colleges -A case study'. From the above you can see that the case study method in educational research may focus either on the study of an educational programme or on a set of processes of an institution, in the context of development of learning skills, language competencies, reading comprehension of students. Guidance counsellors or social workers conduct case studies for diagnosing a particular condition or problem and recommending therapeutic measures. They gather data from a particular individual and confine their interest to the individual as a unique personality.

### **Steps of Case Study**

Case study is conducted by adopting the following steps:

- i) First, the researcher identifies the uniqueness of the case to be studied or a number of cases to be studied. Whether a normal case is to be studied or an unusual and typical case is to be studied is decided at this stage.
- ii) Keeping in view different contexts and perspectives of the problem, the researcher delimits what is to be studied within the scope of investigation.
- iii) Once the dimensions of a unit are identified for investigation, the researcher locates different mechanisms of gathering different varieties of evidences from different sources. After identification of the case and

## Notes

content, the present status of the case is determined through direct observation or record. Here the researcher goes far beyond casual observation or superficial description. For example, to make a case study of a delinquent child, the first thing the researcher has to do is to survey the present status of the child by making an assessment of his physique, cognitive and non-cognitive factors. The mechanisms to be chosen may be observation, interviews, document surveys, audiovisual recordings, projective and non-projective tests, etc. The observation can be participant in nature as employed in ethnography. It can be direct observations. Data may be gathered through various sources like gathering evidences from participants, functionaries and stake holders of a programme, institutional records and documents, observation of events taking place etc.

iv) Data gathered through multiple sources are subject to qualitative analysis. It can be a holistic analysis of the entire case or analysis of a specific aspect of a case.

v) Data analysis involves descriptive procedures. The themes or issues are identified. Interpretations are made about the case in its given contexts. It also involves narration of events chronologically and giving an account of events in totality. The detailed perspectives of significant events concerning a case are highlighted through such narrations.

vi) In case of study of number of cases the analysis is done in two phases. First is analysis of data concerning the themes of each case separately i.e., within-case analysis which follows the first is doing a thematic analysis across the case. This is called cross-case analysis.

vii) In the final step the researcher reports about the experiences and findings of the case. For example, in a study first the researcher gathers evidence through participant observations and interviews. He notices certain trends and puts them into different cluster themes. Then he tries to place these themes in different groups and makes more abstract categories at the last stage of inquiry. Hence with the help of an inductive

approach the researcher's initial level ideas and questions are refined in due course of investigation.

Different research methods have their uniqueness with regard to the focus of studies. As stated above, Biography is used to explore the life of an individual whereas Phenomenological Inquiry is used to understand the essence of experiences of persons about a phenomena. Grounded theory approach is adopted to develop a theory grounded in data from the field. Ethnography aims at describing and interpreting a cultural and social group. Case study is used to develop an in-depth analysis of a single case or multiple cases. However, all these methods adopt qualitative techniques of data collection such as interviews, observation, study of documents, relevant records, etc. and incorporate descriptive and narrative approaches of data analysis.

**Check Your Progress 2**

Notes: a) Space is given below for writing your answers.

b) Compare your answers with those given at the end of the unit.

1. Give one example of each of the following:

- i) Biography
- ii) Phenomenological study
- iii) Ethnography
- iv) Case study
- v) Grounded theory approach

.....  
.....  
.....  
.....

2. What are the common points of different kinds of qualitative studies?

.....  
.....

.....  
.....  
.....

---

## **11.6 CONTEXT OF QUALITATIVE METHODS**

---

All the kinds of qualitative methods have a common framework of ontological, epistemological, axiological, rhetorical and methodological perspectives. i) From the point of view of assumptions regarding nature of reality, the qualitative inquiry values subjective and multiple social realities. It presumes that reality cannot not have existence being separated from the perceiver 1 researcher. Rather, reality is constructed by individuals involved in research situations. Unlike assumptions of single reality as emphasized by scientific inquiry the concept of multiple realities is emphasized in the qualitative paradigm of the researcher. Through qualitative inquiry, multiple realities are narrated by representing diverse perspectives on the phenomenon. Hence diversified constructed realities are the focus of qualitative model of inquiry. ii) In qualitative methods the role of researcher is visualized as inseparable from the problem under study. The qualitative researcher interacts with the situation he/she studies. The researcher becomes a part and parcel of situations studied. He almost tries to find the meaning by being very close to the reality. Hence the gap between observer and observed is minimised. He shifts his status from objective and external observer to insider. iii) From an axiological point of view, qualitative research emphasizes on value loaded inquiry. Since the researcher is the sole instrument of investigation it is natural that his values and biases influence the process of inquiry. Hence, field studies conducted by qualitative researchers reflect their values and biases in presentation of data and interpretation. iv) The language of qualitative research is artistic and literary in nature. Unlike scientific reporting where an impersonal third person presentation is emphasized, the qualitative studies are reported in informal style using the personal voice. Qualitative terms are presented in narrative form rather than presenting a definition along with its explanations instead of using the terms like internal validity, external

validity, generalizability, the qualitative terms like credibility, transferability, dependability and conformability are used in reporting of the studies. v) The qualitative research paradigm emphasizes an inductive logical approach. The researcher studies a given situation. He identifies categories from data gathered through interaction with the situation rather than confirming predecided categories through gathering evidences.

---

## 11.7 COMMON STEPS OF CONDUCTING QUALITATIVE STUDIES

---

The steps in conducting qualitative studies are

- i) identifying problems and research questions
- ii) designing sources of data, sample and data gathering techniques
- iii) conducting field study or collecting data and
- iv) data analysis and reporting.

i) Identifying Problems and Research Questions At the first stage the researcher identifies a problem to be investigated. It emerges from a thorough analysis and review of literature and the experiences of the researcher and experts on several issues and problems. The major process of research, at this stage, is identified with building the central question of the study and subsequently linking it with a number of sub-questions. Qualitative Research While identifying a central question the researcher puts in open-ended and nondirectional queries. The queries focus on the "What" and "How" aspects of the phenomena under investigation. The exploratory questions highlighting the "process" aspects or "meaning" aspects are considered as the main theme of the investigation.

Following the central question, a small number of sub-questions are identified at the first stage of research. There can be a number sub-questions, keeping in view the scope and focus of the study. For instance, in the context of a phenomenological study of understanding professional ethics of teachers, the researcher may ask sub-questions like:

## Notes

1. What does professional ethics mean to an experienced teacher?
2. What do the teachers do to upgrade professional ethics?
3. Describe one person with high professional ethics. These kind of questions become the guiding force for the next steps of research

ii) Design of the study At this stage the researcher develops an open sketch of sampling, tools and techniques to be adopted, sources of data etc. Unlike a scientific investigation where the design is built prior to conduct of the study, the qualitative researcher develops an open and flexible approach. At this stage of inquiry, the following considerations are kept in mind by the researcher.

a) Sampling Purposive sampling techniques are used for identification of informants' responses of the study. For instance, in the case of a phenomenological study the researcher identifies all participants who experience the phenomena being studied. Moreover, the researcher examines the individuals who can contribute significantly to develop a theory in the context of a grounded theory. In the case of a case study, the researcher likes to gather evidence from diverse situations which display multiple facets of the case. Hence, he picks up a heterogeneous group of respondents who are stake holders of the system. As a whole, it can be said that keeping in view the nature of questions involved in the study, the researcher identifies the appropriate sample. The sample size may vary from problem to problem.

b) Forms of data The researcher must be clear about what kinds of data are to be gathered for investigation. Mainly, they can be categorized under the following four heads:

1. evidences about the surroundings, situations, activities, events and performances, etc. through observations;

2. evidences or experiences as explored through interviews;
3. evidences on developments; performances, demography, rate of participating, nature of progress recorded or noticed in various kinds of documents, publications, write ups etc; and
4. physical trace of evidence, behavioural dynamics, expression, feelings, emotional outbursts, etc. to be trapped by audio-visual data gathering devices. The investigator must keep an open mind while employing various data gathering devices keeping in view the nature of data that emerge during the data collection process /field work.

iii) Data collection process As stated earlier, mainly four kinds of data gathering devices are adopted by qualitative researchers viz., observation, interviews, document analysis and audio-visual gadgets. The researcher must be competent enough to deal with these different qualitative techniques of data collection.

a) Interviewing The researcher adopts informal interviews keeping in view the sub-questions related to the issues and themes involved in the study. Identifying sample respondents to be interviewed purposefully is the first step. Whether a researcher goes for one-on-one interview or focus group interview is to be kept in mind. For smoothening the process of the interview, the researcher may make use of an interview protocol / format with open- ended questions and ample space between the questions to write responses of the interviewers comments. The researcher may opt for using audio cassette recorders or a telephonic device for interviews.

b) Observing The qualitative method mainly makes use of participant observation technique where the researcher gets an opportunity to witness the events being a complete insider. While ethnography insists on this kind of observation, in the case of study the researcher may opt for direct observation. In employing the observation technique the researcher must be careful about identifying the site to be observed,

## Notes

identifying who or what to observe and how long, the observation protocol to be used, the role of the observer i.e., participant or outsider are to be thought about during the data collection stage. Moreover, the recording of observations i.e., portraits of the informant, physical setting, particular events and activities, researcher's reactions must be kept in mind while adopting observation techniques.

c) Documentation While adopting this technique the researcher must examine which materials are relevant to the study. He will also have to explore the means of how to trap evidences or having accessibility to relevant documents. As in case of historical studies, the researcher will have to be careful about authenticity of documents and the validity of information presentable therein.

d) Visual devices As a field worker the qualitative researcher must be acquainted with using various kinds of audio-visual recording practices. While using them he must be careful about the sensitivity of situations, ethical aspects, uniformity, etc. All these devices are used to gather evidences in natural setting with a view to studying the emerging total picture of the situation.

iv) Analysis of data and reporting Data collection and data analysis go on side by side in qualitative studies. Different techniques of data analysis include:

a) Review of data and representing data by case, by subject and by themes in the form of description, diagrams, tables and graphs.

b) Initial codification and categorization and further minimizing categories.

v) Frequent appearance of codes, development of categories, development of analytical frame work for theorization.



a) In this approach the researcher makes 'a thorough review of information as noted in the observation formats or interview schedules, notes, etc. He points out key points and writes reflective notes on the basis of an initial review. Sometimes it motivates the researcher to verify certain points by making further interaction with the informants. The researcher translates the ideas of informants into metaphors. Furthermore, the researcher presents descriptions in the form of tables, graphs, pictures, figures, diagrams, etc. He represents the case by themes or by specific subject areas.

b) The process of codification or categorization involves developing a short list of tentative codes that match a text segment. Elaborate lists of codes are arrived at in the initial stage. Then categorization of facts and themes takes place by identification of common codes. Furthermore, the categories are ruled to identify major themes and present them in narrative forms.

c) Another data analysis technique indicates preliminary counts of data and determines how frequently each code appears in the total data. Then the researchers identifies categories and develop analytical frameworks with a view to generating theories. In the whole data analysis procedure, a spiral approach is followed to identify categories, reflecting, cross-questions, reading, interpreting etc.

**Check Your Progress 3**

Notes: a) Space is given below for writing your answers.

b) Compare your answers with those given at the end of the unit.

1. How does a researcher identify research questions?

.....  
.....  
.....  
.....

.....  
.....

2. What are the different sources of data used in qualitative studies?

.....  
.....  
.....  
.....  
.....

3. How is participant observation different from general observation?

.....  
.....  
.....  
.....  
.....

4. What are the different data analysis techniques used in qualitative studies?

.....  
.....  
.....  
.....  
.....

---

## **11.8 VERIFICATION OF TRUSTWORTHINESS OF QUALITATIVE RESEARCH**

---

Unlike standard procedures followed in establishing the objectivity and validity of quantitative research, you will come across many qualitative approaches to verify the trustworthiness of qualitative research. Generally, the following verification procedures are employed for studying trust worthiness of such kind of studies :

i) Prolonged engagement and persistent observation: Researcher's prolonged interaction with informants and participatory observation wins

the confidence of the informants. Hence, this procedure minimizes the possibility of vagueness of data, prolonged engagement in the field helps the researcher to reflect on the genuineness and authenticity of data, on gathering relevant data and on making purposeful use of evidences.

ii) Rural explanations: Once the researcher (after qualitative analysis) has described the patterns and their explanations, it is important to look for rival or competing themes and explanations both inductively and logically. Inductively, it implies looking for other ways of organizing the data that might lead to different results. Logically, it involves searching for other logical possibilities and then finding if those possibilities can be supported by the data. However, while considering rival hypotheses and competing explanations, the strategy to be employed by the researcher is not attempting to disprove alternatives, but to look for data that support alternative explanations. In this strategy, the researcher should give due weightage to supporting evidence look for the best 'fit' between data and analysis.

iii) Triangulation: It is another significant technique of verification where the researcher makes use of multiple and different sources, methods and theories in order to provide corroborating evidence. In other words, through triangulation, evidences are corroborated from different sources to throw light on a theme or perspective. It involves comparing and cross checking consistency of data derived by different means at different times using qualitative methods. It means

- (i) comparing observational data with interview data;
- (ii) comparing observational data with questionnaire data;
- (iii) what participants of a programme say in public with what they say in private;
- (iv) checking for consistency the opinion of the participants about a programme over a period of time and
- (v) comparing the opinion of the participants of a programme with others who are associated with programme in one capacity or the other. The triangulation of data sources within

qualitative methods will seldom lead to a singly totally consistent picture. But such triangulation is helpful to study and understand when and why there are differences. –

- (vi) Peer view or debriefing: The methodology followed in the study with emphasis on subjective approach gets cross-examined by the peer researchers. The peers review the study, present their views in debriefing sessions. Such an operation helps the researcher introspect on the study in the context of peers reflection on the procedural dimensions and reporting.
- (vii) Design checks: The nature of research design and methodology also contribute to distortion in results. Sampling gives rise to three type of errors. The errors may be due to :
  - a. distortion in situations that were sampled for observation;
  - b. distortion introduced by the time periods during which observations took place;
  - c. distortion because of selectivity in the people who were sampled either for observation or interviews. Thus the researcher must be careful to limit results of
- (viii) Negative case analysis: Unlike quantitative study where the hypothesis is set prior to conduct of the study the, qualitative researcher keeps on refitling the working hypotheses in due course of conducting the study. Hence context specificity is kept in mind for refining the initial hypotheses till the completion of the data collection process and its analysis. The search for negative cases and instances that do not fit within the identified pattern and their understanding is also competent in the verification and validation of results.
- (ix) Clarifying researcher's bias: It is presumed that the researcher himself recognizes his bias and subjectivity. Hence in view of making it transparent the researcher comments on past experiences, prejudices, and orientations that influenced the interpretation and approach to study.
- (x) Member checks: Through this technique the entire studies repart with its prime data base, analysis and interpretations i. presented before the participant respondents who can judge

the accuracy and credibility of the account. They are asked to provide critical observations on the research work.

- (xi) Rich and thick descriptive presentation: The researcher describes in detail the participants or setting under study. Hence, the readers get a clear picture of the setting. They get chance to examine the findings that can be applicable or transferred to other similar settings. Different mes of Studies in Educational Research
- (xii) External audits: Through this process the expert consultant who is purely external and independent to the study examines both the process and the product of the account, thus assessing its accuracy. The above indicators are employed to examine the authenticity and credibility of qualitative studies through involvement of human experiences in the process. What is more significant is that the researcher must maintain transparency throughout the study so that the scope of verifying its trustworthiness becomes wider and accessible.

---

## 11.9 LET US SUM UP

---

In the foregoing sections, you were shown that main features of qualitative studies are needed mainly with the purpose of understanding process dynamics of a phenomenon in contextual frameworks. Such understanding may lead towards theorizations. Different kinds of qualitative methods are used in different context, such as biographies, ethnography, grounded theory, phenomenology and case study. However, these methods follow common approaches like interviews, observation, and study of documents and narrative techniques of data analysis. These methods adopt subjective human and participant experience-based techniques in contrast to that of generic methods of inquiry. However, as a researcher, you will have to be mature enough to adopt various checks and balances for trustworthiness of findings of such studies.

---

## 11.10 KEY WORDS

---

**Case Study:** In the social sciences and life sciences, a case study is a research method involving an up-close, in-depth, and detailed examination of a subject of study, as well as its related contextual conditions. Case studies can be produced by following a formal research method.

**Qualitative Research:** Qualitative research is a scientific method of observation to gather non-numerical data. This type of research "refers to the meanings, concepts, definitions, characteristics, metaphors, symbols, and description of things" and not to their "counts or measures".

---

### 11.11 QUESTIONS FOR REVIEW

---

1. Identify and state different problems / topics in the area of education belonging to each of the categories of qualitative researcher:
  - a) Biography
  - b) Ethnography
  - c) Phenomenology
  - d) Grounded theory
  - e) Case study
2. Chalk out the strategy for conducting any of these studies.
3. Differentiate between the strategies adopted in conducting qualitative studies and those in quantitative studies.
4. What is the relevance of qualitative research in education?
5. What are the limitations of qualitative research?
6. How can we judge trustworthiness of qualitative research?
7. What should be the role of a researcher in conducting qualitative research?
8. Can we integrate qualitative studies with quantitative studies? Discuss.

---

### 11.12 SUGGESTED READINGS AND REFERENCES

---

- Best John W. (1992): Research in Education. SLh Indian Reprint, New Delhi: Prentice-Hall India.

- Borg, W. R. and Gall, M.D. (1986): Educational Research: An Introduction.
- Fourth Edition, New York: Longman.
- Creswell, John W. (1998): Qualitative Inquiry and Research Design. New Delhi: Sage.

---

## 11.13 ANSWERS TO CHECK YOUR PROGRESS

---

### Check Your Progress 1

1. Quantitative studies are mainly of explanatory nature focusing on the 'why' components. Qualitative studies focus on 'how' and 'what' components with a view to understanding in process dynamics in specific situations.
2. Natural setting, researchers involvement, visual and verbal images, process dynamics, inductive logic and participant perspectives.
3. To study complex human and social problems in totality, to observe the phenomena without making any manipulations or controls on variable, to study and explore the nature of subjective quality.

### Check Your Progress 2

1. (i) The contribution of an innovative teacher. (ii) Existing curriculum transaction in a school. (iii) Parents involvement in school management. (iv) Study of an ashram school. (v) Professionalism in teaching.
2. Using qualitative data collection techniques like interviews, observation, document analysis, photographic techniques, and descriptive and narrative approaches of data analysis.
1. Experiencing first-hand exposures, study of literature, identify central question and link sub-questions to the theme. 7. Empirical evidences as

## Notes

revealed through interviews and observations, demography, documents, physical trace of evidences, visual and verbal images.

2. A participant observer plays the role of being part and parcel of locality. General observer is alien to the reality.

3. Reviews, descriptions, codification, categorizations.



---

# **UNIT 12: RESEARCH PROCESS: REVIEW OF LITERATURE, IDENTIFYING RESEARCH PROBLEMS**

---

## **STRUCTURE**

- 12.0 Objectives
- 12.1 Introduction
- 12.2 Need and Purpose
- 12.3 Functions
- 12.4 Types
- 12.5 Structure
  - 12.5.1 Title
  - 12.5.2 Introduction
  - 12.5.3 The Problem
  - 12.5.4 Scope
  - 12.5.5 Objectives
  - 12.5.6 Hypotheses
  - 12.5.7 Review of Related Literature
  - 12.5.8 Research Design
  - 12.5.9 Tentative Chapterisation
  - 12.5.10 Limitations
  - 12.5.11 Operational Definitions
- 12.6 Funding
- 12.7 Monitoring
- 12.8 Ethics Managing Resources
- 12.9 Let us sum up
- 12.10 Key Words
- 12.11 Questions for Review
- 12.12 Suggested readings and references
- 12.13 Answers to Check Your Progress

---

## **12.0 OBJECTIVES**

---

After reading this Unit, you will be able to:

know what is a research plan;  
appreciate its need and purpose;  
understand its functions;  
know its different types and structure;  
comprehend the tools and techniques used in a plan; and  
prepare a research plan.

---

## 12.1 INTRODUCTION

---

We have discussed in Unit 12 on Research Design, the importance of organisation and planning in research. Research design involves formulating a strategy for your research. There are a number of factors involved while taking a decision in this regard. Once you have formulated a research design, the next step is to prepare a research plan. As a Research Plan with a research design, research plan also gives your research a sound footing. This Unit is devoted to a discussion of the different aspects of research plan in detail.

Research is a systematic effort towards unraveling the mysteries surrounding us. It involves putting in persistent efforts to know the unknown. In doing so, we start from where we are, what is our present state of knowledge? We tend to interpret this present state by studying the work done by others in the field. We have used the word systematic to characterise research. There are reasons for that; there is a specific purpose while carrying out research; the effort required is monumental in view of the fact that a comprehensive study of what has already been known about it has to be understood. It has to be continued further, from conception to observation to analysis to interpretation and finally reporting, which involves considerable intellectual efforts. Robson has characterised research as systematic, skeptic, and ethic. By systematic, he means that we should be clear about what we are doing, why we are doing, and how we are doing. By skepticism he conveys that the researcher should check, cross-check, and verify his/her views before finalizing them. Researcher should not violate ethics while conducting research, whether it is conceptualisation of topic, data collection, analysis or presentation of results. Research is of two types, viz., pure and

applied. The trend today is towards applied research. That is not to belittle pure research, as both are complimentary to each other. Pure research is done to improve upon our existing state of knowledge. Applied research is done to find out new knowledge to be put into Research Design application. It is carried out with an aim of developing something better for the benefit of the society. We have seen that research needs to be organised and systematic. Research Design is a step towards carrying out research in a planned way. In this Unit, we shall study about what is research design and what is its purpose. We shall also discuss the factors affecting research design and the types of research design.

Research is an important activity affecting the society as a whole therefore, it involves a lot of decision making. Research design also involves a lot of decision - making. It provides a structure and shape to your research project. After finalising your topic, you decide about how you are going to conduct your study. It involves formulation of strategy for all the stages starting from formulation of hypotheses to the analysis of data. Kerlinger defines research design as a plan, structure, and strategy of investigation so conceived as to obtain answers to research questions or problems. The plan is the complete scheme or programme of research. It includes an outline of what the investigator will do from writing the hypotheses and their operational implications to the final analysis of data. Thyer has defined research design as a blueprint or detailed plan for a research study - starting from operationalising variables so that they can be measured, to selecting a sample of interest to study, collecting data to be used as a basis for testing hypotheses, and finally analysing the results. Thus, we can conclude that research design provides us a base on which we conduct our research.

---

## **12.2 NEED AND PURPOSE**

---

### Definition

Research plan has been used synonymously with proposal and synopsis. They refer to a blue print of your research. Both are used for documents that describe in detail:

## Notes

- What are you going to do?
- Why are you going to do?
- How are you going to do?
- In what resources (time, money, infrastructure, etc.) are you going to do?
- What are you not going to do?

Research is a systematic endeavour towards quest for new knowledge. The word systematic is important and needs clarification. Studying informally and casually may also yield new knowledge but that is not research. In research, organised efforts are put right from thinking of a topic to the presentation of results. The magnitude of the study and the efforts involved require systematisation. Research also involves spending public money therefore, it requires proper planning for effectiveness and efficiency. The methods and techniques involved in data collection and analysis may involve subjectivity. Adequate planning needs to be done to minimise this subjectivity. Validity is another important issue for which a proper design is required. Validity ensures that what we are measuring is what we intend to measure.

A research proposal is thus, a document to your plans and ideas of carrying out your research. You may be wondering, how is it different from research design? Research design is also a strategy of how you are you going to conduct your research. But it precedes research plan. Research design is a plan of the technical decisions regarding the what, why, and how of your research. Here we take decisions regarding the nature of investigation, data collection methods to be adopted, and number of contacts to be made with the subjects, and the period of reference with the subjects of study. Research plan can be formulated only when research design has been decided. In fact, research design helps to formulate research plan. It is the document that describes all the decisions that have been taken in the design stage plus the administrative decisions concerned with your research. It presents systematically

everything starting with the title of your project to the tentative structure of your thesis.

As stated earlier, we repeat that research is a systematic endeavour towards quest for new knowledge. There is a role of intuition also in research but a systematic step-bystep approach needs to be followed in research. Studying informally and casually may also yield new knowledge but that is not research. In research, organised efforts need to be put right from thinking of a topic to the presentation of results. Research is a long journey from knowing and understanding the present state of knowledge in a particular field to exploring raw areas where additions and improvements can be done and finally making those improvements. It is a project, which involves investment of considerable resources. And thus, needs planning as a proposal. The purpose of a research plan is to:

- present for him the proposed plan of action;
- present for the supervisor and other authorities also the plan of action for their approval;

The purpose of a research design is to provide information regarding:

- What is the study?
- Why is the study being carried out?
- Where will the study be carried out?
- How will the study be carried out?
- What will be the processes and tasks involved?
- What will be the data?
- How will the data be collected?
- What methods of sampling will be used?
- How will the analysis be done?

---

## **12.3 FUNCTIONS**

---

The functions of a research plan are to:

give directions on what needs to be done, when and how and in what order;  
provide a route from stating the topic to finalising the results;  
enable to evaluate your progress during research;  
define your topic to limit its scope; and  
prove to your supervisor that you have gone into the fine details of your topic and will be able to conduct research.

**Check your progress 1**

Note: i) Write your answer in the space given below. ii) Check your answer with the answers given at the end of the Unit

1) Define research plan. How is it different from a research design?

.....  
.....  
.....  
.....  
.....

2) Describe what would you discuss in the introduction to the plan of your Study.

.....  
.....  
.....  
.....  
.....

---

**12.4 TYPES**

---

There are two types of research plans. These are: quantitative and qualitative proposals. Quantitative proposal is given for experimental and descriptive research whereas qualitative proposal is given for descriptive and exploratory research. Though, there is no hard line of demarcation, generally qualitative research doesn't have hypothesis and operationalisation of concepts. Operationalisation of concepts refers to

giving operational definitions, finding out the independent and dependent variables, sample selection, and finalising the measuring instruments and their reliability and validity. Instead the qualitative proposal will have research procedures which quantitative research doesn't have.

The nature of investigation can be exploratory, descriptive, causal/experimental, semi or quasi- experimental, non-experimental, and field research. Exploration is an important characteristic of research. Any research begins with it when the researcher dives into the unknown and unsolved terrains. He/She starts with a quest for knowing more through exploration. It is an initial foray into the densities of the unknown. Exploration starts with a vague idea of what is intended to be researched. It forms the basis of research. It is not very systematic to the order of research to be undertaken otherwise. It is a flexible approach to undertaking research where the sampling is generally non-probability and the data collection methods are unstructured. It involves a study and analysis of the literature and discussions with peers and fellow colleagues to know their views on the topic. Descriptive research is carried out to provide information about a person, thing or process. It describes the characteristics of an individual, group, organisation, or phenomena, conditions, or a situation. The characteristics are described in terms of the dependent variables. Description may be limited to events of past or present but not of the future. In that case it becomes experimental research.

In descriptive studies data collection is done through structured methods. Samples are selected by random sampling. The nature of investigation moves systematically from exploratory towards experimental. The degree of investigation goes on increasing as we move ahead. Casual investigation in exploration, to description and finally causal investigation in experimental research. It aims to find cause and effect relations between the dependent and independent variables. Experimental research studies the effect of independent variables on dependent variables. The researcher identifies the two different kinds of variables and the relationship between them. For this, he/she reviews the literature on the subject and also related studies done by others. Discussions with

## Notes

peers and other professionals also help in finding out the relationship. Hypotheses are framed for verifying the relationships. The research is conducted under controlled conditions so that the changes in the dependent variables can be attributed solely to the changes in the independent variables. Semi - experimental studies are different from experimental studies in that the sampling in experimental studies is random sampling compared to non - random sampling in semi - experimental or quasi-experimental studies. Non - experimental studies also find out causal relations but they follow the reverse approach. Experimental studies explain the cause - effect relation by identifying the independent variables and later inducing changes in them to find out the resultant effect on the dependent variables. Non- experimental studies ascribe the changes that have Research Process taken place in the dependent variables to some independent variables. They do not induce changes in the independent variables to see the effect on dependent variables. This is generally done in the social sciences and the reason for doing so is the population that are human beings compared to physical and chemical entities in sciences. Let us consider an example to clarify the difference. We want to see the effect of use of IT in the classroom on the performance of students. In experimental studies, we would take use of IT in the classroom as the independent variable and the results of students as dependent variable. We would compare the scores of students after introduction of IT to the scores obtained by them earlier and find out the relation. In non-experimental study, we would check the scores of students after IT has been introduced and find out the relation between them by studying the coefficient of correlation. Let us take another example to understand the difference, where we are studying the effect of OPAC on the use of catalogue. We would divide the users randomly into two groups. One group of users is provided the facilities of a traditional catalogue for access to the literature. The other group is provided the facility of an OPAC to access the literature. We would measure the use of catalogue in the two cases and ascribe the difference to OPAC. Field research is done in the natural surroundings in real life situations. Here the main criterion is doing research in social settings rather than on the techniques of research. Let us discuss some observations on field



research: “Field research is the design, planning and management of scientific investigations in real-life settings” (Fielder) Kaplan comments that Field research involves direct or indirect observation of behaviour in the circumstances in which it occurs without any significant intervention on the part of the observer. We can conclude that field research is conducted in real life settings without any modifications done to the settings. There is little stress that the techniques applied are scientific. Importance is given to the fact that the observer collects data while being on the site along with those observed. He is trained to be part of the observed group and objective in recording the observations. Such research is carried out particularly in subjects like sociology or social work. Field studies have been divided into field research and field experiments. Field experiments are different from field research in that the former involve studying the effect of varying independent variables on dependent variables in real life natural settings. The difference between experimental research and field experiment is that the former are conducted in laboratory settings whereas the latter are conducted in natural settings. Thus, the control in the observations is not possible in field experiments, which is possible in laboratory experiments.

---

## **12.5 STRUCTURE**

---

A research proposal is presented in the following structure:

### **12.5.1 Title**

The title of your study (dissertation or project) is the first part of your plan. We should ensure that the title is self- explanatory. It should convey what we intend to do. There should be no ambiguity. It should be clear, precise, and grammatically correct. It should not be broad or more specific than what we plan to study in the research. If we wish to Research Plan study the impact of personality development programmes on the tackling of users by staff, then the topic could be: Handling of users by staff in libraries: Impact of personality development programmes

## 12.5.2 Introduction

The introduction provides background information to the topic of your study. It includes a thorough review of what is available related to your area of study. Try to clarify the conceptual area zeroing down towards the topic. Ranjit Kumar enumerates the following list of the aspects to be covered in the introduction:

An overview of the main area of study;

A historical perspective (development, growth, etc.) pertinent to the area of study;

Philosophical or ideological issues related to the topic;

Trends in terms of prevalence, if appropriate;

Major theories, if any;

The main issues, problems and advances in the subject area under study;

Important theoretical and practical issues relating to the central problem under study; and

The main findings related to the core issue(s). Let us clarify with an example. If the topic of study is, “Impact of automation on the services of academic libraries in India”. The introduction should discuss:

Brief discussion on the concept and history of automation;

Attempts at automation of libraries in India;

Trends of automation of academic libraries in India; and

Impact of automation on libraries and their services.

### 12.5.3 The Problem

Introduction and the problem could be visualised as occurring together in continuation as background to the study. The difference is that introduction is more general as compared to the problem. The problem starts with where introduction has left the topic. It is continuation of introduction, focusing specifically on the topic. Why have you chosen the topic? Is there any need to conduct such studies? What is the rationale of your topic? What gaps in the existing knowledge does it intend to fill? Ranjit Kumar enumerates the following issues that need to be discussed in the problem: Identify the issues that are the basis of your study; Specify the various aspects of/perspectives on these issues; Identify the main gaps in the existing body of knowledge; Raise some of the main research questions that you want to answer through your study; Identify what knowledge is available concerning your questions, specifying difference of opinion in the literature regarding these questions if differences exist; and Develop a rationale for your study with particular reference to how your study will fill the identified gaps.

Let us discuss what needs to be discussed in the problem if our topic of research is: “Impact of automation on the services of academic libraries in India”. It should include the following based on a review of literature: categorise the different kinds of services and the impact of automation on them; discuss the views and theories propounded by experts in this regard; present case studies of attempts by professionals towards studying impact of automation on services; separate the impact on service providers and users; bring to light issues that remain to be studied or point towards a need for more investigation; and state the relevance of your study, why is it needed, how is it going to fill the gaps, if anything exists. If such studies have already been conducted, you could justify the need for your study if some aspects have not been studied, or some issues have cropped up from earlier studies, or your study is being conducted in a new environment or conditions, or it is a longitudinal study. You would recall, we have discussed about

longitudinal studies in Unit 15 on Research Design. These are studies conducted again on the same population after a gap of time to know the change in dependent variables over a period of time. It helps to provide a pattern of change in the dependent variable

### **12.5.4 Scope**

After the problem has been stated, it is important to explain its scope also. In the scope one should indicate to what extent one intends to probe the topic. He should clarify the scope as far as the subject content is concerned as well as the geographical area is concerned. There should also be a submission of the scope regarding the coverage of the time period. The scope must give an indication of the limitations of the topic.

### **12.5.5 Objectives**

After the problem, the plan should state the objectives of your study. This is one of the most important parts of your plan. It helps to know what you intend to do. Anyone interested in your study gets to know the whole picture from just the objectives. You can very well judge the importance of objectives. Therefore, it is important that they are stated in a crisp language so that are clear and unambiguous. Moreover, they should convey what is the intended outcome of your study. Where do you reach to when you start from here? For that you need to use action verbs like: to know, to find out, to Research Plan evaluate, to automate, to design, etc. Another thing to be borne in mind while stating objectives is that they should not be broad. One should clearly make out the intentions of one's study that could help him/her or even you while evaluating. It is a pointer to measure how far you have succeeded in your project.

Let us discuss the same example again, "Impact of automation on the services of academic libraries in India". The objectives of the study could be stated as to study the impact of automation on:

Change in the form of services;

Any new services introduced;

The service providers;

The frequency of use of services; and

Ease of use of services.

After the problem, the plan should state the objectives of your study. This is one of the most important parts of your plan. It helps to know what you intend to do. Anyone interested in your study gets to know the whole picture from just the objectives. You can very well judge the importance of objectives. Therefore, it is important that they are stated in a crisp language so that are clear and unambiguous. Moreover, they should convey what is the intended outcome of your study. Where do you reach to when you 257 start from here? For that you need to use action verbs like: to know, to find out, to Research Plan evaluate, to automate, to design, etc. Another thing to be borne in mind while stating objectives is that they should not be broad. One should clearly make out the intentions of one's study that could help him/her or even you while evaluating. It is a pointer to measure how far you have succeeded in your project. Let us discuss the same example again, "Impact of automation on the services of academic libraries in India". The objectives of the study could be stated as to study the impact of automation on:

- Change in the form of services;
- Any new services introduced;
- The service providers;
- The frequency of use of services; and

- Ease of use of services

**Check your progress 2**

Note: i) Write your answer in the space given below.

ii) Check your answer with the answers given at the end of the Unit

1. Taking a hypothetical example, state the objectives of your study.

.....  
.....  
.....  
.....  
.....

2. Define hypothesis. Differentiate between null and alternative hypothesis.

.....  
.....  
.....  
.....  
.....

**12.5.6 Hypotheses**

The next section of your plan should be devoted to hypotheses. Hypothesis is an assumption, presumption, or in simple words guess towards some situation or condition. It is an assumption of relationships between the dependent and independent variables. Hypotheses provide the researcher a line of action along which he/she moves to find out answers to the problem he/she is working on. It is not necessary that every research study has a hypothesis. Though we have stated that hypotheses formulation involves guess that does not mean there is no base to it. It involves a thorough review of literature to understand the concept and the relationships existing between the various variables existing. It is on the basis of this review and discussion with peers and others in the field that we propose a relationship between the variables.

This is tested during the course of study. It is pertinent to recall what we have studied in the earlier units, that hypotheses if proved true later results into theories and finally into laws. There are two kinds of hypotheses, null and alternative. Hypothesis is stated as null hypothesis, which is a negative statement. We state the relations expected between the variables in a negative way, e.g., if we expect that automation has resulted in increased use of library services. We would frame the null hypothesis as:  $H_0$  = There is no increase in use of library services due to their automation. The alternative hypothesis would then be:  $H_1$  = There is increase in use of library services due to library automation. Thus, if null hypothesis is proved to be false, the alternative hypothesis is proved to be true.

### **12.5.7 Review of Related Literature**

After the hypotheses are stated, a brief review of literature is presented. It helps the researcher to know and assimilate what others have already done in the field. It gives him the direction of movement into his research. He comes to know the ripe areas for research. Whenever we plan to research, we have to start from somewhere, which we come to know by literature review. Literature review is conducted at two different stages of research work. This is the first stage when it is a brief review of related literature. It is done again after the proposal is approved and the study is conducted. Later the scope of review of literature is greater covering all aspects and a period of coverage is decided which depends on the subject and topic of study.

### **12.5.8 Research Design**

The plan of research from the point of operationalisation of hypotheses to the analysis of data is presented as research design. The research design is the blueprint of your nature of investigation, data collection methods to be used, number of contacts to be made with the subjects, and the reference period of study. Broadly your nature of investigation can be either exploratory, or descriptive, or experimental. It depends upon your

## Notes

topic which one you choose. A study can adopt more than one also. For example, any research starts with an exploratory investigation where you tend to explore your topic. It is a stage where you tend to formulate your topic based on the review of literature and discussion with others in the field. The research here does not rigidly follow research methodology. The data collection methods also do not follow strictly probability sampling. As in social sciences, in library science also descriptive methods of research are followed. It describes the situation of a n object, phenomenon, or a process or an event in the past or present. If it does in the future it becomes experimental research. If our topic of research is: “Impact of automation on the services of academic libraries in India” it has to be descriptive research. Similarly you have to describe the data collection design, whether it will be survey, or case study, or content analysis. It will be survey in case of the topic, “Impact of automation on the services of academic libraries in India”. A mention of the number of contacts to be made with the subjects and period of reference also needs to be made here. In the research design you should also mention the population under study clearly. Describe the data that you intend to collect in terms of the dependent and independent variables. Then you also need to clarify whether you will collect data from the whole of the population or from only a sample. If it has to be a sample study, you need to specify what methods of sampling would be used. Whether it would be probability sampling or non- probability sampling. Within probability and non- probability sampling 259 also state which type of sampling method would be used. The size of the sample should Research Plan be mentioned here. Next in the research design, the researcher should mention the techniques and tools of data collection. In the techniques, he/she should mention, whether it would be observation, questioning, or interview. What tools of data collection would be used, questionnaire (mailed or self- administered), interview schedule, etc. The research design should also state when and where the data collection will be done. After the data collection methods are discussed, the methods of its presentation and analysis should be described. Here a mention should be made of the types of tables and graphs that you intend to use to present your data. Also explain how you would analyse your data. What



statistics do you intend to use? If you plan to do the analysis using some software package like SPSS, MS-Access, or MS-Excel, etc. mention in the design.

### **12.5.9 Tentative Chapterisation**

The physical structure of the research report is also presented in the research plan as tentative chapterisation. What will be the chapters in your thesis are presented here. Normally, the final report also contains the chapters: Introduction, which comprises the introduction, the problem, scope, objectives, hypothesis, limitations, and operational definitions. The other chapters are review of literature, conceptual structure, research design, data collection, analysis of data, discussion, and conclusion.

### **12.5.10 Limitations**

Every research study has limitations. These could be from the point of view of the contents (coverage of the subject), geographical area, time period of study, etc. The researcher should very earnestly admit the limitations in his/her study. This is no drawback for the study. The limitations vary according to the level of study and that of the researcher. If the limitations are presented, the evaluation of the thesis is done keeping in view of these.

### **12.5.11 Operational Definitions**

The research like any other work has to deal with vocabulary. We know the nature of any language is such that there is a lot of flexibility. The occurrence of synonyms, homonyms, and other such concepts may result in confusion. Thus, in the beginning itself, terms that are to be used are standardised in the form of operational definitions given in the research plan

---

## 12.6 FUNDING

---

Research involves financial investment. The researcher can get support for this investment from different agencies. In India, these agencies are University Grants Commission, Council of Scientific and Industrial Research, Indian Council of Agricultural Research, Indian Council of Medical Research, Indian Council of Social Science Research, etc. Researchers in library and information science can get their research proposals funded from any of these agencies depending upon their area of study. Why should an agency fund for your research? It would fund:

to encourage you to do research;

research adds to the intellectual wealth of a nation; and

your research topic lies within the scope of the institution and would be beneficial for it.

The research plan in case of a research for funding should also submit a financial estimate. The financial estimate should be presented under different heads, i.e., the purpose for which you are asking for grants. These may be:

space/rooms;

equipment;

books/ journals;

stationary;

travel; and

publishing.

The agency may provide grants for some or even all of these expenses. As far as publishing is concerned this grant is provided after the completion of the project/ theses.

---

## **12.7 MONITORING**

---

Monitoring is an important component of research like any other project. Research plan should also indicate your plans to monitor your research work. You need to apply time management and prepare a work plan to be submitted along with your research plan. The work plan should indicate when (month and year) would you: start your research; review related literature; prepare data collection instruments; collect data; do coding of data; do statistical analysis of data; prepare draft of your report; and finally submit report. It is not possible to follow the time schedule exactly as specified. It is a plan therefore there may be some fluctuations here and there. It is important while planning to make adjustments for any possible disturbances in the time schedule. Monitoring should be done regularly to avoid any delays in fulfilling the schedule. The researcher should be regular in his/her work and subdivide his plan into short time plans. He/She should translate the plan into tasks to be undertaken everyday. And should evaluate his/her days work to focus the cause of delays and overcome the reasons of delay. In case of academic research for the award of degrees the progress report is to be submitted quarterly or half- yearly approved by the supervisor to be submitted to the authorities. Similarly the progress report has also to be submitted to the funding authorities to prove accountability.

---

## **12.8 ETHICS MANAGING RESOURCES**

---

Researcher should follow a code of ethics in research. We need not overemphasise that it is important in research as in any other activity. Following ethics in work implies doing the work in the right way. Right way implies that it is done in a way acceptable to the society. Ethics has been used interchangeably with morals. Ethics is concerned with what one ought to do. To analyse and make it convenient to understand Chris Hart divides the ethical issues in research as those concerned with:

Research; Research Plan Researcher; Subjects; and Sponsoring Body/ University. While doing the right things the researcher should understand that the above four are the parties affected which should not suffer through your decisions while doing research. Let us discuss by way of examples the areas where ethics have to be taken into consideration: State the topic in an unambiguous way. It should not be the topic of some other study. The research should add to the repertoire of knowledge. Acknowledge the works of others that have been used in the research. Plagiarism should not be resorted to. Ask for financial grants only where required for research. In data collection the anonymity of the subjects should be maintained. The privacy of subjects should also be maintained. In case the subjects are animals or plants as in the case of sciences their treatment should be ethical. Tampering of data to achieve results should not be done.

**Check your progress 3**

Note: i) Write your answer in the space given below.  
ii) Check your answer with the answers given at the end of the Unit.

- 1) Discuss the areas where ethical considerations are important in a research plan.

.....  
.....  
.....  
.....  
.....

---

**12.9 LET US SUM UP**

---

In this Unit we have discussed the various aspects of research plan. We introduced the concept and discussed the need for, purpose and functions of research plan. The Unit provides comprehensive details on the structure of a plan. A step-by-step approach has been followed to explain

all the parts of a plan. Adequate examples have been given at all the stages to make all the steps clear. Research can be funded through various agencies, which has also been described in the Unit. There is a discussion of the heads under which funds can be applied. Researcher should take care that he/she follows ethics while conducting research. What are these areas has been stated for the help of the researcher. It has been discussed in this Unit so that the researcher gets conscious of ethics since the start of research. The Units ends with some model plans given as examples.

---

## **12.10 KEY WORDS**

---

**Hypothesis :** Assumption regarding the relations between the variables in a study.

**Problem :** A description of the problem you shall study in your research work.

**Research Design :** The strategy that a researcher adopts to undertake his research. It concerns the operationalisation of hypothesis, data collection, and data analysis.

**Research Plan :** The blue print of your research that states everything from the title to the limitations of your study. It also includes the financial estimates of your project and the work plan.

**Work Plan :** The plan of your research in terms of the tasks to be done along with the time at which they will be done.

---

## **12.11 QUESTIONS FOR REVIEW**

---

- 1) Define research plan. How is it different from a research design?
- 2) Describe what would you discuss in the introduction to the plan of your Study
- 3) Discuss the areas where ethical considerations are important in a research plan.

---

## 12.12 SUGGESTED READINGS AND REFERENCES

---

- Ackoff, Russell. (1953). *The Design of Social Research*. Chicago: University of Chicago.
- Adams, Gerald R. and Schvaneveldt, Jay. D. (1985). *Understanding Research Methods*. New York: Longman.
- Hart, Chris (2005). *Doing your Masters Dissertation*. Delhi: Vistaar.
- Kothari, C.R. (2004). *Research Methodology: Methods and Techniques*. 2nd ed. New Age: Delhi.
- Krishan Kumar (1998). *Research Methods in Library Science*. 2nd ed. Delhi: HarAnand
- Neuman, W. Lawrence (1997). *Social Research Methods*. 3rd ed. Boston: Allyn and Bacon.
- Ranjit Kumar (1999). *Research Methodology: A Step- By- Step Guide for Beginners*. Delhi: Sage.

---

## 12.13 ANSWERS TO CHECK YOUR PROGRESS

---

### Check your progress 1

1) A research plan is a document to your plans and ideas of carrying out your research. It is the document that describes all the decisions that have been taken in the design stage plus the administrative decisions concerned with your research. It presents systematically everything starting with the title of your project to the tentative Research Process structure of your thesis. Research design is also a strategy of how you are you going to conduct your research. But it precedes research plan. Research design is a plan of the technical decisions regarding the what, why, and how of your research. In research plan, we take decisions regarding the nature of investigation, data collection methods to be adopted, and number of contacts to be made with the subjects, and the period of reference with the subjects of study. Research plan can be

formulated only when research design has been decided. In fact, research design helps to formulate research plan.

2) The introduction to the plan of study provides a contextual background to the study. It discusses the basic theoretical and philosophical issues related to your topic. It also discusses the latest trends in the subject and area of discussion.

### **Check your progress 2**

1) Let us assume that the title of the study is “Use of textbook collection in college libraries”. The objectives of the study could be stated as to know: a. The adequacy of the text book collection; b. Subjects whose collection needs to be strengthened; and c. Improving upon the collection.

2) Hypothesis is an assumption, presumption, or in simple words guess towards some situation or condition. It is an assumption of relationships between the dependent and independent variables. It provides the researcher a line of action along which he moves to find out answers to the problem he is working on. Null hypothesis is stated as a negative relationship between dependent and independent variables. The positive relation between the variables is called the alternative hypothesis.

### **Check your progress 3**

1) In a research plan, ethical considerations demand that one should not take a topic of research that has already been studied. If it has to be studied, it should be done from some different perspective or context. Plagiarism should not be resorted to and work of others should be acknowledged. Research grants should be asked for only if genuinely required.

---

# UNIT 13: HYPOTHESIS AND VARIABLES

---

## STRUCTURE

- 13.0 Objectives
- 13.1 Introduction
- 13.2 Meaning and Characteristics of Hypothesis
- 13.3 Formulation of Hypothesis
- 13.4 Possible Difficulties in Formulation of a Good Hypothesis
- 13.5 Types of Hypotheses
  - 13.5.1 Null Hypothesis
  - 13.5.2 Alternative Hypothesis
- 13.6 Errors in Testing a Hypothesis
- 13.7 Importance of Hypothesis Formulation
- 13.8 Variables
- 13.9 Meaning of Variables
- 13.10 Types of Variables
  - 13.10.1 Stimulus, Organism and Response Variables
  - 13.10.2 Independent and Dependent Variables
  - 13.10.3 Extraneous and Confounded Variables
  - 13.10.4 Active and Attribute Variables
  - 13.10.5 Quantitative and Categorical Variables
  - 13.10.6 Continuous Variables and Discrete Variables
- 13.11 Let us sum up
- 13.12 Key Words
- 13.13 Questions for Review
- 13.14 Suggested readings and references
- 13.15 Answers to Check Your Progress

---

## 13.0 OBJECTIVES

---

After reading this unit, you will be able to:

- Define and describe hypothesis and its characteristics;
- explain formulation of hypothesis;



- Enumerate the possible difficulties in formulating hypothesis;
- Explain types of hypotheses;
- Identify in hypotheses testing;
- Define variables
- Meaning of Variables
- Identify different types of variables i.e. independent variable, dependent variable, extraneous variables etc. in a research study.

---

## 13.1 INTRODUCTION

---

Scientific process or all empirical sciences are recognised by two inter-related concepts, namely; (a) context of discovery (getting an idea) and (b) context of justification (testing and results). Hypotheses are the mechanism and container of knowledge moving from the unknown to known. These elements form techniques and testing ground for scientific discovery. Hypotheses are tentative explanation and potential answer to a problem. Hypothesis gives the direction and helps the researcher interpret data. In this unit, you will be familiarised with the term hypothesis and its characteristics. It is, then, followed by the hypothesis formulation and types of hypothesis. Errors in hypothesis testing are also highlighted. Further, In order to test the hypothesis, researcher rarely collects data on entire population owing to high cost and dynamic nature of the individual in population. Therefore, they collect data from a subset of individual – a sample - and make the inferences about entire population. This leads us to what we should know about the population and sample. So, researcher plans sample design and uses various method of sampling. This unit will acquaint you with the meaning of sampling and basic terminology which is used in sampling design. Now, it will move to purpose of sampling. And finally, various probability and non-probability sampling methods along with advantages and disadvantages are described.

When a possible correlation or similar relation between phenomena is investigated, such as whether a proposed remedy is effective in treating a disease, the hypothesis that a relation exists cannot be examined the same way one might examine a proposed new law of nature. In such an

## Notes

investigation, if the tested remedy shows no effect in a few cases, these do not necessarily falsify the hypothesis. Instead, statistical tests are used to determine how likely it is that the overall effect would be observed if the hypothesized relation does not exist. If that likelihood is sufficiently small (e.g., less than 1%), the existence of a relation may be assumed. Otherwise, any observed effect may be due to pure chance.

In statistical hypothesis testing, two hypotheses are compared. These are called the null hypothesis and the alternative hypothesis. The null hypothesis is the hypothesis that states that there is no relation between the phenomena whose relation is under investigation, or at least not of the form given by the alternative hypothesis. The alternative hypothesis, as the name suggests, is the alternative to the null hypothesis: it states that there is some kind of relation. The alternative hypothesis may take several forms, depending on the nature of the hypothesized relation; in particular, it can be two-sided (for example: there is some effect, in a yet unknown direction) or one-sided (the direction of the hypothesized relation, positive or negative, is fixed in advance).

Conventional significance levels for testing hypotheses (acceptable probabilities of wrongly rejecting a true null hypothesis) are .10, .05, and .01. The significance level for deciding whether the null hypothesis is rejected and the alternative hypothesis is accepted must be determined in advance, before the observations are collected or inspected. If these criteria are determined later, when the data to be tested are already known, the test is invalid.

The above procedure is actually dependent on the number of the participants (units or sample size) that are included in the study. For instance, to avoid having the sample size be too small to reject a null hypothesis, it is recommended that one specify a sufficient sample size from the beginning. It is advisable to define a small, medium and large effect size for each of a number of important statistical tests which are used to test the hypotheses.

---

## 13.2 MEANING AND CHARACTERISTICS OF HYPOTHESIS

---

A hypothesis (plural hypotheses) is a proposed explanation for a phenomenon. For a hypothesis to be a scientific hypothesis, the scientific method requires that one can test it. Scientists generally base scientific hypotheses on previous observations that cannot satisfactorily be explained with the available scientific theories. Even though the words "hypothesis" and "theory" are often used synonymously, a scientific hypothesis is not the same as a scientific theory. A working hypothesis is a provisionally accepted hypothesis proposed for further research, in a process beginning with an educated guess or thought.

A different meaning of the term hypothesis is used in formal logic, to denote the antecedent of a proposition; thus in the proposition "If P, then Q", P denotes the hypothesis (or antecedent); Q can be called a consequent. P is the assumption in a (possibly counterfactual) What If question.

The adjective hypothetical, meaning "having the nature of a hypothesis", or "being assumed to exist as an immediate consequence of a hypothesis", can refer to any of these meanings of the term "hypothesis".

In conducting research, the important consideration after the formulation of a research problem is the construction of hypothesis. As you know, any scientific inquiry starts with the statement of a solvable problem, when the problem has been stated, a tentative solution in the form of testable proposition is offered by the researcher. Hypothesis is often considered a tentative and testable statement of the possible relationship between two or more events / variables under investigation. According to Mcguigan (1990), 'a testable statement of a potential relationship between two or more variables, i.e. advance as potential solution to the problem'. Kerlinger (1973) defined 'a hypothesis is a conjectural statement of the relation between two or more variables'. In order to be useful in any study, the hypothesis needs to be stated in such a way that it might be subjected to empirical testing. The researcher is responsible to

suggest or find some way to check how the hypothesis stands against empirical data. When a hypothesis is formulated, the investigator must determine usefulness of the formulated hypothesis. There are several criteria or characteristics of a good research hypothesis. A good hypothesis is one which meets such criteria to a large extent. Some of these characteristics are enumerated below:

- 1) Hypothesis should be conceptually clear;
- 2) Hypothesis must be testable;
- 3) Hypothesis should be related to the existing body or theory and impact;
- 4) Hypothesis should have logical unity and comprehensiveness;
- 5) Hypothesis should be capable of verification; and
- 6) Hypothesis should be operationisable.

---

### **13.3 FORMULATION OF HYPOTHESIS**

---

Science proceeds with observation, hypothesis formulation and hypothesis testing. After testing the hypothesis, through various statistical tests, researcher can accept or reject the hypothesis. If the hypothesis is accepted then researcher can replicate the results, if hypothesis is rejected then researcher can refined or modify the results. By stating a specific hypothesis, the researcher narrows the focus of the data collection effort and is able to design a data collection procedure which is aimed at testing the plausibility of the hypothesis as a possible statement of the relationship between the terms of the research problem. It is, therefore, always useful to have a clear idea and vision about the hypothesis. It is essential for the research question as the researcher intends to verify, as it will direct and greatly help to interpretation of the results.

---

## 13.4 POSSIBLE DIFFICULTIES IN FORMULATION OF A GOOD HYPOTHESIS

---

There are three major possible difficulties; a researcher could face during formulation of hypothesis. First, the absence of knowledge of a theoretical framework is a major difficulty in formulating a good research hypothesis. Second, if detailed theoretical evidences are not available or if the investigator is not aware of the availability of those theoretical evidences, a research hypothesis cannot be formulated. Third, when the investigator is not aware of the scientific research techniques, she/he will not be able to frame a good research hypothesis. Despite these difficulties, the investigator attempts in her/his research to formulate a hypothesis. Usually the hypothesis is derived from the problem statement. The hypothesis should be formulated in a positive and substantive form before data are collected. In some cases additional hypothesis may be formulated after collection of data, but they should be tested on a new set of data and not on the old set which has suggested it. The formulation of a hypothesis is a creative task and involves a lot of thinking, imagination and innovation. Reichenbach (1938) has made a distinction between the two processes found commonly in any hypothesis formulation task. One is the context of discovery and another is the context of justification. The manner or the process through which a scientist arrives at a hypothesis illustrates the context of justification. A scientist is concerned more with a context of justification in the development of a hypothesis. He never puts his ideas or thoughts as they nakedly occur in the formulation of a hypothesis. Rather, he logically reconstructs his ideas or thoughts and draws some justifiable inferences from those ideas and thoughts. He never cares to relate how he actually arrived at a hypothesis. He does not say, for example, that while he was shaving, this particular hypothesis occurred to him. He usually arrives at a hypothesis by the rational reconstruction of thoughts. When a scientist reconstructs his thoughts and communicates them in the form of a hypothesis to others, he uses the context of justification. When he arrives at a hypothesis, he extensively as well as intensively surveys a mass of data, abstracts them, tries to find out similarities among the abstracted

## Notes

data and finally makes a generalisation or deduces a proposition in the form of a hypothesis.

Here is an important distinction to be made between formulating a hypotheses and choosing one. Although a researcher often becomes interested in a question about human behaviour for personal reasons, the ultimate value of research study depends on the researcher bringing methodological criteria to bear on the selection of the hypothesis to be tested. In other words, Good hypothesis are made, not born. Hypothesis plays a key role in formulating and guiding any study. The hypotheses are generally derived from earlier research findings, existing theories and personal observations and experience. For instance, you are interested in knowing the effect of reward on learning. You have analysed the past research and found that two variables are positively related. You need to convert this idea in terms of a testable statement. At this point you may develop the following hypothesis. Those who are rewarded shall require lesser number of trails to learn the lesson than those who are not rewarded. A researcher should consider certain points while formulating a hypothesis:

- i) Expected relationship or differences between the variables.
- ii) Operational definition of variable.
- iii) Hypotheses are formulated following the review of literature  
The literature leads a researcher to expect a certain relationship. Hypotheses are the statement that is assumed to be true for the purpose of testing its validity.

As suggested by Russell and Reichenback (1947), the hypotheses should be stated in the logical form on the general implications. A hypothesis can be put in the form of an if ..... then statement; if A is true then B should follow. For example, verbal development theory of amnesia states that childhood amnesia caused by the development of language. To test this theory, researcher can make a hypothesis like this – if the lack of verbal ability is responsible for childhood amnesia, then the children

should not be able to verbally recall events usually words that they did not know at the time of events.

**Check Your Progress 1**

Notes: a) Space is given below for writing your answers.  
b) Compare your answers with those given at the end of the unit.

1) Discuss the Meaning and Characteristics of Hypothesis.

.....  
.....  
.....  
.....  
.....

2) Discuss the Formulation of Hypothesis.

.....  
.....  
.....  
.....  
.....

3) Possible Difficulties in Formulation of a Good Hypothesis.

.....  
.....  
.....  
.....  
.....

---

**13.5 TYPES OF HYPOTHESES**

---

As explained earlier, any assumption that you seek to validate through investigation is called hypotheses. Hence theoretically, there should be one type of hypotheses on the basis of the investigation that is, research hypothesis. However, because of the conventions in scientific enquiries and wording used in the constructions of the hypothesis, Hypotheses can be classified into several types, like; universal hypotheses, existential

hypotheses, conceptual hypotheses etc. Broadly, there are two categories of the hypothesis:

- i) Null hypothesis
- ii) Alternative hypothesis

### 13.5.1 Null Hypothesis

Null hypothesis is symbolised as  $H_0$ . Null hypothesis is useful tool in testing the significance of difference. In its simplest form, this hypothesis asserts that there is no true difference between two population means, and the difference found between sample means is, accidental and unimportant, that is arising out of fluctuation of sampling and by chance. Traditionally null hypothesis stated that there is zero relationship between terms of the hypothesis. For example, (a) schizophrenics and normal do not differ with respect to digit span memory (b) There is no relationship between intelligence and height. The null hypothesis is an important component of the decision making methods of inferential statistics. If the difference between the samples of means is found significant the researcher can reject the null hypothesis. It indicates that the differences have statistically significant and acceptance of null hypothesis indicates that the differences are due to chance. Null hypothesis should always be specific hypothesis i.e. it should not state about or approximately a certain value. The null hypothesis is often stated in the following way:  $H_0 : \mu_{HV}$ .

### 13.5.2 Alternative Hypothesis

Alternative hypothesis is symbolised as  $H_1$  or  $H_a$ , is the hypothesis that specifies those values that are researcher believes to hold true, and the researcher hopes that sample data will lead to acceptance of this hypothesis as true. Alternative hypothesis represents all other possibilities and it indicates the nature of relationship. The alternative hypothesis is stated as follows:  $H_1 : \mu_{HV} > \mu_{LV}$  The alternative hypothesis is that the mean of population of those who have the vocabulary is greater than the mean of those to lack the vocabulary. In this



example the alternative hypothesis is that the experimental population had higher mean than the controls. This is called directional hypothesis because researcher predicted that the high vocabulary children would differ in one particular direction from the low vocabulary children. Sometimes researcher predicts only that the two groups will differ from each other but the researcher doesn't know which group will be higher. This is non directional hypothesis. The null and alternative hypothesis in this case would be stated as follows:  $H_0 : \mu_1 = \mu_2$   $H_1 : \mu_1 \neq \mu_2$  Thus, the null hypothesis is that mean of group 1 equals the mean of group 2, and the alternative hypothesis is that the mean of group 1 does not equal the mean of group 2.

---

## 13.6 ERRORS IN TESTING A HYPOTHESIS

---

You have already learned that hypotheses are assumptions that may be prove to be either correct or incorrect. It is possible to arrive at a incorrect conclusion about a hypothesis for the various reasons if –

- Sampling procedure adopted faulty
- Data collection method inaccurate
- Study design selected is faulty
- Inappropriate statistical methods used
- Conclusions drawn are incorrect Two common errors exist when testing a hypothesis.

Type I error – Rejection of a null hypothesis when it is true.

Type II error - Acceptance of a null hypothesis when it is false.

In statistics, a null hypothesis is a statement that one seeks to nullify (that is, to conclude is incorrect) with evidence to the contrary. Most

## Notes

commonly, it is presented as a statement that the phenomenon being studied produces no effect or makes no difference. An example of such a null hypothesis might be the statement, "A diet low in carbohydrates has no effect on people's weight." An experimenter usually frames a null hypothesis with the intent of rejecting it: that is, intending to run an experiment which produces data that shows that the phenomenon under study does indeed make a difference (in this case, that a diet low in carbohydrates over some specific time frame does in fact tend to lower the body weight of people who adhere to it). In some cases there is a specific alternative hypothesis that is opposed to the null hypothesis, in other cases the alternative hypothesis is not explicitly stated, or is simply "the null hypothesis is false" — in either event, this is a binary judgment, but the interpretation differs and is a matter of significant dispute in statistics.

A type I error (or error of the first kind) is the rejection of a true null hypothesis. Usually a type I error leads to the conclusion that a supposed effect or relationship exists when in fact it does not. Examples of type I errors include a test that shows a patient to have a disease when in fact the patient does not have the disease, a fire alarm going on indicating a fire when in fact there is no fire, or an experiment indicating that a medical treatment should cure a disease when in fact it does not.

A type II error (or error of the second kind) is the failure to reject a false null hypothesis. Some examples of type II errors are a blood test failing to detect the disease it was designed to detect, in a patient who really has the disease; a fire breaking out and the fire alarm does not ring; or a clinical trial of a medical treatment failing to show that the treatment works when really it does.

In terms of false positives and false negatives, a positive result corresponds to rejecting the null hypothesis, while a negative result corresponds to failing to reject the null hypothesis; "false" means the conclusion drawn is incorrect. Thus a type I error is a false positive, and a type II error is a false negative.

When comparing two means, concluding the means were different when in reality they were not different is a type I error; concluding the means were not different when in reality they were different is a type II error. Various extensions have been suggested as "type III errors", though none have wide use[according to whom?].

All statistical hypothesis tests have a probability of making type I and type II errors. For example, all blood tests for a disease will falsely detect the disease in some proportion of people who do not have it, and will fail to detect the disease in some proportion of people who do have it. A test's probability of making a type I error is denoted by  $\alpha$ . A test's probability of making a type II error is denoted by  $\beta$ . These error rates are traded off against each other: for any given sample set, the effort to reduce one type of error generally results in increasing the other type of error. For a given test, the only way to reduce both error rates is to increase the sample size, and this may not be feasible. A test statistic is robust if the Type I error rate is controlled.

These terms are also used in a more general way by social scientists and others to refer to flaws in reasoning.

Fire alarm analogy for the types of errors

Sign that says fire alarm. A fire alarm provides a good analogy for the types of hypothesis testing errors. Preferably, the alarm rings when there is a fire and does not ring in the absence of a fire. However, if the alarm rings when there is no fire, it is a false positive, or a Type I error in statistical terms. Conversely, if the fire alarm fails to ring when there is a fire, it is a false negative, or a Type II error.

Using hypothesis tests correctly improves your chances of drawing trustworthy conclusions. However, errors are bound to occur.

Unlike the fire alarm analogy, there is no sure way to determine whether an error occurred after you perform a hypothesis test. Typically, a clearer

## Notes

picture develops over time as other researchers conduct similar studies and an overall pattern of results appears. Seeing how your results fit in with similar studies is a crucial step in assessing your study's findings.

Now, let's take a look at each type of error in more depth.

### Type I Errors: False Positives

When you see a p-value that is less than your significance level, you get excited because your results are statistically significant. However, it could be a type I error. The supposed effect might not exist in the population. Again, there is usually no warning when this occurs.

Why do these errors occur? It comes down to sample error. Your random sample has overestimated the effect by chance. It was the luck of the draw. This type of error doesn't indicate that the researchers did anything wrong. The experimental design, data collection, data validity, and statistical analysis can all be correct, and yet this type of error still occurs.

Even though we don't know for sure which studies have false positive results, we do know their rate of occurrence. The rate of occurrence for Type I errors equals the significance level of the hypothesis test, which is also known as alpha ( $\alpha$ ).

The significance level is an evidentiary standard that you set to determine whether your sample data are strong enough to reject the null hypothesis. Hypothesis tests define that standard using the probability of rejecting a null hypothesis that is actually true. You set this value based on your willingness to risk a false positive.

Using the significance level to set the Type I error rate

When the significance level is 0.05 and the null hypothesis is true, there is a 5% chance that the test will reject the null hypothesis incorrectly. If you set alpha to 0.01, there is a 1% of a false positive. If 5% is good,

then 1% seems even better, right? As you'll see, there is a tradeoff between Type I and Type II errors. If you hold everything else constant, as you reduce the chance for a false positive, you increase the opportunity for a false negative.

Type I errors are relatively straightforward. The math is beyond the scope of this article, but statisticians designed hypothesis tests to incorporate everything that affects this error rate so that you can specify it for your studies. As long as your experimental design is sound, you collect valid data, and the data satisfy the assumptions of the hypothesis test, the Type I error rate equals the significance level that you specify. However, if there is a problem in one of those areas, it can affect the false positive rate.

Warning about a potential misinterpretation of Type I errors and the Significance Level

When the null hypothesis is correct for the population, the probability that a test produces a false positive equals the significance level. However, when you look at a statistically significant test result, you cannot state that there is a 5% chance that it represents a false positive.

Why is that the case? Imagine that we perform 100 studies on a population where the null hypothesis is true. If we use a significance level of 0.05, we'd expect that five of the studies will produce statistically significant results—false positives. Afterward, when we go to look at those significant studies, what is the probability that each one is a false positive? Not 5 percent but 100%!

That scenario also illustrates a point that I made earlier. The true picture becomes more evident after repeated experimentation. Given the pattern of results that are predominantly not significant, it is unlikely that an effect exists in the population.

Type II Errors: False Negatives

## Notes

When you perform a hypothesis test and your p-value is greater than your significance level, your results are not statistically significant. That's disappointing because your sample provides insufficient evidence for concluding that the effect you're studying exists in the population. However, there is a chance that the effect is present in the population even though the test results don't support it. If that's the case, you've just experienced a Type II error, which is also known as beta ( $\beta$ ).

What causes Type II errors? Whereas Type I errors are caused by one thing, sample error, there are a host of possible reasons for Type II errors—small effect sizes, small sample sizes, and high data variability. Furthermore, unlike Type I errors, you can't set the Type II error rate for your analysis. Instead, the best that you can do is estimate it before you begin your study by approximating properties of the alternative hypothesis that you're studying. When you do this type of estimation, it's called power analysis.

To estimate the Type II error rate, you create a hypothetical probability distribution that represents the properties of a true alternative hypothesis. However, when you're performing a hypothesis test, you typically don't know which hypothesis is true, much less the specific properties of the distribution for the alternative hypothesis. Consequently, the true Type II error rate is usually unknown!

Type II errors and the power of the analysis

The Type II error rate (beta) is the probability of a false negative. Therefore, the inverse of Type II errors is the probability of correctly detecting an effect. Statisticians refer to this concept as the power of a hypothesis test. Consequently,  $1 - \beta =$  the statistical power. Analysts typically estimate power rather than beta directly.

If you read my post about power and sample size analysis, you know that the three factors that affect power are sample size, variability in the population, and the effect size. As you design your experiment, you can

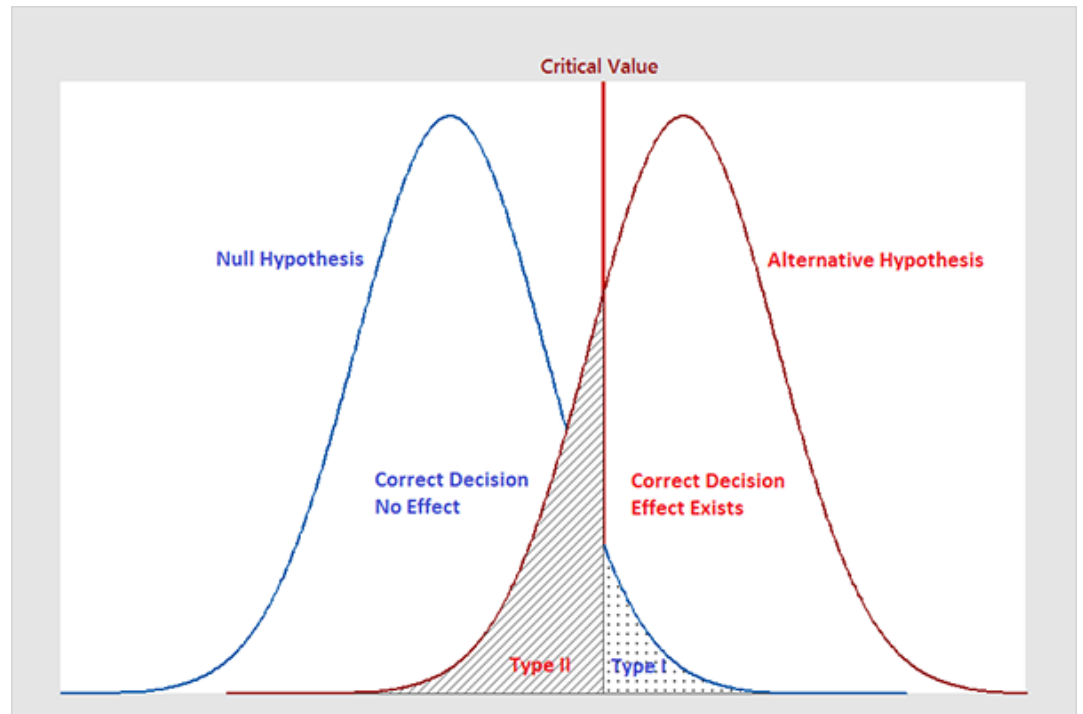
enter estimates of these three factors into statistical software and it calculates the estimated power for your test.

Suppose you perform a power analysis for an upcoming study and calculate an estimated power of 90%. For this study, the estimated Type II error rate is 10% ( $1 - 0.9$ ). Keep in mind that variability and effect size are based on estimates and guesses. Consequently, power and the Type II error rate are just estimates rather than something you set directly. These estimates are only as good as the inputs into your power analysis.

Low variability and larger effect sizes decrease the Type II error rate, which increases the statistical power. However, researchers usually have less control over those aspects of a hypothesis test. Typically, researchers have the greatest control over sample size, which makes it the critical way to manage your Type II error rate. Holding everything else constant, increasing the sample size reduces the Type II error rate and increases power.

#### Graphing Type I and Type II Errors

The graph below illustrates the two types of errors using two sampling distributions. The critical region line represents the point at which you reject or fail to reject the null hypothesis. Of course, when you perform the hypothesis test, you don't know which hypothesis is correct. And, the properties of the distribution for the alternative hypothesis are usually unknown. However, use this graph to understand the general nature of these errors and how they are related.



The distribution on the left represents the null hypothesis. If the null hypothesis is true, you only need to worry about Type I errors, which is the shaded portion of the null hypothesis distribution. The rest of the null distribution represents the correct decision of failing to reject the null.

On the other hand, if the alternative hypothesis is true, you need to worry about Type II errors. The shaded region on the alternative hypothesis distribution represents the Type II error rate. The rest of the alternative distribution represents the probability of correctly detecting an effect—power.

Moving the critical value line is equivalent to changing the significance level. If you move the line to the left, you're increasing the significance level (e.g.,  $\alpha$  0.05 to 0.10). Holding everything else constant, this adjustment increases the Type I error rate while reducing the Type II error rate. Moving the line to the right reduces the significance level (e.g.,  $\alpha$  0.05 to 0.01), which decreases the Type I error rate but increases the type II error rate.

Is One Error Worse Than the Other?



As you've seen, the nature of the two types of error, their causes, and the certainty of their rates of occurrence are all very different.

A common question is whether one type of error is worse than the other? Hypothesis tests are designed to be able to control Type I errors while Type II errors are much less defined. Consequently, many statisticians state that it is better to fail to detect an effect when it exists than it is to conclude an effect exists when it doesn't. That is to say, there is tendency to assume that Type I errors are worse.

However, reality is more complex than that. You should carefully consider the consequences of each type of error for your specific test.

Suppose you are assessing the strength of a new jet engine part that is under consideration. Peoples lives are riding on the part's strength. A false negative in this scenario merely means that the part is strong enough but the test fails to detect it. This situation does not put anyone's life at risk. On the other hand, Type I errors are worse in this situation because they indicate the part is strong enough when it is not.

Now suppose that the jet engine part is already in use but there are concerns about it failing. In this case, you want the test to be more sensitive to detecting problems even at the risk of false positives. Type II errors are worse in this scenario because the tests fail to detect the problem and leave these problematic parts in use for longer.

Using hypothesis tests effectively requires that you understand their error rates. By setting the significance level and estimating the power of your test, you can manage both error rates so they meet your requirements.

---

## **13.7 IMPORTANCE OF HYPOTHESIS FORMULATION**

---

Hypothesis is the basic function of the scientific research. If simple, brief and clear scientific hypothesis has been formulated, there shall be no

## Notes

problem for the investigator to proceed in the research field. Its utility or importance for and research may be studied as under.

Accordingly to Goode and Hatt ('without' hypothesis formulation the research is unfocussed, a random empirical wandering. The results cannot be studied as facts with clear meaning. Formulation of hypothesis links between theory and investigation which lead to discovery of addition to knowledge.

### Check Your Progress 2

Notes: a) Space is given below for writing your answers.

b) Compare your answers with those given at the end of the unit.

1. Discuss the Types of Hypotheses.

.....  
.....  
.....  
.....  
.....

2. What do you mean by Errors in Testing a Hypothesis?

.....  
.....  
.....  
.....  
.....

3. Importance of Hypothesis Formulation.

.....  
.....  
.....  
.....  
.....

---

## 13.8 VARIABLES

---

In the process of formulating a research problem there are two important considerations; the use of constructs/concepts and the construction of hypotheses. Constructs/concepts are highly subjective as their understanding varies from person to person and therefore, as such, may not be measurable. In a research study, it is important that the concepts used should be operationalised in measurable terms so that the extent of variation in respondents understanding is reduced if not eliminated. Knowledge about constructs and variables are very important to understand conceptual clarity and quantitative accuracy as they provide the ‘fine tuning’ to research. This unit attempts to acquaint you with the term variables and constructs which are used by the psychologists in gaining knowledge about the behaviour and mental processes. It begins with definition of variables then you will find the details about the types of variables along with the examples. Further, you will be exposed to the nature of the scientific concept or construct and the way in which behavioural scientist travel from the construct level to observation level. Finally, types of constructs are described.

---

## 13.9 MEANING OF VARIABLES

---

A variable, as the name implies, is something that varies. This is the simplest way of defining a variable. Webster says that a variable is “a thing that is changeable” or “a quantity that may have a number of different values.” True, a variable is something that has at least two values: however, it is also important that the values of the variable be observable. Thus, if what is being studied is a variable, it has more than one value and each value can be observed. For example, the outcome of throwing a dice is a variable. That variable has six possible values (each side of the dice having one to six dots on it), each of which can be observed. However, a behavioural scientist attempts to define a variable more precisely and specifically. Kerlinger (1986) defined variable ‘a property that taken as different values’. According to D’Amato (1970) variables may be defined as those attributes of objects, events, things and beings, which can be measured. According to Postman and Egan (1949)

## Notes

a variable is a characteristic or attribute that can take on a number of values, for example, number of items that an individual solves on a particular test, the speed with which we respond to a signal, IQ, sex, level of anxiety, and different degree of illumination are the examples of variables that are commonly employed in psychological research.

Before discussing the types of variables, it is important to know how the variables of study related to theoretical concepts. Because the variables exist in the world but the theory is an idea, researcher makes certain assumption to relate the two. These assumptions are guide ropes that tie a theory to the real world. The variables are tangible: duration, frequency, rate, or intensity of bar presses; items checked on a questionnaire; murders committed; books written. The theoretical concept is intangible: hunger, motivation, anxiety. The variables are related to the theoretical concepts by means of the operational definitions used to measure the concepts.

Suppose a theory reveals that increasing anxiety will increase the affiliation motive. To test out this theory, you may take the theoretical concepts of anxiety and affiliation motive and relate them to variables by means of operational definitions. The theory is an abstract statement. For example, anxiety can be measured by the Anxiety Scale and affiliation by how close people sit to each other in the experiment. These two measures constitute the variables of the study. The scores on the variables of anxiety and distance apart are related to one another as test of the hypothesis. The relationship between the variable is taken as providing support for or against the particular theory that generated the experiment.

---

### **13.10 TYPES OF VARIABLES**

---

To understand how variables are used and discussed in psychological researches, you would like to understand several distinctions that are made among the type of variables. The descriptions of different types of variables are given below:

### 13.10.1 Stimulus, Organism and Response

#### Variables

Psychologists are interested in studying the behaviour or causes of behaviour as variables. Many psychologists have adopted a theoretical viewpoint or model called the S-O-R model to explain all behaviour. The symbols S, O, and R represent different categories of variables. S is the symbol of stimuli, and the category may be referred to in general as stimulus variables. A stimulus variable is some form of energy in the environment, such as light, to which the organism is sensitive. O is the symbol for organism variables, that is the changeable physiological and psychological characteristics of the organisms being observed. Examples of such variables are anxiety level, age and heart rate etc. Finally, R is the symbol for response and, in general, response variables, which refer to some behaviour or action of the organism like pressing a lever, and reaction to any stimulus, are the examples of responses variables. You can understand an application of S-O-R model through the following example. Suppose that an experiment is conducted in which a rat is placed on a metal grid floor, the grid is electrified, and the length of time it takes the rat to jump from the grid to a platform is measured. Using the S-O-R model, the electrical shock would be called a stimulus variable. The intensity of shock would be the value of the variable. The particular state of the organisms would be measured by the organismic variables. For example, the skin resistant of the rat at the time of shock was introduced would be an organismic variables. A response variable would be the latency (i.e. the elapsed time between the onsets of the shock and when the rat reaches the platform).

#### 13.10.2 Independent and Dependent Variables

An independent variable or stimulus variable (as Underwood calls it) is that factor manipulated or selected by the experimenter in his attempt to ascertain its relationship to an observed phenomenon. Dependent upon the mode of manipulation, some experts have tried to divide the

## Notes

independent variable into 'Type E' independent variable and 'Type S' independent variable (D'Amato, 1970). Type E independent variable is one of which is directly or experimentally manipulated by the experimenter and type S independent variable is one which is manipulated through the process of selection only. For example the experimenter wants to study the effect of noise upon the task performance in an industry. Here the IV (Independent Variable) is the noise and the DV (Dependent Variable) is the task performance. He may manipulate the noise by dividing into three categories — continuous noise, intermittent noise and no noise and examine its effect upon the task performance. Here the noise is being directly manipulated by the experimenter and hence, it constitutes the example of Type-E independent variable. Suppose, for the time being, that the experimenter is interested in answering the question: Is the rate of production dependent upon the age of the workers? Age is here the independent variable. For investigating this problem, the experimenter will have to select groups of workers on the basis of their age in a way by which he can get an appropriate representation from different age groups ranging from say, 16 to 55 years. Subsequently, he will compare the rate of production obtained by each age group and finally, conclude whether or not age is a factor in enhancement of the performance. Hence this constitutes the examples of S-independent variables. A dependent variable is the factor that appears, disappears, or varies as the experimenter introduces, removes or varies the independent variable. (Townsend, 1953). The dependent variable is a measure of the behaviour of the subject. The dependent variable is the response that the person or animal makes. This response is generally measured using at least one of several different dimensions (Alberto & Troutman 2006). The dimensions are – (a) frequency – Number of times that a particular behaviour occurs, ( b) duration - the amount of time that a behaviour lasts. (c) latency –the amount of time between and when the behaviour is actually performed (d) force – the intensity or strength of a behaviour. Here, you can examine the relationship between independent and dependent variables. The relationship is that of dependence. One variable depends upon the other. Suppose you find a relationship between

meaningfulness of the learning material and speed of learning. Speed of learning then depends upon meaningfulness; the greater the meaningfulness, the faster the learning. The speed of learning is, therefore, called dependent variable; meaningfulness is independent variable. Similarly, rest between work periods is independent variables; output of work is dependent variable. Sudden noise is independent variable; change in breathing is dependent variable. In an experiment one discovers and confirms a relationship between an independent variable and a dependent variable.

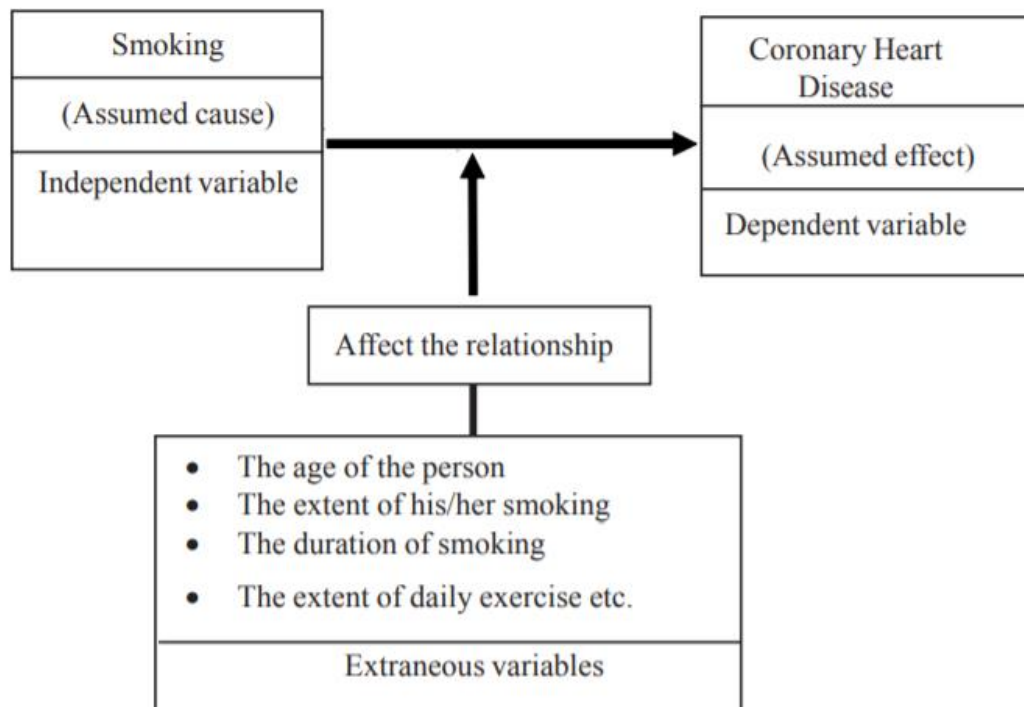
### **13.10.3 Extraneous and Confounded Variables**

Any and all other variables that may 'mask' the relationship between independent variable and dependent variable are known as extraneous variables. Extraneous variables may directly affect the dependent variable or may combine with the independent variable to produce an affect. Therefore, extraneous variables must be controlled so that the experimenter can determine whether the dependent variable changes in relation to variation in the independent variable. Several others factors operating in a real life situation may affect changes in the dependent variable. These factors, not measured in the study, may increase or decrease the magnitude or strength of the relationship between independent and dependent variables. Extraneous variables are relevant in nature, and in experimental studies, they belong to three major types i.e., organismic variables, situational variables and sequential variables. The subject related variables include age, sex, intelligence, personality etc. are organismic variables. The situational variables include environmental variables operating in the experimental setting (e.g. noise, temperature, humidity) and variables related to the experimental task. The sequence related variables deal with sequence effects. They arise when participants in experiments are required to be tested in several conditions. Exposure to many conditions may result in adaptation, fatigue or practice effects which, if allowed to operate, may make the results difficult to interpret.

## Notes

Confounding variables Variables and Constructs is one that varies with the independent variable. While doing a study if we are not careful then two variables may get combined so that the effect of one cannot be separated from the effect of other. This is known as confounding. For instance, if you conducted a study of the effect of television viewing on perception of violence and the experimental group contained only adolescents, whereas the control group only adults, the age of participants would be confounded with the independent variable under study. Confounding makes the conclusions of the study doubtful. It is ,therefore, necessary that effort should be made to unconfound the variables. To explain these variables let us take one example. Suppose you want to study the relationship between smoking and coronary heart disease. You assume that affecting this relationship, such as a number of cigarettes or the amount of tobacco smoked every day; the duration of smoking; the age of the smoker; dietary habits; and the amount of exercise undertaken by the individuals. All of these factors may affect the extent to which smoking might cause coronary heart disease. These variables may either increase or decrease the magnitude of the relationship. In this example, the extent of smoking is the independent variable, coronary heart disease is the dependent variable and all the variables that might affect this relationship, either positively or negatively, are extraneous variables. Independent, dependent & extraneous variables in a causal relationship





### 13.10.4 Active and Attribute Variables

Any variable that is manipulated is called active variables. Examples of active variables are reward, punishment, methods of teaching, creating anxiety through instructions and so on. Attribute variable is that variable which is not manipulated but measured by the experimenter. Variables that are human characteristics like intelligence, Aptitudes, sex, socio economic status, education, field dependence and need for achievement are the example of attributes variables. The word 'attribute' is more accurate enough when used within animated objects or references. Organisations, institutions, groups, population and geographical areas have attributes. Organisations are variably productive; groups differ in cohesiveness; geographical areas vary widely in resources.

### 13.10.5 Quantitative and Categorical Variables

Quantitative variables is one that varies in amount whereas categorical variables varies in kind. Speed of response, intensity of sound, level of Illumination, intelligence etc. are the example of quantitative variables and gender, race, religion are the example of categorical variables.

## Notes

Precise and accurate measurement are possible with the quantitative variables because they can be easily ordered in terms of increasing and decreasing magnitude. Categorical variables can be of three types: Constant, dichotomous and polytomous. When a variable can have only one value or category, for example taxi, tree and water, it is known as a constant variable. When a variable can have only two categories as in yes/no, good/bad and rich/poor, it is known as dichotomous variables. When variables can be divided into more than two categories, for example: religion (Christian, Muslim, Hindu); political parties (Labor, Liberal, Democrat); and attitudes (strongly favorable, favorable, uncertain, unfavorable, strongly unfavorable), it is called a polytomous variable.

### **13.10.6 Continuous Variables and Discrete Variables**

Quantitative variables are further divided into two categories, namely, continuous variables and discrete variables. A distinction between continuous and discrete variables is especially useful in planning of research and analysis of data. A continuous variable is one which is capable of being measured in any arbitrary degree of fineness or exactness. Age, height, intelligence, reaction time, etc., are some of the examples of a continuous variable. The age of the person can be measured in years, month and days. Thus, all such variables which can be measured in the smallest degree of fineness are called continuous variable. The discrete variables are those variables which are not capable of being measured in any arbitrary degree of fineness or exactness because the variables contain a clear gap. For example, the number of members in a family, no. of females in particular group, no of books in library and so on constitutes the examples of a discrete variable.

#### **Activity -1**

Check whether the following are continuous or discrete variables: C

D

a) the bar presses that a rat makes in a Skinner box ( ) ( )

b) the height of six-year-old boys and girls in Chicago ( ) ( )

c) the score you make on a true-false exam ( ) ( )

d) the distance various people can travel in 5 hours ( ) ( )

**Activity -2**

Identity Types of Variables A researcher wants to administer an intelligence test to 30 college students. After collecting information on subjects’ age, sex, height, weight, political preference, career goals, and socioeconomic status, the researcher administers and attitude survey on current issues to all 30 subjects.

Required: Identify examples of the following types of variables in the paragraph and the scales by which they would be measured: a) discrete b) continuous c) categorical d) quantitative.

**Check Your Progress 3**

Notes: a) Space is given below for writing your answers.  
b) Compare your answers with those given at the end of the unit.

1. What is Variables?

.....  
 .....  
 .....  
 .....  
 .....

2. Discuss about Meaning of Variables.

.....  
.....  
.....  
.....  
.....

3. What are the Types of Variables?

.....  
.....  
.....  
.....  
.....

---

### **13.11 LET US SUM UP**

---

Concepts in Hempel's deductive-nomological model play a key role in the development and testing of hypotheses. Most formal hypotheses connect concepts by specifying the expected relationships between propositions. When a set of hypotheses are grouped together they become a type of conceptual framework. When a conceptual framework is complex and incorporates causality or explanation it is generally referred to as a theory. According to noted philosopher of science Carl Gustav Hempel "An adequate empirical interpretation turns a theoretical system into a testable theory: The hypothesis whose constituent terms have been interpreted become capable of test by reference to observable phenomena. Frequently the interpreted hypothesis will be derivative hypotheses of the theory; but their confirmation or disconfirmation by empirical data will then immediately strengthen or weaken also the primitive hypotheses from which they were derived."

Hempel provides a useful metaphor that describes the relationship between a conceptual framework and the framework as it is observed and perhaps tested (interpreted framework). "The whole system floats, as it were, above the plane of observation and is anchored to it by rules of interpretation. These might be viewed as strings which are not part of the network but link certain points of the latter with specific places in the

plane of observation. By virtue of those interpretative connections, the network can function as a scientific theory." Hypotheses with concepts anchored in the plane of observation are ready to be tested. In "actual scientific practice the process of framing a theoretical structure and of interpreting it are not always sharply separated, since the intended interpretation usually guides the construction of the theoretician." It is, however, "possible and indeed desirable, for the purposes of logical clarification, to separate the two steps conceptually."

Knowledge of different types of variables and constructs play a crucial role in research. Variables and constructs are important in bringing clarity and specificity to the conceptualising of a research problem, to formulation of hypothesis and to the development of a research instrument. They affect how the data can be analysed, what statistical test can be applied to the data, what interpretation can be made and what conclusion can be drawn. A variable is some property of an event that takes on different values. There are different kinds of variables such as independent variables, dependent variables, quantitative variables and categorical variables, active and attribute variables, continuous and discrete variables, extraneous and intervening variables and so on. A construct is a concept. It has an added meaning and it is adopted for a special scientific purpose. Constructs are of two types; intervening and hypothetical constructs. Intervening variables is a term which is internal and directly unobservable psychological processes that, in turn, inferred from behaviour. A hypothetical construct is a theoretical term which is employed to describe something "real." That is, it is an intermediary which has tangible characteristics.

---

## **13.12 KEY WORDS**

---

Variable : A variable is a property that taken as different values.

Independent variables : The condition manipulated or selected by the experimenter to determine its effect on behaviour.

Dependent variables : A measure of the subject's behaviour that reflects that independent variable's effects.

Quantitative variable : One that varies in amount.

## Notes

Categorical variable : One that varies in kind.

Continuous variable : One that falls along a continuum and is not lifted to a certain number of values.

Discrete variable : One that falls into separate bins with no intermediate values possible.

Active variables : Manipulated variables are active variables.

Attribute variables : Measured variables are attribute variables.

Constructs : Is a concept, used for scientific purpose, is a part of theoretical framework.

Intervening variables : Is a construct which is utilised as a summary term for a group of other constructs.

Hypothetical constructs : Is a theoretical term which is employed to describe something real.

---

### **13.13 QUESTIONS FOR REVIEW**

---

- 1) Define variable and discuss the various kinds of variable.
- 2) Explain Intervening variables and Hypothetical constructs in your own words.
- 3) Write short notes on any two:
  - i) Independent & dependent variables.
  - ii) Quantitative & categorical variables.
  - iii) Active & attribute variables.
- 4) Is physical attractiveness related to friendship?
- 5) Does meaningful material affect the rate of learning?
- 6) Does reinforcement improve the learning for solving simple discrimination task?

---

### **13.14 SUGGESTED READINGS AND REFERENCES**

---

- Hilborn, Ray; Mangel, Marc (1997). *The ecological detective: confronting models with data*. Princeton University Press. p. 24. ISBN 978-0-691-03497-3. Retrieved 22 August 2011.
- "In general we look for a new law by the following process. First we guess it. ...", —Richard Feynman (1965) *The Character of Physical Law* p.156
- Wilbur R. Knorr, "Construction as existence proof in ancient geometry", p. 125, as selected by Jean Christianidis (ed.), *Classics in the history of Greek mathematics*, Kluwer.
- Gregory Vlastos, Myles Burnyeat (1994) *Socratic studies*, Cambridge ISBN 0-521-44735-6, p. 1
- "Neutral hypotheses, those of which the subject matter can never be directly proved or disproved, are very numerous in all sciences." — Morris Cohen and Ernest Nagel (1934) *An introduction to logic and scientific method* p. 375. New York: Harcourt, Brace, and Company.
- "Bellarmine (Ital. Bellarmino), Roberto Francesco Romolo", *Encyclopædia Britannica*, Eleventh Edition.: 'Bellarmine did not proscribe the Copernican system ... all he claimed was that it should be presented as a hypothesis until it should receive scientific demonstration.' This article incorporates text from a publication now in the public domain: Chisholm, Hugh, ed. (1911). "Hypothesis". *Encyclopædia Britannica*. 14 (11th ed.). Cambridge University Press. p. 208.
- Crease, Robert P. (2008) *The Great Equations* ISBN 978-0-393-06204-5, p.112 lists the conservation of energy as an example of accounting a constant of motion. Hypothesized by Sadi Carnot, truth demonstrated by James Prescott Joule, proven by Emmy Noether.
- *Compilers: Principles, Techniques, and Tools*, pp. 26–28
- Knuth, Donald (1997). *The Art of Computer Programming*. 1 (3rd ed.). Reading, Massachusetts: Addison-Wesley. p. 3-4. ISBN 0-201-89683-4.
- *How Not To Pick Variables*, Retrieved July 11, 2012 [DEAD LINK]

- Edsger Dijkstra, To hell with “meaningful identifiers”!

---

## **13.15 ANSWERS TO CHECK YOUR PROGRESS**

---

### **Check Your Progress 1**

1. See Section 13.2
2. See Section 13.3
3. See Section 13.4

### **Check Your Progress 2**

1. See Section 13.5
2. See Section 13.6
3. See Section 13.7

### **Check Your Progress 3**

1. See Section 13.8
2. See Section 13.9
3. See Section 13.10



---

# **UNIT 14: RESEARCH METHOD, PRIMARY AND SECONDARY DATA, STYLE AND REFERENCE, RESEARCH REPORT**

---

## **STRUCTURE**

- 14.0 Objectives
- 14.1 Introduction
- 14.2 Research method
- 14.3 Primary and secondary data
- 14.4 Style and reference
- 14.5 Research report
- 14.6 Let us sum up
- 14.7 Key Words
- 14.8 Questions for Review
- 14.9 Suggested readings and references
- 14.10 Answers to Check Your Progress

---

## **14.0 OBJECTIVES**

---

After this unit 14, we can able to understand:

- To discuss about the Research method;
- To know about Primary and secondary data;
- To know the Style and reference;
- To discuss the process of Research report.

---

## **14.1 INTRODUCTION**

---

Data (singular datum) are individual units of information. A datum describes a single quality or quantity of some object or phenomenon. In analytical processes, data are represented by variables.

Although the terms "data", "information" and "knowledge" are often used interchangeably, each of these terms has a distinct meaning. In

## Notes

popular publications, data is sometimes said to be transformed into information when it is viewed in context or in post-analysis.. In academic treatments of the subject, however, data are simply units of information. Data is employed in scientific research, businesses management (e.g., sales data, revenue, profits, stock price), finance, governance (e.g., crime rates, unemployment rates, literacy rates), and in virtually every other form of human organizational activity (e.g., censuses of the number of homeless people by non-profit organizations).

Data is measured, collected and reported, and analyzed, whereupon it can be visualized using graphs, images or other analysis tools. Data as a general concept refers to the fact that some existing information or knowledge is represented or coded in some form suitable for better usage or processing. Raw data ("unprocessed data") is a collection of numbers or characters before it has been "cleaned" and corrected by researchers. Raw data needs to be corrected to remove outliers or obvious instrument or data entry errors (e.g., a thermometer reading from an outdoor Arctic location recording a tropical temperature). Data processing commonly occurs by stages, and the "processed data" from one stage may be considered the "raw data" of the next stage. Field data is raw data that is collected in an uncontrolled "in situ" environment. Experimental data is data that is generated within the context of a scientific investigation by observation and recording. Data has been described as the new oil of the digital economy.

Data, information, knowledge and wisdom are closely related concepts, but each has its own role in relation to the other, and each term has its own meaning. According to a common view, data is collected and analyzed; data only becomes information suitable for making decisions once it has been analyzed in some fashion. One can say that the extent to which a set of data is informative to someone depends on the extent to which it is unexpected by that person. The amount of information content in a data stream may be characterized by its Shannon entropy.

Knowledge is the understanding based on extensive experience dealing with information on a subject. For example, the height of Mount Everest is generally considered data. The height can be measured precisely with an altimeter and entered into a database. This data may be included in a book along with other data on Mount Everest to describe the mountain in a manner useful for those who wish to make a decision about the best method to climb it. An understanding based on experience climbing mountains that could advise persons on the way to reach Mount Everest's peak may be seen as "knowledge". The practical climbing of Mount Everest's peak based on this knowledge may be seen as "wisdom". In other words, wisdom refers to the practical application of a person's knowledge in those circumstances where good may result. Thus wisdom complements and completes the series "data", "information" and "knowledge" of increasingly abstract concepts.

Data is often assumed to be the least abstract concept, information the next least, and knowledge the most abstract. In this view, data becomes information by interpretation; e.g., the height of Mount Everest is generally considered "data", a book on Mount Everest geological characteristics may be considered "information", and a climber's guidebook containing practical information on the best way to reach Mount Everest's peak may be considered "knowledge". "Information" bears a diversity of meanings that ranges from everyday usage to technical use. This view, however, has also been argued to reverse the way in which data emerges from information, and information from knowledge. Generally speaking, the concept of information is closely related to notions of constraint, communication, control, data, form, instruction, knowledge, meaning, mental stimulus, pattern, perception, and representation. Beynon-Davies uses the concept of a sign to differentiate between data and information; data is a series of symbols, while information occurs when the symbols are used to refer to something.

Before the development of computing devices and machines, people had to manually collect data and impose patterns on it. Since the

development of computing devices and machines, these devices can also collect data. In the 2010s, computers are widely used in many fields to collect data and sort or process it, in disciplines ranging from marketing, analysis of social services usage by citizens to scientific research. These patterns in data are seen as information which can be used to enhance knowledge. These patterns may be interpreted as "truth" (though "truth" can be a subjective concept), and may be authorized as aesthetic and ethical criteria in some disciplines or cultures. Events that leave behind perceivable physical or virtual remains can be traced back through data. Marks are no longer considered data once the link between the mark and observation is broken.

Mechanical computing devices are classified according to the means by which they represent data. An analog computer represents a datum as a voltage, distance, position, or other physical quantity. A digital computer represents a piece of data as a sequence of symbols drawn from a fixed alphabet. The most common digital computers use a binary alphabet, that is, an alphabet of two characters, typically denoted "0" and "1". More familiar representations, such as numbers or letters, are then constructed from the binary alphabet. Some special forms of data are distinguished. A computer program is a collection of data, which can be interpreted as instructions. Most computer languages make a distinction between programs and the other data on which programs operate, but in some languages, notably Lisp and similar languages, programs are essentially indistinguishable from other data. It is also useful to distinguish metadata, that is, a description of other data. A similar yet earlier term for metadata is "ancillary data." The prototypical example of metadata is the library catalog, which is a description of the contents of books.

---

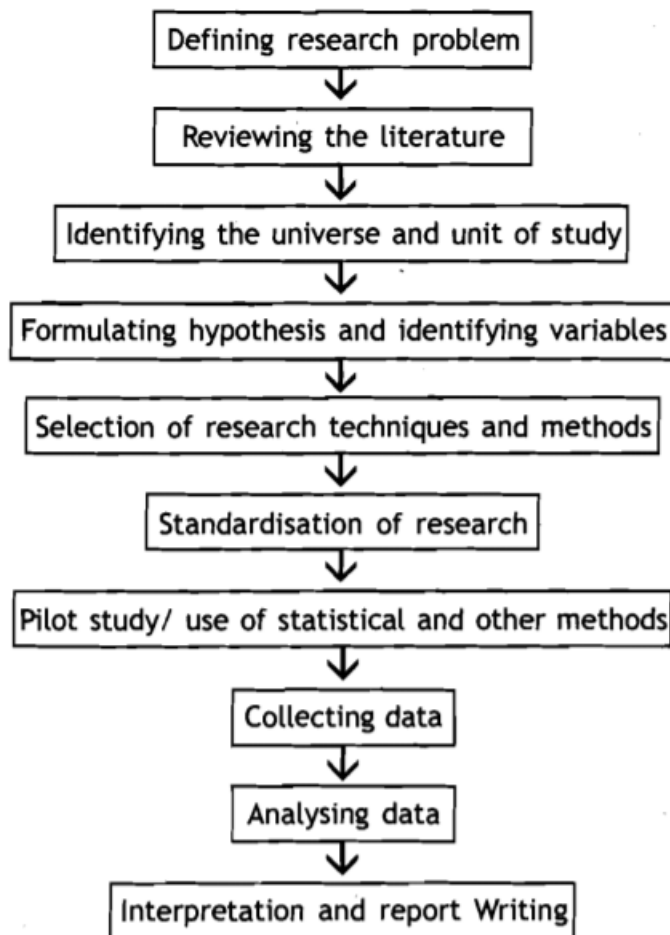
## **14.2 RESEARCH METHOD**

---

The structuring of the research process is an essential part of science. However this does not mean that these steps are always in a sequence. In fact, the various phases of research overlap. At times the first step determines the nature of the last step. The steps involved are not

mutually exclusive, nor are they separate and distinct. Figure 14.1 depicts the broad steps researchers usually take in the process.

Figure 14.2 Ten Steps in Research Design



Research design is the strategic plan of the project that sets out the broad structure of the research. See Box 14.1 for a brief mention of functions and purposes of preparing a research design.

#### Box 14.2 Function and Purpose of Research Design

Black and Champion (1976: 76-77) have pointed out the following three functions of research design.

- A research design provides a blueprint for operationalizing the research activity.
- It defines the limit and scope of the research.
- It provides an opportunity to the researcher to foresee possible areas of problems in the process carrying out the research.

Manheim (1977: 142) identified the following five purposes of preparing a research design.

- To gather sufficient evidence to support one's hypothesis and to disprove alternative hypotheses.
- To carry out a research that can be repeated in terms of its subject matter and research procedure. In other words, it does not pertain only to unique situation that has no relevance to society at large.
- To be able to work out correlations among variables in a manner that produces interrelated propositions.
- To make out the need for a pilot study in order to carry out the future plans of a full-fledged research-project.
- To be able to economise on time and resources by selecting appropriate techniques of data collection.

### **Defining your research problem**

In the research process, the first step is to select and clearly define the problem to be researched. You need to find the problem and formulate it so that it can be subjected to research. A research problem in general refers to some enquiry, which a researcher undertakes in the context of either a theoretical or practical situation or wants to obtain an explanation of the same. The formulation of a general topic into specific research problems constitutes the first step in a scientific inquiry. Essentially two steps are involved in formulating the research problems, that is, understanding the problems thoroughly and rephrasing it into meaningful terms from ' the analytical point of view. You need to select the subject that is familiar and feasible so that research material or sources of research are within your reach. It is better to select the research problem before a preliminary study of the existing literature. Formulating or defining a research problem is an important step in the research process and a clearly stated problem is research half done. You need to clearly state the research questions in the light of the topic of research and the theoretical foundations on which it rests. Next, you need to spell out the aims and objectives as per the requirements of your research questions.

This gives the research process a well-defined focus and direction. Unless one has a clear idea of the objectives, the course of the research will not be smooth and the data will not have the I desired consistency because it is possible for you to approach a topic from the viewpoint of different perspectives, each addressing a different I set of issues. For instance, research on the sociology of development can have many research questions like women's role in development, the role of caste and kinship in development, or social consequences of development on family and community life of people. While preparing a research design, demarcate the focus of research by jotting down the outline and features of the topic and the aims and objectives of &search.

### **Choice of field site(s)**

Embarking on a research, give as much emphasis to the area of research as the topic. To some extent, the choice of the area determines the success of your research. The relevance of a research topic depends on its usefulness to the problems of the area either in terms of a practical purpose or obtaining a theoretical understanding - of the epistemological issues. For instance, a study on communal relations cannot be carried out in a tribal village. In such a study, you would need to observe the interaction between different religious groups and therefore you would choose an area inhabited by people of several religious communities. It is desirable to have two or three sub-areas in mind within the broad area. For instance, sometimes you may encounter some unforeseen and unmanageable problems at the district level or the village level and then you would need to find an alternative to fall back upon. You should first spell out your choice of field site(s) and then start gathering information on it. This would help you gain an understanding of the geographical and socio-political conditions of the area, which would have a bearing on the collection of data. This would help you frame your research strategies and questions in a manner suited to the area and its people.

### **Consideration of time and resources**

## Notes

You need to be fully aware of the limits of your resources and also clearly define the time frame while designing your research. Unless you draw up a schedule of the different steps of your research, it is likely to become a long drawn process, which is bad for both quality and relevance. Imagine researches on cholera epidemics taking years to complete. The delay would mean poor quality research and an unchecked death rate. We also know that unless you get liberal time your research would fail because you cannot subject social reality to overnight machine tests in the laboratory to obtain quick results. You need to evaluate the time requirement in a realistic manner and plan the strategy accordingly. Careful planning and sticking to a time schedule will help you use your resources effectively and complete the research in time. Besides, you should be aware of the limitations of your resources and plan the strategy in a realistic and cost effective manner. If the resources are exhausted midway, it will be a severe blow to research. If an agency is funding the research, your credibility is at stake. Hence, you need to clearly state in your research design the time and resources that you have for research. You need also to foresee and note down the effects of your resource constraints on the research process. After identifying the affects you would develop strategies to counter those you can possibly do. The account of those, which cannot be managed, will help future researchers to be familiar.

### **Reviewing secondary material**

The purpose of reviewing the existing literature on your research theme is to help you assess the feasibility of the project but also to formulate an effective methodology. You would need to consult academic journals, conference proceedings, government reports, books, etc. (see Box 14.2).

Box 14.2 The Use of Computer in Literature Searching For the purpose of making use of the Internet facilities to search related literature, you need to read previous unit. If your research topic is precisely defined, the search can be a very fast and efficient way of obtaining relevant references in a number of bibliographic tools.



You may review two types of literature, literature concerning the concepts and theories, and the empirical literature consisting of studies made earlier. You may come across even such studies that contain both theoretical as well as substantive aspects of your research. The outcome of the review will be that you will know about available data and other materials on the theme of your research. A more sophisticated and clearer statement of specific research questions is likely to emerge after the literature review. When researchers prepare a research design they draw an outline of the entire research process. They need to have a clear picture of the nature of data that would help tackle the research questions. For instance, researchers decide in advance how many case studies would help them draw meaningful conclusions or the number of life histories that they need to collect and of which categories of persons. A lot of hard work and insightful thinking goes into the process. Researchers review the past studies on their topics and work upon their research questions to arrive at a realistic research design. Scheduling the time and events to observe in the field forms an important component of your research design. This provides you a sense of direction while collecting data. This does not imply that you have to strictly follow your schedule regardless of the situation in the field. The actual field conditions do guide you and correspondingly your research design may face unanticipated changes. Yet, you cannot just land up in the field unprepared and bewildered and hence you need to plan out the various stages and strategies of research. At the same time you have to be ready to make adjustments according to the field exigencies.

**Hypothesis**

After extensive literature survey, you need to state in clear terms the working hypothesis or hypotheses. The hypothesis is a tentative assumption made in order to test its logical or empirical consequences. You may define a hypothesis as a proposition or a set of propositions set forth as an explanation for the occurrence of some specified phenomena either asserted merely as a provisional conjecture to guide some

## Notes

investigation or accepted as highly probable in the light of established facts. A hypothesis may seem contrary to the real situation. It may prove to be correct or incorrect. In any event, it leads to an empirical test. Your hypothesis needs to be clear and precise and capable of being tested. It is to be limited in scope and consistent with known or established facts and should be amenable to testing within the stipulated time. It needs to explain what it claims to explain and should have empirical reference. A hypothesis may have variables and it may be looking for the nature of the relationship between the variables. The variables are empirical properties that take two or more values. For the purpose of research, you need to make a distinction between dependent and independent variables. The variables that you wish to explain are regarded as dependent variables (or criterion variables). The other variable expected to explain the change in the dependent variable is referred to as an independent variable (or predictor variable). The dependent variable is the expected outcome of the independent variable and independent variables produce dependent variables. Variables can have three types of relationships among them. A positive relationship is one where an increase in one variable leads to an increase in the other. A negative relationship is one where an increase in one variable leads to a decrease in the other. Finally, a zero relationship is one which shows no significant relationship between two variables. Such a distinction between dependent and independent variables is analytical and relates only to the research purpose. It needs to be mentioned that the formulation of hypothesis is not always a part of the research process. You may carry out exploratory research when you do not have sufficient knowledge of the situation to prepare a hypothesis.

### **Theoretical orientation**

Your research design needs to clearly spell out the data collection methods to be employed. Your methodological and philosophical orientations govern your choice of methods. Your research design would elucidate the methodological and theoretical basis of research and help you identify appropriate methods and techniques of data collection. For instance, if you have positivistic orientation, you would rely on

observational method because for you social reality would be an observable entity. On the other hand if you adopt a phenomenological model, you would employ various kinds of interviews to unravel the logico-mathematical model of culture. A researcher conforming to the post-modernist approach would view social reality as multidimensional and record multiple voices and interpretations. An action research with limited time resources would employ triangulation' (comprising multiple methods and multiple investigators), and focus group discussions. You have to carefully choose from the vast repertoire of sociological/anthropological methods, the ones that suit your research purpose most.

### **Universe and unit of study**

Before starting with data collection you have to identify the universe and the unit of study. The identification of universe implies demarcation of the physical area and social unit of study. The universe consists of the population within a well-defined area where' the study is to be conducted. However, such a group is usually too large and not possible to be covered by a single investigator. Therefore, a smaller and more manageable group may be selected by sampling. The outlines of the universe and its attributes may be delineated more clearly by a taking a census and then making the choice of the group(s) to work on. Within the broad universe further specification of the possible units that could be studied makes up the actual or effective universe. The group(s) selected as focus of study is called the unit of study.

**Pilot study** The pilot study is the leading study in your research area. The pilot study leads the researcher to the full-length investigation depending on the size of the population and the amount of time. In other words, a pilot study is an exploratory study done before the actual work starts in the field. It is a pre-testing of your research methods and techniques in order to perfect them. Pilot study will ensure that right questions have been put in the questionnaires for making the fieldwork fruitful. It makes you aware of the difficulties beforehand and provides you an opportunity of modifying your techniques to suit field conditions. Pilot study depends

## Notes

upon the size of the population, the time available and the availability of funds.

Sampling A universe is often too large for an individual to work upon. A sample is the smaller representation of a larger whole. Sampling allows the researcher to work scientifically and saves time. Analysing large quantities of material is wasteful and an intensive analysis of fewer cases is k economical. You need to be cautious and careful while sampling. As , explained earlier, the universe refers to a defined population size. Such a universe may be further divided depending on the specifications required. This is known as stature or subpopulation. A stature is a divisible ) category which depends upon the kind of problem in which one is interested. A sampling frame includes all the elements of a - population from which the sample is drawn. The determination of an error while / sampling, statistically or qualitatively is known as sampling error. The sample must be a true representative of the universe, as well as being adequate in size

Data collection After obtaining some idea and understanding your field and working out your methods and techniques of data collection, you may plan how to access the field. Quite often social research requires the study of the Survey Methods 'other' community and researchers need to make extensive preparations to gain entry into the society under study. One needs to plan who would facilitate access and how one would contact such persons. It is also possible to study one's own society. Some scholars (for example see Madan 1975) may not adhere to the idea of studying the 'other' community only. In Box 14.3 we bring you an excerpt from Madan (2004: 203)' who 'questioned the requirement of the personal study of an alien culture on the part of every anthropologist'.

### **Box 14.3**

On Studying One's Own Community Instead, I emphasised the importance of bridging the gap, or conversely, creating it, between the observer and the observed. I described fieldwork as the feat of 'living

intimately with strangers' (Madan 1975). I might have added: 'or strangely with intimates', which was what I had done during my fieldwork among the Pandits of rural Kashmir. The anthropologist studying his own culture I wrote, 'is an insider who takes up the posture of an outsider, by virtue of his training as-an anthropologist or a sociologist, and Looks at his own culture, hoping to be surprised. If he is, only then may he achieve new understandings' (Madan 1975: 149)

can use one or more methods to collect the data, taking into consideration the nature of investigation, objectives and scope of inquiry, financial resources and time available and degree of accuracy. The data to be collected would need to be adequate and dependable. Analysis and report writing After data collection, you would turn to their analysis. Analysis requires a number of closely related operations such as establishing categories and their application to raw data through coding, tabulation so that you can draw statistical inferences.. Tabulation is a part of the technical procedure wherein you are able to put your classified data in the form of tables (see Box 14.4).

Box 14.4 Classifying and Coding the Data Classifications facilitate rapid, accurate and comprehensive searches of stored field material, but a poor classification or careless retrieval may be worse than having none at all. In connection with this, particular attention should be paid to classifications which separate data which are otherwise related. For example, if "name-giving ceremonies" are indexed only under RITUAL, a search intended to assemble all data on KINSHIP may fall short of the mark. Notes must, in the first instance, be coded so that they can be subsequently located in a mass of material. You will probably wish to refer back to earlier notes quite frequently in the field, to check up on certain matters and test \$ informal hypotheses. At the very least, all sheets should be numbered sequentially.

You need to clearly delineate the form of analysis you wish to eventually adopt. Although often the nature of data collected by you determines the nature of analysis, yet at the stage of opting for certain methods of data

## Notes

collection you would have some idea of the analytical tools you are likely to employ. If you plan to adopt certain computer packages you would need to collect data keeping that in mind. While analysis may depend on the nature of data, you need to be careful to avoid the reverse situation, that is, the pre-determined mode of analysis solely determines the methods of data collection. You may face getting a one-sided picture of the social reality if you were to adopt computer-based methods only, because computer packages offer analysis of a particular dimension of reality while social research requires as broad a picture of reality as possible. You would be better off collecting data covering as many dimensions of reality as possible. In any case, you need to be quite clear about the mode of analysis to employ to interpret the data collected. Your research design is meant to reflect your theoretical orientation. In this way, you are actually planning every stage of research, from identifying the topic of research and method of data collection to report writing.

**Elements of Research Design** Your research design would be complete if you spell out the manner in Quantitative and Survey Methods which you would present the results of the research. It is an equally important step because you would need to keep in mind the ethics of representation, especially if the research deals with sensitive issues. While you seek to unravel social reality, you cannot play with the privacy of the people who are more than just the subjects of research. It is your responsibility to do justice both to the research and to the people. There is a practice of presenting data with pseudonyms and modification of identities, events and location. You need to always elucidate in your research design the manner in which you would report the results. Presentation of research findings for publication implies their distribution among the public, including those you studied. This is the point when you achieve the aim of making a contribution to the general body of literature related to the subject of your research.

---

## **14.3 PRIMARY AND SECONDARY DATA**

---

### **1. MEANING OF DATA:**

The word 'data' is Latin in origin, and literally, it means anything that is given. Different sources have defined the word in different ways. Webster's Third 28 29 Data: Definition, Types, Nature, Properties and Scope New International Dictionary defines data as "something given or admitted; facts or principles granted or presented; that upon which an inference or argument is based, or from which an ideal system of any sort is constructed". According to Oxford Encyclopaedic English Dictionary data are "known facts or things used as a basis for inference or reckoning". These dictionaries also state that even though data is the plural form of datum, it is often treated as a singular collective noun. Hence, its treatment as a singular noun is equally acceptable. For the sake of consistency, however, the word is used in this Unit as the plural form of datum. UNESCO defines data as 'facts, concepts or instructions in a formalised manner suitable for communication, interpretation or processing by human or automatic means'. Robert A. Arnold, in his 'Modern Data Processing' [Wiley, 1972], has defined the terms in the context of commerce as a function of business and accounting. Dictionary of Modern Economics defines data as "observations on the numerical magnitude of economic phenomena such as national income, unemployment, or the retail price". Data are defined in McGraw-Hill Encyclopaedia of Science and Technology as 'numerical or qualitative values derived from scientific experiments'. While another definition of data in Sciences is obtainable from CODATA (Committee on Data for Science and Technology) as quoted by Luedke and others in ARIST, 12, 119-181. CODATA defines data as a "crystallised presentation of the essence of scientific knowledge in the most accurate form". According to this definition, clarity and accuracy are two essential attributes of data. One also learns of yet another attribute of data from the CODATA definition. That is to say, data are the essence of the matter. The phrase 'essence of scientific knowledge' in this definition is synonymous with 'qualitative values derived from scientific experiments' as given in the McGraw-Hill definition. In social sciences, data are stated as values or facts, together with their accompanying study design, code books, research reports, etc. and are used by researchers for the purpose of secondary analysis. At one extreme, economics and demography have

## Notes

been heavily quantitative materials or observations. Sociology and, more recently, political science, fall between these two extremes. The change in research orientation in the subject can be seen with changing data, especially with data relating to public opinion. In humanities, the text such as Biblical materials or Shakespeare's drama deals with a fixed quantity of data represented by a finite amount of text to be interpreted. Clashing interpretations may be irresolvable, since each interpretation views the text differently, while the text to be interpreted may be finite and fixed. However, in sciences, the total text is to be interpreted and the text of data is not fixed before interpretation. The text of fact is constantly expanding. Scientists not only observe facts but also use instruments to generate more systematic data. In Information Science, Shuman [BASIS, 1975, 1(7), 11-12,34] defines data as "quantitative facts derived from experimentation, calculation, or direct observation". Shuman opines that a more meaningful definition of data is "the symbolisation of knowledge".

To understand further, we can say that data or facts have no shape that is relevant to a particular viewpoint. It must be given relevance, arrangement, coherence, usefulness within a definite framework of meaning, intent or interest.

### **TYPES OF DATA:**

In order to understand the nature of data it is necessary to categorise them into various types. Different categorisations of data are possible. The first such categorisation may be on the basis of disciplines, e.g., Sciences, Social Sciences, etc. in which they are generated. Within each of these fields, there may be several ways in which data can be categorised into types. For the sake of convenience we shall discuss the types as present in sciences and then in social sciences.

### **Types of Data in Social Sciences**



As in sciences, data in social sciences are also organized into different types so that their nature can be easily understood. The following categorization is normally observed in social sciences:

i) Data with reference to scale of measurement: Based on the scale of measurement, data can be categorized as follows:

a) Nominal data – The nominal scale is used for assigning numbers as the identification of individual unit. For example, the classification of journals according to the discipline they belong to, may be considered as nominal data. If numbers are assigned to describe the categories, the numbers represent only the name of the category.

b) Ordinal data – It indicates the ordered or graded relationship among the numbers assigned to the observations made. These numbers connote ranks of different categories having relationship in a definite order. For example, to study the responsiveness of library staff a researcher may assign ‘1’ to indicate poor, ‘2’ to indicate average, ‘3’ to indicate good and ‘4’ to indicate excellent. The numbers 1, 2, 3 and 4 in this case are set of ordinal data which indicate that 4 is better than 3 which in turn is better than 2 and so on. The ordinal data show the direction of the difference and not the exact amount of difference.

c) Interval data – Interval data are ordered categories of data and the differences between various categories are of equal measurement.

For example, we can measure the IQ (Intelligence Quotient) of a group of children. After assigning numerical value to the IQ of each child, the data can be grouped with interval of 10, like 0 to 10, 10 to 20, 20 to 30 and so on. In this case, ‘0’ does not mean the absence of intelligence and children with IQ ‘20’ are not doubly intelligent than children with IQ ‘10’.

d) Ratio data – Ratio data are the quantitative measurement of a variable in terms of magnitude. In ratio data, we can say that one thing is twice or

## Notes

thrice of another as for example, measurements involving weight, distance, price, etc.

ii) Data with reference to continuity: Data with reference to continuity can be categorised as follows:

a) Continuous data – Continuous data are an infinite set of possible values. Between a range there are infinite possible values. For example, height of an individual is not restricted to values like 155 cm. and after that to 156 cm. It can be 155.59 cm. or 155.99 cm. – continuous value.

b) Discrete data – The discrete data are finite or potentially countable set of values. For example, the number of members in a library. It can be 2,575 or 2,599 but certainly not  $2,599\frac{1}{2}$ . Similarly, the number of citizens in a country, the number of vehicles registered is the examples of discrete data.

iii) Data with reference to number of characteristics: Data can also be categorised on the basis of number of variables considered. These are:

a) Univariate data – Univariate data are obtained when one characteristic is used for observation, e.g., the performance of student in a given class.

b) Bivariate data – Bivariate data result when instead of one, two characteristics are measured simultaneously, e.g., height and weight of tenth class students.

c) Multivariate data – Multivariate data consist of observations on three or more characteristics, e.g., family size, income and savings in a metropolitan city in India.

iv) Data with reference to time: There are two types of data under this category. These are:

a) Time series data – Data recorded in a chronological order across time are referred to as time series data. It takes different values at different

times, e.g., the number of books added to a library in different years, monthly production of steel in a plant, yearly intake of students in a university.

b) Cross-sectional data – This refers to data for the same unit or for different units at a point of time, e.g., data across sections of people, region or segments of the society.

v) Data with reference to origin: Data under this category can be put as follows:

a) Continuous data – The data obtained first hand from individuals by direct observation, counting, and measurement or by interviews or mailing a questionnaire are called primary data. It may be complete enumeration or sampling, e.g., data collected from a market survey.

b) Secondary data – The data collected initially for the purpose and already published in books or reports but are used later on for some other purpose are referred to as secondary data. For example, data collected from census reports, books, data monographs, etc.

vi) Data with reference to characteristic: Data can be categorised on the basis of the characteristics as follows:

a) Quantitative data – When the characteristic of observation is quantified we get quantitative data. Quantitative data result from the measurement of the magnitude of the characteristic used. For example, age of a person, price of a commodity, income of a family, etc.

b) Qualitative data – When the characteristic of observation is a quality or attribute, we get qualitative data. For example, sex or colour of a person, or intelligence of a student.

### **PRIMARY DATA:**

## Notes

Primary data is data that is collected by a researcher from first-hand sources, using methods like surveys, interviews, or experiments. It is collected with the research project in mind, directly from primary sources.

The term is used in contrast with the term secondary data. Secondary data is data gathered from studies, surveys, or experiments that have been run by other people or for other research.

Typically, a researcher will begin a project by working with secondary data. This allows time to formulate questions and gain an understanding of the issues being dealt with before the more costly and time consuming operation of collecting primary data.

Primary Data:

These are the data which are collected from some primary sources i.e., a source of origin where the data generate.

These are collected for the first time by an investigator or an agency for any statistical analysis.

“Data which are gathered originally for a certain purpose are known as primary data.” — Horace Secrist

Merits:

1. It has high degree of accuracy.
2. For some enquiries, secondary data is not available.
3. These are more reliable.
4. It needs no extra precautions.

Demerits:

1. It requires lot of time.
2. It needs much money.
3. These data can be obtained through skilled persons only.
4. Sometimes, these data are not available altogether.

### **SECONDARY DATA:**

Secondary data refers to data that is collected by someone other than the user. Common sources of secondary data for social science include censuses, information collected by government departments, organizational records and data that was originally collected for other research purposes. Primary data, by contrast, are collected by the investigator conducting the research.

Secondary data analysis can save time that would otherwise be spent collecting data and, particularly in the case of quantitative data, can provide larger and higher-quality databases that would be unfeasible for any individual researcher to collect on their own. In addition, analysts of social and economic change consider secondary data essential, since it is impossible to conduct a new survey that can adequately capture past change and/or developments. However, secondary data analysis can be less useful in marketing research, as data may be outdated or inaccurate

These are the data which are collected from some secondary source i.e. the source of reservation storage where the data is collected by one person and used by other agency. These are collected as primary data and used by other as secondary data.

## Notes

“The data which are used in an investigation, but which have been gathered originally by someone else for some other purpose are known as secondary data.” — Blair

Merits:

1. It is easy to collect.
2. Time and money is saved.
3. Sometimes primary data cannot be obtained.
4. Some data are more reliable than primary.

Demerits:

1. These are not reliable as primary data.
2. Extra caution is needed to use these data.
3. All types of data are not available.
4. Purpose of original collection may have been different.

### **Sources of secondary data**

Secondary data can be obtained from different sources:

information collected through censuses or government departments like housing, social security, electoral statistics, tax records

internet searches or libraries

Gps, remote sensing

km progress reports

Administrative data and census

Government departments and agencies routinely collect information when registering people or carrying out transactions, or for record keeping – usually when delivering a service. This information is called administrative data.

It can include:

personal information such as names, dates of birth, addresses

information about schools and educational achievements

information about health

information about criminal convictions or prison sentences

tax records, such as income

A census is the procedure of systematically acquiring and recording information about the members of a given population. It is a regularly occurring and official count of a particular population. It is a type of administrative data, but it is collected for the purpose of research at specific intervals. Most administrative data is collected continuously and for the purpose of delivering a service to the people.

Advantages and disadvantages of secondary data

Secondary data is available from other sources and may already have been used in previous research, making it easier to carry out further research. It is time-saving and cost-efficient: the data was collected by someone other than the researcher. Administrative data and census data may cover both larger and much smaller samples of the population in detail. Information collected by the government will also cover parts of the population that may be less likely to respond to the census (in countries where this is optional).

A clear benefit of using secondary data is that much of the background work needed has already been carried out, such as literature reviews or case studies. The data may have been used in published texts and statistics elsewhere, and the data could already be promoted in the media or bring in useful personal contacts. Secondary data generally have a pre-

## Notes

established degree of validity and reliability which need not be re-examined by the researcher who is re-using such data. Secondary data is key in the concept of data enrichment, which is where datasets from secondary sources are connected to a research dataset to improve its precision by adding key attributes and values.

Secondary data can provide a baseline for primary research to compare the collected primary data results to and it can also be helpful in research design.

However, secondary data can present problems, too. The data may be out of date or inaccurate. If using data collected for different research purposes, it may not cover those samples of the population researchers want to examine, or not in sufficient detail. Administrative data, which is not originally collected for research, may not be available in the usual research formats or may be difficult to get access to.

Secondary analysis or re-use of qualitative data

While 'secondary data' is associated with quantitative databases, analysis focused on verbal or visual materials created for another purpose, is a legitimate avenue for the qualitative researcher. Actually one could go as far as claim that qualitative secondary data analysis “can be understood, not so much as the analysis of pre-existing data; rather as involving a process of re-contextualizing, and re-constructing, data.

### Check Your Progress 2

Notes: a) Space is given below for writing your answers.

b) Compare your answers with those given at the end of the unit.

1. Discuss the Research method.

.....  
.....  
.....



- .....
- .....
2. Describe the Primary and secondary data.
- .....
- .....
- .....
- .....
- .....

---

## 14.4 STYLE AND REFERENCE

---

Research reports present references and bibliography. A bibliography is a list of published works, although by common usage both published and unpublished materials are listed in a bibliography. Many researchers use these two terms references and bibliography interchangeably, but the two terms have definite meanings. A bibliography is a list of titles - books, research reports, articles, papers etc. that may or may not have been referred to in the text of the research report. References include only such studies, books, articles or papers that have been actually referred to in the text of the research report. In short REFERENCES consists of all documents, including journal articles, books, chapters, technical reports, computer programmes and unpublished works that are mentioned in the text of the manuscript. A bibliography contains everything that would be in the reference section plus other publications which were consulted by the researcher but were not cited in the manuscript. After having clarity about references and bibliography, let us understand the need and importance of referencing and footnotes. Articles, papers, books, research reports (Dissertations/thesis) monographs etc. quoted inside the text of the report should find a place in the reference section. In the text of the report, the author's surname along with the year of publication is given e.g. (Glatthorn, 1998). When few sentences are quoted from a source, the page number too is noted, e.g. (Glatthorn 1998 :137-138). Full length reference be placed at the end of the chapter or at the end of the thesis/report or at the foot of that page as footnote. The traditional style of giving references is to place them as the footnotes on the relevant page(s). The footnotes are serialized inside the text and in the footnotes

## Notes

of each chapter. In some cases footnotes are generally avoided, instead full reference is given at the end of the report. Footnotes and reference perform many functions. As the name implies, footnotes are usually found at the foot of a page, although in some manuscripts they appear at the end of each chapter or at the end of a paper. Footnotes and references are used to ;

- i) Validate a point, statement or argument. The original source or authority is acknowledged through the use of a footnote or reference.
- ii) provide the reader with sufficient information to enable him/her to consult the sources independently.
- iii) provide cross-references to material appearing in other parts of the report.
- iv) explain, supplement or amplify material that is included in the main body of chapter paper
- v) acknowledge a direct quotations. Thus, it is very clear that researchers acknowledge their indebtedness to other authors not only as a matter of courtesy but also as means of confirming their work. By now you might have understood the concept of footnotes and references along with their importance. Now, let us see how to use footnotes and references in the report.

### VARIOUS STYLES OF REFERENCING

There are mainly two style manuals used for referencing. These are: American Psychological Association, Publication Manual, 3rd edition. Washington, DC : American Psychological Association, (1983). The Chicago Manual of style, 13<sup>th</sup> revised edition, Chicago University of Chicago Press, 1982. Generally, references are arranged in alphabetical order where the researcher has cited the name of the author and the year of publication of the work in the text. Another practice followed is references are arranged in a sequence as they appear in the text of the research report. Here related statement in the body of the text is numbered. : However, most research reports use alphabetical listing of

references. Now, let us see how to use footnotes:- \* Footnotes are always double-spaced between each other, though each footnote is typed single-spaced. It is usual to give the full name of the author in its normal order, i.e. first name (or initial) and second name precede surname. e.g. 6 John, W. Best. (1993). Research in Education. New Delhi : Prentice Hall of India, P. 148 here '6' indicates number given in the text, "John" is first name, "W" is second name and Best is surname and P. 148 indicates that matter or direct sentence or quotation is taken from that page. Ibid in the footnote refers to the same work and the reference that precedes it. Here the succeeding references to a work immediately follow the first full citation. Ibid in latin means the same.e.g. 6 John, W. Best. (1993). Research in Education. New Delhi : Prentice Hall of India. P. 148 7 Ibid. P.148 (This indicates the same work and the same page as above i.e. '6' here). 8 Ibid, p. 149 (This indicates the same work as above but a different page) Op. cit :-Op.cit. in Latin means the work cited. It is used in a footnote to the same work as a preceding but not immediately preceding reference, so here another reference to the same work is made but not consecutively.

For example.

Allan, A Glatthorn (1998): Writing the Winning Dissertations: A StepbyStep Guide. California: Corwin Press Inc. P.189. 6. Fred, N. Kerlinger. (1 973): Foundations of Behavioural Research. NewYork: Holt, R. Inehart & Winston. P. 259. 7. Glatthorn, op.cit. P. 191. Here reference 7 refers to the same reference as 5 except the pages differ in the two cases Loc. Cit. Loc. Cit, is used when reference is made to the same page as a preceding but not immediately preceding reference, the last name of the author and phrase loc. Cit. are used. e.g. 8. Kerlinger, loc, cit. here this refers to same work as in '6' on the same 5 8 Page.

A number of other abbreviations appear in research reports. While writing a research report, abbreviations, may be used to condense space in references or footnotes. If a researcher is not familiar, shehe should consult the relevant literature as and when required. In the following

table, a comprehensive list of abbreviations has been given for ready reference.

Thus, we have seen how to use footnotes/references in the report. Here our discussion is limited to only references/footnotes. Note the following points while using footnotes. Having adopted a method of footnoting is consistent throughout the whole report. Footnotes should be concise, but clarity and readability should not be sacrificed for brevity. All footnotes regardless of length are terminated by a full stop. The same bottom margin should be maintained on each page of the typescript, regardless of the number of footnotes.

---

## **14.5 RESEARCH REPORT**

---

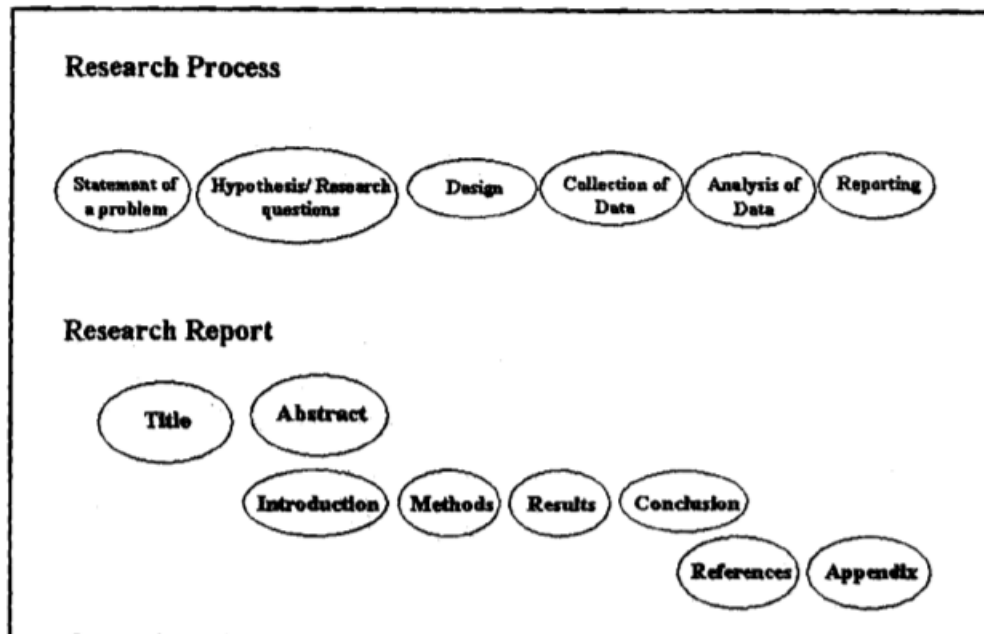
Writing a research report is a valuable experience for a researcher. It is an essential part of the research process. Most research reports are either in the form of research articles or abstracts or thesis and dissertations or project reports. These reports are the vehicle for researchers to communicate the results of an investigation to others across space and time. The research journal articles, master's degree dissertations doctoral thesis and project reports, all have the common objective -to disseminate research results and findings, ideas and information. There are, of course, other ways of communicating research results, may be through oral presentation in a seminar or conference or on-line journals in a website. Reporting research findings and results are of paramount importance in all areas of research. Because, it is hardly worth doing research if it is not disseminated. The purpose of writing a research report is to communicate the ideas and information with other people.

So, communication of research results should take place through research reports with a number of different audiences in mind: fellow researchers, peers, practitioners, teachers, curriculum planners and developers and the general public.

In this unit, we will discuss the meaning of the research report, how to prepare a research report and its various components. Besides this, the

significance of a research report and different types forms of research report will also be discussed. You will gain a deeper understanding about the format of a research article, an abstract, a thesis and dissertation and a project report. It is hoped that these formats would help you in writing articles for research journals and to prepare a complete report after conducting a research work project.

### Box 1: Research Process and Research Report



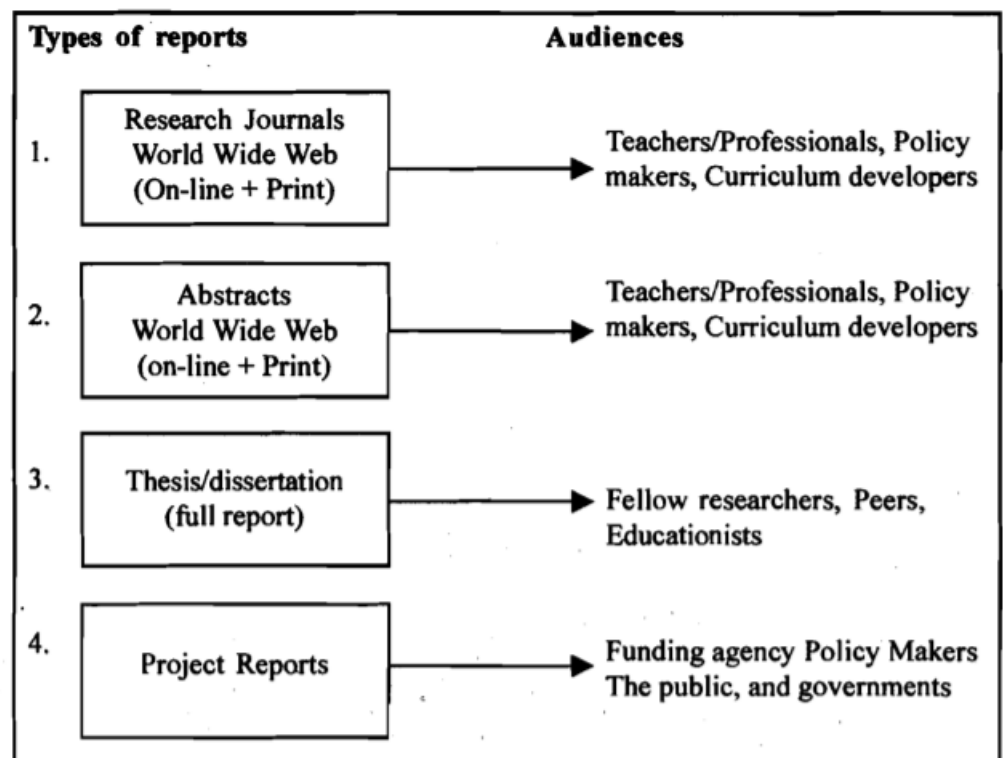
### TYPES OF RESEARCH REPORTS

Researchers, disseminate their research findings through research articles prepared for print and on-line journals, conference papers, theses and dissertation, and project reports. Research reports are usually written for a special group of audience such as one's fellow researchers or peer groups, practitioners, policy makers, curriculum planners and developers, teachers or the general public. Therefore research reports are prepared keeping in view different purposes and different audiences (see Box 2) and also of different length for different audiences. For example, a research study on learning environments in schools their impact on academic. This could be disseminated via an academic journal article focusing on the theory emerging from the research, another journal article

## Notes

concentrating on the pupils' perspectives including case study material. A student of a post graduate programme with a title "A study of learning environments for primary schools teachers trained through distance education in Eastern Ethiopia" could prepare a thesis consisting of quantitative and qualitative data, a discussion on methodology and conclusions. A concise summary or an abstract of an article or a report could be prepared by a researcher so that the audience/reader can learn the rationale behind the study methodology, pertinent results and important recommendation.

### Box 2: Disseminating research findings for different audiences



There are debates and discussions pertaining to different format and style of writing a research article, an abstract, a thesis and dissertation and a project report. Before drafting a research report, you must ask yourself who will read it. It is useful to think about a variety of reports that results from a research. This section discusses different varieties or types of reports. They are:

- i) Research articles
- ii) Abstracts

- iii) Theses and dissertations
- iv) Project reports

### 14.5.1 Research Articles

The purpose of research articles is to inform the readers of what you investigated, why and how you conduct the investigation, the results and conclusions. These articles are usually reports of empirical studies, review articles or theoretical articles.

**Reports of empirical studies** These are the reports of original research. They consist of distinct sections that reflect the stages in the research process and appears in the sequence of four stages (APA, 1983 pg. 21). They are: i) Introduction: development of the problem under investigation and statement of the purpose of the investigation, ii) Method: description of the method used to conduct the investigation. iii) Results: report of the results that were found, and iv) Discussion: interpretation and discussion of the implications of the results.

**Review articles** Review articles are critical evaluations of research material that has already been published. The author of a review article considers the progress of current research toward clarifying a problem by organizing, integrating and evaluating previous published material (APA, 1983). A review article appears in the sequence of the following stages. Defines and clarifies the problem; Summarises previous investigations in order to inform the reader about the current research; o Identifies relations, contradictions, gaps and inconsistencies in the literature; and 0 Suggests the steps in solving the problem.

**Theoretical articles** The author presents empirical information incorporating the theoretical issues of a problem. Here the author finds out the development of theory in order to expand and refine theoretical constructs (APA, 1983). The sections of the theoretical articles are

## Notes

usually arranged by relationship rather than by chronology. The sections or the stages of writing this type of report is like review articles. Sometimes, you may find brief research reports, comments and replies, discussions of different types of methods (qualitative/ethnographic studies), case histories and monographs published in journal articles. These are also reported according to the guidelines discussed for journal articles.

Various writing style for preparing research articles for different journals are described in manuals by Campbell, Ballou and Slade (1982), Turabian (1973) and Modern Language Association (MLA) of America (Gibaldi and Achtext, 1988) and 32 manual of American Psychological Association (1983)

### 14.5.2 Abstracts

An abstract is a comprehensive summary of the contents of the article or a thesis/dissertation submitted for evaluation. It allows the audience or readers to go through the contents of a journal article or a research report quickly. These abstracts serve as one of the most useful reference guides to the researcher and keep him/her abreast of the work being done in his/her own field and also in the related fields (Koul, 1986 p. 94). These abstracts are published in journals and educational periodicals. For example, i) Psychological Abstracts (1927 - Onwards) ii) Education Abstracts (1949 - Onwards) iii) Sociology of Education Abstracts (1965 - Onwards) iv) Dissertation Abstracts International (1952 - Onwards) In other words, an abstract is a summarized form of a research report (within one or two paragraphs of about 150-200 words). It includes the problem hypothesis or research questions, procedures, principal results and conclusions of a research work.

When to write an abstract? A researcher writes an abstract when his/her dissertation for a higher degree is accepted as a part of dissertation and thesis. For example, preparing an abstract of a doctoral (Ph.D) thesis while submitting an article for publication in a journal. when a research



report is presented to an audience. An abstract communicates the scope of a research article. It also presents the summarized version of a topic to be discussed by the readers/audience. It facilitates academic discussion pertaining to a specific research problem. It helps the researchers to identify the issues while going through the abstract relevant to their research from the published articles.

In other words, a good abstract summarises the key information from every major section in the body of the report. It provides the key issues and conclusions from the report precisely. The characteristics of a good abstract are: Accurate Self contained Concise and specific Coherent and readable

An abstract should have the following main sections.

- i) Introduction - purpose of the study, research problem and hypotheses/research questions
- ii) Main Body - brief description of the methodology
- iii) Results - specific data considered for analysis
- iv) Conclusion - important conclusion or recommendation of the research study.

### 14.5.3 Thesis and Dissertation

A Thesis or a dissertation is a record of research activities. It is usually produced in partial fulfillment of the requirements of a course/programme or for an advanced degree. It involves presenting a research problem with an argument or point of Restarch Report: Various view. The methods or procedures adopted are substantiated with reasoned argument Components md Structure and evidence. This is written to share the issues and concerns related to a specific research problem with fellow researchers supported by discussions besides presenting the outwmedfindings. This record is submitted to an institutionJexamining committee for awarding degrees to the student. The reading audiences are committee members, fellow researchers, peer

group, teachers. These reports in the form of theses and dissertations are usually preserved by the universities that award the authors their doctoral and masters degrees. Sometimes these research studies are published in whole or in part in various educational periodicals or journals. Because the reports of many research studies are never published, a check of the annual list of theses and dissertations issued by various agencies is necessary for a thorough coverage of the research literature (Koul, 1986 p. 96). The following discussion describes the sections of a typical thesis or a dissertation.

### **Differences between a thesis or a dissertation and a research journal**

The major difference between a thesis or dissertation and a research article is the length of the document and the contents covered. For example, researchers who publish articles are limited by the established publishing criteria of a particular journal. Suppose a research article of six or eight pages, as prescribed by a specific journal, cannot include all the information contained in a 150-200 page thesis or dissertation. The author of a thesis or a dissertation produces a "final" manuscript; but the author of a journal article produces a "copy" manuscript. The requirements of a thesis and a dissertation are not necessarily identical to the requirements of manuscripts submitted for publication of a journal (APA, 1983 p. 189). The manuscripts of research articles are read by editors, reviewers and compositors for publication. They must conform to the format and other policies of the journal to which they are submitted. The theses and dissertations reach their audiences in the exact form in which they are prepared. They have been prepared for a research-productive career. These theses and dissertations are submitted to the institutions/examining committee as a part of a course/programme. Therefore, they must satisfy the specific requirements prepared by an institution. Sometimes, the requirements/style mentioned by standard manuals may be or may not be considered. Universities/institutions/schools who have launched a course/programme should provide written guidelines and a format which explain all modification to APA style. The thesis or dissertation in its

original I form is not acceptable to journals but the condensed versions of doctoral dissertations I may appear as journal articles.

### 14.5.4 Project Reports

In the light of the varied types and purposes of projects, the format of a project report will depend upon the level at and audience for which the research is done. For example, the academic research project for a degree and projects funded by funding agencies like UNICEF, World Bank or UNESCO differ in their formats. I Public and private educational funding agencies sponsor research projects either to, ' an individual or to a team or group of researchers through an institution. These agencies require researchers applying for financial help to carry out a project, to submit a research proposal at the outset and a project report at the end/or after completion of a project within a specified time.-The final report (a large scale or a small scale) of a research project hnded by an agency is a written document that the researcher sends to the funding agency. It may take the form (greatly reduced in content and length) of an article in a professional journal. The organization of khe , content and structure of a project report and academic theses might look alike. Ruurcb Reports and Appliutioas These research reports may vary in length. While preparing a project report one should bear in mind the audience for the report. For example, scientific or general report is prepared as per the theme and audience of a project work. Possible formats for a project report are as follows:

Example 1 Executive Summary - a synopsis of the research focusing on its practical implications Aims and objectives - as specified by the funder or researcher. Context - a discussion of the organization and its work and the reasons for undertaking the project work Results - an account or description of what the research project discovered. Recommendations - a list of actions to be implemented.

Example 2 Titlepage List of contents Tables and figures Project objectives Methods, procedures used for collection of data Budget.

**Check Your Progress 2**

Notes: a) Space is given below for writing your answers.  
b) Compare your answers with those given at the end of the unit.

1. Discuss the Style and reference.

.....  
.....  
.....

2. Write about the steps of Research report.

.....  
.....  
.....

---

**14.6 LET US SUM UP**

---

Unit 14 has enumerated the various steps for undertaking a sociological research with the aim of preparing you to carry out one such research. This will be a practical exercise for you to COIII>'-' AS a compulsory requirement of completing MSO 002. You will need to prepare a research design before actually carrying out your mini research project. The Reflection and Action exercise for Unit 14 is that you prepare a research design for your proposed research. You may of course modify it as you come to learn in more detail about the various steps needed in the research process.

In this unit, we have discussed the need for chapterisation and its functions. Each chapter in the report has its own functions. We have also seen that there is no uniformity in the scheme of chapterisation. Diversity in chapterisation mainly depends on, (i) field of study, (ii) nature of research area, (iii) requirement of a department or agency. We have also seen some of the cases where diversity enters. Footnotes in the report, though it is a traditional concept, has many functions in the research report. It is always advisable to maintain consistency in writi.~g footnotes. We have discysed how to make use of footnotes in the report.

---

## 14.7 KEY WORDS

---

**References:** Reference is a relationship between objects in which one object designates, or acts as a means by which to connect to or link to, another object. The first object in this relation is said to refer to the second object. It is called a name for the second object

**Data:** Data are individual units of information. A datum describes a single quality or quantity of some object or phenomenon. In analytical processes, data are represented by variables. Although the terms "data", "information" and "knowledge" are often used interchangeably, each of these terms has a distinct meaning.

---

## 14.8 QUESTIONS FOR REVIEW

---

1. Discuss the Research method
2. Describe the Primary and secondary data
3. Discuss the Style and reference
4. Write about the steps of Research report

---

## 14.9 SUGGESTED READINGS AND REFERENCES

---

- Singleton, Jr Royce A. and Bruce C. Straits 1999. Approaches to Social Research. Oxford University Press: New York
- Sarantakos, 5.1998 (first published in 1993). Social Research. Macmillan: London
- Driscoll & Brizee. What is Primary Research? Purdue Online Write Lab. Retrieved from <https://owl.english.purdue.edu/owl/resource/559/01/> on June 24th, 2017
- BYU FHSS Research Support Center. Data Types and Sources. Retrieved from <https://fhssrsc.byu.edu/Pages/Data.aspx> on June 24th, 2018.

## Notes

- Schutt, R. Investigating the Social World. Sage Publications, 2006. p423-426,412-416
- McCaston, M. Katherine. Tips for Collecting, Reviewing, and Analyzing Secondary Data. Partnership & Household Livelihood Security Unit(PHLS), February 1998. <https://web.archive.org/web/20070709112209/http://www.livelihoods.org/info/pcdl/docs/work/SL%20Nepal/Reference%20Sheets/Tips%20for%20Using%20Secondary%20Data.doc>
- 696 Research Methods, Secondary Data Analysis <http://www.csulb.edu/~msaintg/ppa696/696scond.htm>
- Sundararajan, V. Ethnicity, discrimination and health outcomes: a secondary analysis of hospital data from Victoria, Australia. Diversity in Health and Social Care, 2007.
- Banta, J.E. Substance Abuse and Dependence Treatment in Outpatient Physician Offices, 1997-2004. American Journal of Drug & Alcohol Abuse.vol 33.aug 2007. p583-593.
- Mochmann, Ekkehard. Data Archiving and the Uses of Secondary Analysis. Central Archives for Empirical Social Research, University of Cologne. [https://web.archive.org/web/20070612083623/http://www.metadater.org/archiving\\_and\\_secondary\\_analysis.htm](https://web.archive.org/web/20070612083623/http://www.metadater.org/archiving_and_secondary_analysis.htm)
- O'Sullivan, E. & Rassel, G. R.. Research Methods for Public Administrators. 3rd Ed. Longman,1999. p265,268-269.
- Kelly, M. Primary and Secondary Data. McKinnon Secondary College, 2005. <http://www.mckinnonsc.vic.edu.au/vceit/infodata/primarysecondary.htm>
- Corti, L. & Bishop, L. (2005) 'Strategies in Teaching Secondary Analysis of Qualitative Data' FQS 6(1) <http://www.qualitative-research.net/index.php/fqs/article/view/509>

---

## 14.10 ANSWERS TO CHECK YOUR PROGRESS

---

### Check Your Progress 1

4. See Section 13.2
5. See Section 13.3
6. See Section 13.4

**Check Your Progress 2**

4. See Section 13.5
5. See Section 13.6