

# Design and analysis of sample surveys

School of Science and Informatics

Department of Mathematics, Statistics and Physical Sciences

Taita Taveta University

[www.ttu.ac.ke](http://www.ttu.ac.ke)

Dr. Noah Mutai

January 23, 2023

## Contents

<b>1</b>	<b>Fundamentals</b>	<b>4</b>
1.1	Introduction . . . . .	4
1.2	Types of Sampling . . . . .	5
1.3	Properties of random sampling . . . . .	5
1.4	Properties of estimators . . . . .	6
1.5	Principal steps involved in planning and execution of a sample survey. . . . .	7
1.6	Advantages of Sampling . . . . .	9
1.7	Exercises . . . . .	9
<b>2</b>	<b>Simple Random Sampling(SRS)</b>	<b>11</b>
2.1	Introduction and Description . . . . .	11
2.2	How to draw a simple random sample . . . . .	11
2.2.1	Lottery method . . . . .	11
2.2.2	Random number tables . . . . .	12
2.2.3	How to use a random number table. . . . .	12
2.2.4	Computer software . . . . .	13
2.3	Simple random sampling with replacement (SRSWR) . . . . .	13
2.3.1	Definition and Estimation of Population Mean, Variance and Total . . . . .	13
2.4	Simple random sampling without replacement . . . . .	15
2.4.1	Definition and estimation of population mean, variance and total . . . . .	16
2.5	Confidence intervals for population mean $\bar{Y}$ and Total $Y'$ . . . . .	19
2.6	Sampling for proportions and percentages . . . . .	20
2.6.1	Estimation of population proportion . . . . .	21
2.6.2	Estimation of population total or total number of count . . . . .	22
2.6.3	Confidence Interval estimation for P . . . . .	22
2.7	Determination of sample sizes . . . . .	23
2.8	Exercises . . . . .	24
<b>3</b>	<b>Stratified Random Sampling</b>	<b>27</b>
3.1	Introduction and Description . . . . .	27
3.1.1	Estimation of Population Mean, Variance and Total . . . . .	28
3.2	Estimation of Variance . . . . .	29
3.3	Allocation problem and choice of sample sizes is different strata . . . . .	30
3.3.1	Equal allocation . . . . .	30
3.3.2	Proportional Allocations . . . . .	30

3.4	Allocation of Sample Sizes(Neymann Allocation) . . . . .	30
3.4.1	Variances under different allocations . . . . .	33
3.4.2	Comparison of variances of sample mean under SRS with stratified mean under proportional and optimal allocation: . . . . .	34
3.5	Estimate of variance and confidence intervals . . . . .	35
3.6	Exercises . . . . .	37
<b>4</b>	<b>Systematic sampling</b>	<b>39</b>
4.1	Introduction and Description . . . . .	39
4.1.1	Advantages of systematic sampling: . . . . .	40
4.2	Estimation of Population Mean, Variance and Total . . . . .	40
4.2.1	Estimation of population mean . . . . .	40
4.3	Exercises . . . . .	43
<b>5</b>	<b>Cluster sampling</b>	<b>45</b>
5.1	Introduction and Description . . . . .	45
5.2	Estimation of Population Mean, Variance and Total . . . . .	46
5.3	Exercises . . . . .	47
<b>6</b>	<b>Ratio and Regression Estimation</b>	<b>48</b>
6.1	Ratio Estimation . . . . .	48
6.2	Bias and mean squared error of ratio estimator . . . . .	49
6.3	Regression Estimation . . . . .	50
6.3.1	Estimate of variance . . . . .	51
6.4	Regression estimates when $\beta$ is computed from sample . . . . .	52
6.4.1	Bias of $\hat{Y}_{reg}$ . . . . .	52
6.5	Exercises . . . . .	53
<b>7</b>	<b>Double sampling (Two phase sampling)</b>	<b>54</b>
7.1	Introduction and description . . . . .	54
7.2	Double sampling in ratio method of estimation . . . . .	55
7.3	Exercises . . . . .	56
<b>8</b>	<b>Varying probability sampling(pps)</b>	<b>57</b>
8.1	Introduction and description . . . . .	57
8.2	PPS sampling with replacement (WR) . . . . .	58
8.2.1	Estimation of Population Mean, Variance and Total . . . . .	58
8.3	Exercises . . . . .	59

<b>9 Two Stage Sampling(Sub-sampling)</b>	<b>60</b>
9.1 Introduction and Description . . . . .	60
9.2 Estimation of population mean . . . . .	62
9.3 Estimate of variance . . . . .	63
<b>10 Sources of Errors in Surveys</b>	<b>64</b>
10.1 Introduction . . . . .	64
10.2 Non-Sampling Errors . . . . .	64
10.3 Sampling errors . . . . .	66
<b>11 Organisation of National surveys, and the Kenya Bureau of Statistics(K.N.B.S)</b>	<b>67</b>
<b>12 Bibliography</b>	<b>68</b>

### Terms

- Definition — an explanation of the mathematical meaning of a word.
- Theorem — A statement that has been proven to be true.
- Proposition — A less important but nonetheless interesting true statement.
- Lemma — A true statement used in proving other true statements (that is, a less important theorem that is helpful in the proof of other results).
- Corollary — A true statement that is a simple deduction from a theorem or proposition.
- Proof — The explanation of why a statement is true.
- Conjecture — A statement believed to be true, but for which we have no proof. (a statement that is being proposed to be a true statement).
- Axiom — A basic assumption about a mathematical situation (a statement we assume to be true).

# 1 Fundamentals

## 1.1 Introduction

**Sample survey, finite population sampling or survey sampling** is a method of drawing an inference about the characteristics of a population or universe by **observing under only a part of the population**. Such methods are extensively used by government bodies throughout the world for assessing, among others, different characteristics of national economy as a required for making decisions and for the planning and projection of future economic structure.

Ideally, total information about the population is obtained through census, where every individual in the population is involved in giving out information. However, most of the times due to certain constraints to be discussed later, it is not always possible to carry out a census.

In a sample survey the purpose of the survey statistician is to estimate some functions of the population parameter,  $\theta(y)$ , say, by choosing a sample(part of the population) and by observing the values of  $y$  only on units selected in the sample. The statistician therefore want to make an inference about the population by observing only a part of it. This is essential and perhaps the only practical method of inference about the characteristics of the population since in many socioeconomic investigations the survey population may be very large, containing say hundreds or thousands of units.

**Definition 1.1.** Survey population — A finite(survey) population is a collection of known number  $N$  of identifiable units labeled  $1, 2, 3, \dots, i, \dots, N$  where  $i$  stands for the label as well as the physical unit labeled  $i$ . The number  $N$  is the size of the population. The parametric functions of general interest for estimation are;

1. Population total,  $Y = \sum_{i=1}^N Y_i$
2. Population mean:  $\bar{Y} = \frac{Y}{N} = \frac{1}{N} \sum_{i=1}^N Y_i$
3. Population variance:  $S_Y^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2$
4. Population coefficient of variance:  $C_Y = \frac{S_Y}{\bar{Y}}$ , where  $S_Y$  is the population variance and  $\bar{Y}$  is the population mean.

**Definition 1.2.** Sample — is a part of the population/subset of the population selected for study. A sample may be drawn from a population either under with replacement(wr) or under without replacement(wor).

After a sample is selected, data are collected from the sampled units. We shall denote by  $y_i$  the value of  $y$  on the unit selected at the  $i^{th}$  draw ( $i = 1, 2, \dots, n$ ). Thus for example if the sample is  $S = \{2, 3, 2\}$ ,  $y_1 = Y_2, y_2 = Y_3, y_3 = Y_2$ . Clearly  $y_i$  is a random variable whose possible values lie in the set  $\{Y_1, Y_2, \dots, Y_N\}$

For a sample  $s$ , we shall denote some statistics as follows;

1. Sample total,  $y = \sum_{i=1}^n y_i$
2. Sample mean,  $\bar{y} = \frac{y}{n} = \frac{1}{n} \sum_{i=1}^n y_i$
3. Sample variance,  $s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$
4. Sample coefficient of variation,  $c_y = \frac{s_y}{\bar{y}}$ , where  $s_y$  is the sample variance and  $\bar{y}$  is the sample mean.

**Definition 1.3.** Sampling units: This refers to the individual items whose characteristics are to be measured in the sample survey.

**Definition 1.4.** Sampling frame: This is the list of all sampling units. It may be a list of units with identification and particulars or a map showing the boundaries of sampling units e.g. a manufacturing firm may want to determine how popular a newly manufactured product is within the community suggests a possible frame for the survey. The firm may decide to concentrate its surveys in urban residential areas only. In this case, you have a complete list of estates in urban areas. The residents in those chosen estates will be interviewed and inferences are made.

**Definition 1.5.** Sampled population: It is the set of individuals in the sampling frame. Its actually the subset of the target population. Note: Sampled population is not necessarily the same as target population.

**Definition 1.6.** Sampling scheme: Its the technique by which the elements which constitute the sampled are obtained from the population.

## 1.2 Types of Sampling

1. Haphazard sampling: No scheme has been used at all it is neither probability nor non-probability sampling.
2. Purposive or judgemental sampling or non probability sampling.
3. Probability/random sampling- statistical theory is used and the kind of inferences made are based on statistical procedures. There is some element of chance associated with selection of items into the sample.

## 1.3 Properties of random sampling

We are able to define the set of distinct samples,  $S_1, S_2, \dots, S_N$ , which the procedure is capable of selecting if applied to a specific population. This means that we can say precisely what sampling units belong to  $S_1$  to  $S_2$  and so on.

1. Each possible sample  $S_i$  has assigned to it a known probability of selection  $\pi_i$ .

2. We select one of the  $S_i$  by a process in which each  $S_i$  receives its appropriate probability  $\pi_i$ , of being selected.
3. The method for computing the estimate from the sample must be stated and must lead to a unique estimate for any specific sample.

The simplest type of sampling is SRS(simple random sampling). We shall also make use of common terms in statistics like; statistic, estimator, point estimation and interval estimation, regression estimation. Design and analysis of sample survey is about knowing those estimators and design procedures that are good.

### 1.4 Properties of estimators

- Precision — how much variation there is in the estimation from sample to sample.
- Trueness — on average how close is the estimate to the population characteristics being estimated.
- Accuracy — combination of precision and trueness. Precision of estimates will be measured by variance e.g. if the estimator is  $X$  then;

$$Var(x) = \sigma_x^2 = E[X - \mu_x]^2 \quad (1.1)$$

where  $\mu_x = E[X]$

Trueness of an estimation will be measured by the bias, which is defined as a difference between the expectation of estimate and population parameter of which is an estimate.

$$Bias(x) = E[X] - \mu_x \quad (1.2)$$

If  $E[X] - \mu_x = 0$ , then it is an unbiased estimator. Accuracy is measured using mean squared error(MSE)

$$MSE(x) = \sigma_x^2 + (Bias(x))^2 \quad (1.3)$$

Generally estimators which have low variance (high precision) and low bias are preferred.

### 1.5 Principal steps involved in planning and execution of a sample survey.

The broad steps to conduct any sample surveys are as follows:

1. **Objective of the survey:** The objective of the survey has to be clearly defined and well understood by the person planning to conduct it. It is expected from the statistician to be well versed with the issues to be addressed in consultation with the person who wants to get the survey conducted. In complex surveys, sometimes the objective is forgotten and data is collected on those issues which are far away from the objectives.
2. **Population to be sampled:** Based on the objectives of the survey, decide the population from which the information can be obtained. For example, population of farmers is to be sampled for an agricultural survey whereas the population of patients has to be sampled for determining the medical facilities in a hospital.
3. **Data to be collected:** It is important to decide that which data is relevant for fulfilling the objectives of the survey and to note that no essential data is omitted. Sometimes, too many questions are asked and some of their outcomes are never utilized. This lowers the quality of the responses and in turn results in lower efficiency in the statistical inferences.
4. **Degree of precision required:** The results of any sample survey are always subjected to some uncertainty. Such uncertainty can be reduced by taking larger samples or using superior instruments. This involves more cost and more time. So it is very important to decide about the required degree of precision in the data. This needs to be conveyed to the surveyor also.
5. **Method of measurement:** The choice of measuring instrument and the method to measure the data from the population needs to be specified clearly. For example, the data has to be collected through interview, questionnaire, personal visit, combination of any of these approaches, etc. The forms in which the data is to be recorded so that the data can be transferred to mechanical equipment for easily creating the data summary etc. is also needed to be prepared accordingly.
6. **The frame:** The sampling frame has to be clearly specified. The population is divided into sampling units such that the units cover the whole population and every sampling unit is tagged with identification. The list of all sampling units is called the frame. The frame must cover the whole population and the units must not overlap each other in the sense that every element in the population must belong to one and only one unit. For example, the sampling unit can be an individual member in the family or the whole family.
7. **Selection of sample:** The size of the sample needs to be specified for the given sampling plan. This helps in determining and comparing the relative cost and time of different sampling plans. The method and plan adopted for drawing a representative sample should also be detailed.



8. **The Pre-test:** It is advised to try the questionnaire and field methods on a small scale. This may reveal some troubles and problems beforehand which the surveyor may face in the field in large scale surveys.
9. **Organization of the field work:** How to conduct the survey, how to handle business administrative issues, providing proper training to surveyors, procedures, plans for handling the non-response and missing observations etc. are some of the issues which need to be addressed for organizing the survey work in the fields. The procedure for early checking of the quality of return should be prescribed. It should be clarified how to handle the situation when the respondent is not available.
10. **Summary and analysis of data:** It is to be noted that based on the objectives of the data, the suitable statistical tool is decided which can answer the relevant questions. In order to use the statistical tool, a valid data set is required and this dictates the choice of responses to be obtained for the questions in the questionnaire, e.g., the data has to be qualitative, quantitative, nominal, ordinal etc. After getting the completed questionnaire back, it needs to be edited to amend the recording errors and delete the erroneous data. The tabulating procedures, methods of estimation and tolerable amount of error in the estimation needs to be decided before the start of survey. Different methods of estimation may be available to get the answer of the same query from the same data set. So the data needs to be collected which is compatible with the chosen estimation procedure.
11. **Information gained for future surveys:** The completed surveys work as guide for improved sample surveys in future. Beside this they also supply various types of prior information required to use various statistical tools, e.g., mean, variance, nature of variability, cost involved etc. Any completed sample survey acts as a potential guide for the surveys to be conducted in the future. It is generally seen that the things always do not go in the same way in any complex survey as planned earlier. Such precautions and alerts help in avoiding the mistakes in the execution of future surveys.
12. **Pilot Survey** In planning a survey efficiently, some prior information about the population under consideration and the operational and cost aspects of data collection will be needed. When such information is not available

## 1.6 Advantages of Sampling

Sample surveys have potential advantages over complete enumeration(census). They include;

1. **Reduced cost** — If data are secured from only a small fraction of the aggregate, expenditures may be expected to be smaller than if a complete census is attempted
2. **Greater speed** — For the same reason, the data can be collected and summarized more quickly with a sample than with a complete count. This may be a vital consideration when the information is urgently needed.
3. **Greater scope** — In certain types of inquiry, highly trained personnel or specialized equipment, limited in availability, must be used to obtain the data. A complete census may then be impracticable: the choice lies between obtaining the information by sampling or not at. Thus surveys which rely on sampling have more scope and flexibility as to the types of information that can be obtained.
4. **Greater accuracy** — Because personnel of higher quality can be employed and can be given intensive training, a sample may actually produce more accurate results than the kind of complete enumeration that it is feasible to take.
5. **Risk** — When a survey involves risky tests such as testing a new drug, sampling should be used.

## 1.7 Exercises

1. Discuss the statement: "The need to collect statistical information arises in almost every conceivable sphere of human activity."
2. Describe briefly each of the following terms:
  - (a) Primary data
  - (b) Secondary data
  - (c) Mail inquiry
  - (d) Questionnaire/schedule
  - (e) Population
  - (f) Census
  - (g) Element
  - (h) Sample
  - (i) Sampling unit

(j) Sampling frame

3. Differentiate between target and sampled population. What problem arises if two populations are not same?
4. What is the primary advantage of probability sampling over the non probability sampling? Cite three situations where non probability sampling is to be preferred
5. Assume a sample survey shall be carried out to find out about how satisfied students are with their faculty.
  - (a) How would you define the population?
  - (b) Would you consider a census of all students or rather a sample survey? (Why?)
  - (c) How would you operationalise? being satisfied with their faculty?
  - (d) What is a sampling frame and how could one be obtained in the example?
  - (e) How could a random sample be obtained?
  - (f) How do you consider the idea of obtaining a sample from alumni?

## 2 Simple Random Sampling(SRS)

### 2.1 Introduction and Description

We shall consider various sampling procedures (schemes) for selection of units in the sample. Since the objective of a survey is to make inferences about the population, a procedure that provides a precise estimator of the parameter of interest is desirable. Many sampling schemes have been developed to achieve this objective. To begin with, simple random sampling, the simplest and the most basic sample selection procedure, is discussed.

**Definition 2.1. Simple random sampling** — The sampling procedure is known as simple random sampling if every population unit has the same chance of being selected in the sample. The sample thus obtained is termed a simple random sample.

For selecting a simple random sample in practice, units from population are drawn one by one. If the unit selected at any particular draw is replaced back in the population before the next unit is drawn, the procedure is called with replacement (WR) sampling. A set of units selected at  $n$  such draws, constitutes a simple random with replacement sample of size  $n$ . In such a selection procedure, there is a possibility of one or more population units getting selected more than once. In case, this procedure is continued till  $n$  distinct units are selected, and all repetitions are ignored, it is called simple random sampling (SRS) without replacement (WOR). This method is equivalent to the procedure, where the selected units at each draw are not replaced back in the population before executing the next draw.

Another definition of simple random sampling, both with and without replacement, could be given on the basis of probabilities associated with all possible samples that can be selected from the population.

**Definition 2.2.** Simple random sampling is the method of selecting the units from the population where all possible samples are equally likely to get selected.

### 2.2 How to draw a simple random sample

The most commonly used procedures for selecting a simple random sample are:

1. Lottery method
2. Random number tables
3. Computer software

#### 2.2.1 Lottery method

In this method, each unit of the population of  $N$  units is assigned a distinct identification mark (number) from 1 to  $N$ . This constitutes the population frame. Each of these numbers is then written on a different

slip of paper. All the  $N$  slips of paper are identical in respect of size, color, shape, etc. Fold all these slips in an identical manner and put them in a container or drum, in which a thorough mixing of the slips is carried out before each blindfold draw. The paper slips are then drawn one by one. The units corresponding to the identification labels on the selected slips, are taken to be members of the sample.

### 2.2.2 Random number tables

A random number table is an arrangement of ten digits from 0 to 9, occurring with equal frequencies independently of each other and without any consistently recurring trends or patterns. Several standard tables of random numbers prepared by Tippet (1927), Fisher and Yates (1938), Kendall and Smith (1939), Rand Corporation (1955), and Rao et al. (1974) are available.

Direct Approach. Again, the first step in the method is to assign serial numbers 1 to  $N$  to the  $N$  population units. If the population size  $N$  is made up of  $K$  digits, then consider  $K$  digit random numbers, either row wise or column wise, in the random number table. The sample of required size is then selected by drawing, one by one, random numbers from 1 to  $N$ , and including the units bearing these serial numbers in the sample.

This procedure may involve number of rejections of random numbers, since zero and all the numbers greater than  $N$  appearing in the table are not considered for selection. The use of random numbers has, therefore, to be modified. Two of the commonly used modified procedures are:

### 2.2.3 How to use a random number table.

Part of a Table of Random Numbers			
61424	20419	86546	00517
90222	27993	04952	66762
50349	71146	97668	86523
85676	10005	08216	25906
02429	19761	15370	43882
90519	61988	40164	15815
20631	88967	19660	89624
89990	78733	16447	27932

Figure 1

1. Let's assume that we have a population of 185 students and each student has been assigned a number from 1 to 185. Suppose we wish to sample 5 students (although we would normally sample

more, we will use 5 for this example).

2. Since we have a population of 185 and 185 is a three digit number, we need to use the first three digits of the numbers listed on the chart.
3. We close our eyes and randomly point to a spot on the chart. For this example, we will assume that we selected 20631 in the first column.
4. We interpret that number as 206 (first three digits). Since we don't have a member of our population with that number, we go down to the next number 899 (89990). Once again we don't have someone with that number, so we continue at the top of the next column. As we work down the column, we find that the first number to match our population is 100 (actually 10005 on the chart). Student number 100 would be in our sample. Continuing down the chart, we see that the other four subjects in our sample would be students 049, 082, 153, and 164.
5. Researchers use different techniques with these tables. Some researchers read across the table using given sets (in our examples three digit sets). For our class, we will use the technique I have described.

#### 2.2.4 Computer software

In practice, the lottery method of selecting a random sample can be quite burdensome if done by hand. Typically, the population being studied is large and choosing a random sample by hand would be very time-consuming. Instead, there are several computer programs that can assign numbers and select  $n$  random numbers quickly and easily. Many can be found online for free.

### 2.3 Simple random sampling with replacement (SRSWR)

#### 2.3.1 Definition and Estimation of Population Mean, Variance and Total

A sample is said to be selected by simple random sampling with replacement(srswr) by  $n$  draws from a population of size  $N$  if the sample is drawn by observing the following rule;

1. At each draw, each unit in the population has the same chance of being selected.
2. A unit selected at a draw is returned to the population before the next draw.

The same unit, therefore might be selected more than once. Thus the probability of getting a sample(sequence),  $i = 1, 2, \dots, i_n$  is;

$$P(\{i = 1, 2, \dots, i_n\}) = \frac{1}{N}, \dots, \frac{1}{N} = \frac{1}{N^n} \quad (2.1)$$

There are  $N^n$  possible samples(sequences) in the sample space  $S$ , for a given  $(N, n)$ . A *srswr* of  $n$  draws from a population of size  $N$  will be denoted by  $srswr(N, n)$ .

**Theorem 2.1.** In  $srswr(N, n)$  sample mean  $\bar{y}$  is an unbiased estimator of the population mean  $\bar{Y}$ .

*Proof.*

$$E(\bar{y}) = E\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = y_1 \quad (2.2)$$

Since  $y_1, y_2, \dots, y_n$  are independently and identically distributed (iid) random variables with;

$$P(y_i = Y_k) = \frac{1}{N}, k = 1, 2, \dots, N, i = 1, 2, \dots, n. \quad (2.3)$$

Now,  $E(y_1) = \frac{1}{N} \sum_{k=1}^N Y_k = \bar{Y}$ . Hence,  $E(\bar{y}) = \bar{Y}$ .

Alternatively, let  $t_i$  be the number of times  $i$  occurs in the sample. Therefore,  $t_i$  follows a multinomial distribution with  $E(t_i) = \frac{n}{N}$ ,  $Var(t_i) = \frac{n}{N} \left(1 - \left(\frac{1}{N}\right)\right)$ ,  $Cov(t_i, t_j) = \frac{-n}{N^2} (i \neq j = 1, 2, \dots, N)$

(Show this)

Now,  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^N t_i Y_i \Rightarrow E(\bar{y}) = E\left(\frac{1}{n} \sum_{i=1}^N t_i Y_i\right) = \frac{1}{n} \sum_{i=1}^N Y_i E(t_i)$ . But  $E(t_i) = \frac{n}{N}$ ,

$\Rightarrow$

$$E(\bar{y}) = \frac{1}{n} \frac{n}{N} \sum_{i=1}^N Y_i = \bar{Y} \quad (2.4)$$

Hence  $E(\bar{y}) = \bar{Y}$ . □

**Corollary 2.1.1.** In  $srswr(N, n)$  and unbiased estimator of  $Y$  is  $\hat{Y} = NV ar(\bar{y})$ . Prove this.

**Theorem 2.2.** In  $srswr(N, n)$ , the sample variance is given by

$$Var(\bar{y}) = \frac{\sigma^2}{n}, \sigma^2 = \frac{1}{N} \sum_{i=1}^N (Y_i - \bar{Y})^2 \quad (2.5)$$

*Proof.*  $Var(\bar{y}) = Var\left(\frac{1}{n} \sum_{i=1}^N t_i Y_i\right)$

$$= \frac{1}{n^2} \sum_{i=1}^N Y_i^2 Var(t_i) + \frac{1}{n^2} \sum_{i \neq j} Y_i Y_j Cov(t_i, t_j)$$

$$\text{but } Var(t_i) = \frac{n}{N} \left(1 - \left(\frac{1}{N}\right)\right)$$

$$\text{and } Cov(t_i, t_j) = \frac{-n}{N^2} (i \neq j = 1, 2, \dots, N)$$

$$\begin{aligned} \Rightarrow Var(\bar{y}) &= \frac{1}{Nn} \left(1 - \left(\frac{1}{N}\right)\right) \sum_{i=1}^N Y_i^2 - \frac{1}{N^2 n} \left[ \left(\sum_{i=1}^N Y_i\right)^2 - \sum_{i=1}^N Y_i^2 \right] \\ &= \frac{N}{N^2 n} \sum_{i=1}^N Y_i^2 - \frac{1}{N^2 n} \sum_{i=1}^N Y_i^2 + \frac{1}{N^2 n} \sum_{i=1}^N Y_i^2 - \frac{1}{N^2 n} \sum_{i=1}^N \left(\sum_{i=1}^N Y_i\right)^2 \\ &= \frac{1}{Nn} \sum_{i=1}^N Y_i^2 - \frac{1}{N^2 n} \left(\sum_{i=1}^N Y_i\right)^2 = \frac{1}{Nn} \sum_{i=1}^N Y_i^2 - \frac{1}{N^2 n} (N^2 \bar{Y}^2) \\ &= \frac{1}{Nn} \sum_{i=1}^N Y_i^2 - \frac{1}{n} (\bar{Y}^2) = \frac{1}{n} \left[ \frac{1}{N} \sum_{i=1}^N Y_i^2 - \bar{Y}^2 \right] \end{aligned}$$

$$= \frac{\sigma^2}{n} \quad (2.6)$$

Note:

$$(x_1 + x_2)^2 = x_1^2 + 2x_1x_2 + x_2^2 \Rightarrow \left(\sum_{i=1}^N x_i\right)^2 = \sum_{i=1}^N x_i^2 + \sum \sum_{i \neq j} x_i x_j$$

□

**Corollary 2.2.1.** In  $srswr(N, n)$  the sample variance is given by;

$$Var(\bar{y}) = \frac{N-1}{Nn} S_y^2; S_y^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2 \quad (2.7)$$

**Corollary 2.2.2.** In  $srswr(N, n)$ ,  $Var(\hat{Y}) = \frac{N^2 \sigma^2}{n}$ . As  $n$  increases,  $Var(\bar{y})$  decreases, even if  $n = N$ ,  $Var(\bar{y})$  does not vanish. Also in  $srswr$ ,  $n$  may be arbitrarily large.

## 2.4 Simple random sampling without replacement

A sample of size  $n$  is said to be selected by simple random sampling without replacement (srswor) if the selection procedure is such that every possible sequence(sample) has the same chance of being selected. Sampling design is achieved by drawing a sample by the following draw-by-draw procedure;

1. At each draw each available unit in the population has the same chance of being selected.
2. A unit selected at a draw is removed from the population before the next draw.

If the population is of size  $N$  and we require a simple random sample without replacement of size  $n$ , then this is chosen at random from  $\binom{N}{n}$  distinct sample. Each of the  $\binom{N}{n}$  samples has the same probability  $\frac{1}{\binom{N}{n}}$  or  $\binom{N}{n}^{-1}$  of being selected.

**Lemma 2.3.** For a  $srswor(N, n)$  design the probability of a specified unit being selected at any given draw is  $\frac{1}{N}$  i.e.

$$P_r(i_k) = \frac{1}{N}, r = 1, 2, \dots, n. \quad (2.8)$$

for any given  $i_k$

**Lemma 2.4.** For a  $srswor(N, n)$  the probability of two specified units being selected at any two given draws is  $\frac{1}{N} \left(\frac{1}{N-1}\right)$ , i.e.

$$P_{r,s}(i_r, i_s) = \frac{1}{N(N-1)}, r < s, r = 1, 2, \dots, n. \quad (2.9)$$



for any given  $i_r \neq i_s$

**Lemma 2.5.** For a  $srswor(N, n)$  the probability that a specified unit is included in the sample is  $\frac{n}{N}$  i.e.

$$P(i \in s) = \pi_i \text{ (say)}, i = 1, 2, \dots, N \quad (2.10)$$

**Lemma 2.6.** For a  $srswor(N, n)$  the probability that any two specified units are included in the sample is  $\frac{n(n-1)}{N(N-1)}$  i.e.

$$P(i \in s, j \in s) = \pi_{i,j} \text{ (say)}, i \neq j, i = 1, 2, \dots, N \quad (2.11)$$

The quantities  $\pi_i$  and  $\pi_{i,j}$  (as defined in 2.5 and 2.6) are respectively the inclusion probabilities of units  $i$  and  $(i, j)$  in the sample. These are called respectively, the first order and second order inclusion probabilities of a design.

#### 2.4.1 Definition and estimation of population mean, variance and total

We consider the problem of estimating,  $\bar{Y}$ ,  $Y$  and  $S^2$  in  $srswor$ . Consider a population of size  $N$  and let  $n$  be the size of the simple random sample drawn from this population without replacement. Now let  $a_i$  equals 1 if the  $i^{th}$  unit is selected and 0 elsewhere,  $i = 1, 2, \dots, N$

Then  $a_i$  is a random variate such that;  $E(a_i) = 1 \times \text{probability of } i^{th} \text{ selected unit.}$

$= 1 \times \frac{n}{N} = \frac{n}{N}$  inclusion probability.

$E(a_i, a_j) = 1 \times \text{probability of the } i^{th} \text{ and } j^{th} \text{ unit selected.}$

$= 1 \times \frac{n}{N} \times \frac{n-1}{N-1} = \frac{n(n-1)}{N(N-1)}$

Therefore the sample total is

$$y = \sum_{i=1}^N a_i Y_i = \sum_{i=1}^n y_i \quad (2.12)$$

The sample mean is given as;

$$\bar{y} = \frac{1}{n} \sum_{i=1}^N a_i Y_i = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y} \quad (2.13)$$

**Theorem 2.7.** In  $srswor(N, n)$  the sample mean  $\bar{y}$  is an unbiased estimator of the population mean  $\bar{Y}$

*Proof.*  $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{n} \sum_{i=1}^N a_i Y_i$

$E(\bar{y}) = E\left(\frac{1}{n} \sum_{i=1}^N a_i Y_i\right) = \frac{1}{n} \sum_{i=1}^N Y_i E(a_i)$

$= \frac{1}{n} \sum_{i=1}^N Y_i \cdot \frac{n}{N} = \frac{1}{N} \sum_{i=1}^N Y_i = \bar{Y}$  □

**Corollary 2.7.1.** For  $srswor(N, n)$ ,  $\hat{Y} = N\bar{y}$  is an unbiased estimator of the population total  $Y$ .

**Theorem 2.8.** In  $srswor(N, n)$ ,  $Var(\bar{y}) = \frac{N-n}{Nn} S^2$ , where  $S^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2$

*Proof.* From  $Var(y) = E(y^2) - (E(y))^2$

it implies that  $Var(\bar{y}) = E(\bar{y}^2) - (E(\bar{y}))^2$

But  $E(\bar{y}) = E\left(\frac{1}{n} \sum_{i=1}^n y_i\right) = E\left(\frac{1}{n} \sum_{i=1}^N a_i Y_i\right) = \frac{1}{n} \sum_{i=1}^N Y_i E(a_i) = \frac{1}{N} \sum_{i=1}^N Y_i$

Next,  $E(\bar{y}^2) = E\left(\frac{1}{n} \sum_{i=1}^N a_i Y_i\right)^2$

$= E\left(\frac{1}{n^2} \sum_{i=1}^N a_i Y_i^2 + \frac{1}{n^2} \sum \sum_{i \neq j} a_i a_j Y_i Y_j\right)$

(Factor in expectation and use the fact that  $E(a_i) = \frac{n}{N}$

and  $E(a_i, a_j) = \frac{n(n-1)}{N(N-1)}$ )

$= \frac{1}{nN} \sum_{i=1}^N Y_i^2 + \frac{n-1}{Nn(N-1)} \sum \sum_{i \neq j} Y_i Y_j$

But  $\sum \sum_{i \neq j} Y_i Y_j = \left(\sum_{i=1}^N Y_i\right)^2 - \sum_{i=1}^N Y_i^2$

Therefore,  $E(\bar{y}^2) = \frac{1}{Nn} \sum_{i=1}^N Y_i^2 + \frac{n(n-1)}{Nn(N-1)} \left[\left(\sum_{i=1}^N Y_i\right)^2 - \sum_{i=1}^N Y_i^2\right]$

$= \left[\frac{1}{nN} - \frac{n-1}{nN(N-1)}\right] \sum_{i=1}^N Y_i^2 + \frac{n-1}{Nn(N-1)} \left(\sum_{i=1}^N Y_i\right)^2$

Now  $(x_1 + x_2)^2 = x_1^2 + 2x_1x_2 + x_2^2 \Rightarrow \left(\sum_{i=1}^N x_i\right)^2 = \sum_{i=1}^N x_i^2 + \sum \sum_{i \neq j} x_i x_j$

Therefore,  $Var(\bar{y}) = \frac{N-n}{Nn(N-1)} \sum_{i=1}^N Y_i^2 + \frac{n-1}{Nn(N-1)} \left(\sum_{i=1}^N Y_i\right)^2 - \left(\frac{1}{N} \sum_{i=1}^N Y_i\right)^2$

$= \frac{N-n}{Nn(N-1)} \sum_{i=1}^N Y_i^2 + \left[\frac{n-1}{Nn(N-1)} - \frac{1}{N^2}\right] \left(\sum_{i=1}^N Y_i\right)^2$

$= \frac{N-n}{Nn(N-1)} \sum_{i=1}^N Y_i^2 - \frac{N-n}{N^2n(N-1)} \left(\sum_{i=1}^N Y_i\right)^2$

$= \frac{N-n}{Nn(N-1)} \left[\sum_{i=1}^N Y_i^2 - N\bar{Y}^2\right]$

But  $S^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2 = \frac{1}{N-1} \left(\sum_{i=1}^N Y_i^2 - N\bar{Y}^2\right)$ .

Therefore;

$$Var(\bar{y}) = \frac{N-n}{Nn} S^2 \quad (2.14)$$

on simplification.

□

**Theorem 2.9.** In  $srswor(N, n)$  an unbiased estimator of  $Var(\bar{y})$  is  $\frac{N-n}{Nn} s^2$  where  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$ . Proof this.

**Theorem 2.10.** In  $srswor(N, n)$  the sample variance is an unbiased estimator of the population variance i.e.

$E(s^2) = S^2$  where  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$  and  $S^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2$

*Proof.*  $s^2 = \frac{1}{n-1} (\sum_{i=1}^n y_i^2 - n\bar{y}^2) = \frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - \frac{1}{n} (\sum_{i=1}^n y_i)^2\right]$

$= \frac{1}{n-1} \left[\sum_{i=1}^n y_i^2 - \frac{1}{n} \left(\sum_{i=1}^n y_i^2 + \sum \sum_{i \neq j} y_i y_j\right)\right]$

$= \frac{1}{n-1} \left[\left(1 - \frac{1}{n}\right) \sum_{i=1}^n y_i^2 - \frac{1}{n} \sum \sum_{i \neq j} y_i y_j\right]$

(open brackets and simplify)

$= \frac{1}{n} \sum_{i=1}^n y_i^2 - \frac{1}{n(n-1)} \sum \sum_{i \neq j} y_i y_j$

Taking expectations on both sides we have,

$E(s^2) = \frac{1}{n} E\left(\sum_{i=1}^n y_i^2\right) - \frac{1}{n(n-1)} E\left(\sum \sum_{i \neq j} y_i y_j\right)$

but

$$E \left( \sum_{i=1}^n y_i^2 \right) = E \left( \sum_{i=1}^N a_i Y_i^2 \right) = \frac{n}{N} \sum_{i=1}^N Y_i^2$$

$$\text{since } E(a_i) = \frac{n}{N}$$

$$\text{and } E(a_i a_j) = \frac{n(n-1)}{N(N-1)}$$

and

$$E \left( \sum \sum_{i \neq j} y_i y_j \right) = E \left( \sum \sum_{i \neq j} a_i a_j Y_i Y_j \right) = \frac{n}{N} \frac{n-1}{N-1} \sum \sum_{i \neq j} Y_i Y_j.$$

Therefore;

$$\begin{aligned} E(s^2) &= \frac{1}{n} \frac{n}{N} \sum_{i=1}^N Y_i^2 - \frac{1}{n(n-1)} \frac{n}{N} \frac{n-1}{N-1} \sum \sum_{i \neq j} Y_i Y_j \\ &= \frac{1}{N} \sum_{i=1}^N Y_i^2 - \frac{1}{N(N-1)} \sum \sum_{i \neq j} Y_i Y_j \\ &= \frac{1}{N} \sum_{i=1}^N Y_i^2 - \frac{1}{N(N-1)} \left( \sum_{i=1}^N Y_i \right)^2 + \frac{1}{N(N-1)} \sum_{i=1}^N Y_i^2 \\ &= \left[ \frac{1}{N} + \frac{1}{N(N-1)} \right] \sum_{i=1}^N Y_i^2 - \frac{1}{N(N-1)} \left( \sum_{i=1}^N Y_i \right)^2 \\ &= \frac{1}{N-1} \left[ \sum_{i=1}^N Y_i^2 - N \bar{Y}^2 \right] = S^2 \end{aligned}$$

Hence;

$$E(s^2) = S^2 \quad (2.15)$$

□

**Corollary 2.10.1.** For  $srs_{wor}(N, n)$ , an unbiased variance estimator of  $Y$  is  $Var(\hat{Y}) = \frac{N(N-n)}{n} s^2$

$$\text{Proof. } Var(\hat{Y}) = Var(N\bar{y}) = N^2 Var(\bar{y})$$

$$= N^2 \frac{N-n}{Nn} s^2$$

$$= \frac{N(N-n)}{n} s^2$$

which completes the proof.

□

**Corollary 2.10.2.** An estimator of error of  $\bar{y}$  is  $\hat{\sigma}(\bar{y}) = \sqrt{\frac{N-n}{Nn}} s$ . An estimator of the coefficient of variation is  $c(\bar{y}) = \sqrt{\frac{N-n}{Nn}} \frac{s}{\bar{y}}$ .

$c(\bar{y})$  is a ratio estimator and biased estimator of  $C(\bar{y})$ .

NOTE: The sample mean in  $srs_{wor}(N, n)$  is a better estimator of  $\bar{Y}$  (in the small variance sense) than sample mean in  $srs_{wr}(N, n)$ .

$$\text{Proof. } Var(\bar{y}|srs_{wr}) - Var(\bar{y}|srs_{wor}) = \frac{n-1}{Nn} S^2 > 0 \text{ for } n > 1. \quad \square$$

In sampling from an infinite population (where each  $Y_i$  is an independently and identically distributed random variable) with variance of each random variable as  $\sigma^2$ ,  $Var(\bar{y}) = \frac{\sigma^2}{n}$ . In simple random sampling with replacement, draws may be made an infinite number of times and  $Var(\bar{y}) = \frac{\sigma^2}{n}$ . In simple random sampling without replacement, however,  $Var(\bar{y}) = \left[ 1 - \left( \frac{n}{N} \right) \frac{S^2}{n} \right]$ . The quantity  $1 - \frac{n}{N}$  appearing in the expression above is a correction factor for the finite size of the population and is called

the finite population correction factor(fpc) or simply the finite multiplier. If  $n$  is very small compared to  $N$ , the fpc is close to unity and the sampling variance of  $\bar{y}$  in srswor will be approximate the same as srswr. If  $N$  is very small say  $N \leq 10$ , then whatever  $n$ ,  $f$  is not negligible and therefore there is considerable gain in using srswor over srswr.

**NOTE:** The finite population correction factor (fpc) is used when you **sample without replacement from more than 5% of a finite population**. It's needed because under these circumstances, the Central Limit Theorem doesn't hold and the standard error of the estimate (e.g. the mean or proportion) will be too big.

## 2.5 Confidence intervals for population mean $\bar{Y}$ and Total $Y$

The sample mean  $\bar{y}$  and the variance  $s^2$  are point estimates of the unknown population mean and variance respectively. An interval estimate of unknown population parameter is a random interval constructed such that it has a given probability of including the parameters. Consider a population with unknown parameter, if one can find an interval  $(a, b)$  such that;

$$P(a \leq \theta \leq b) = 0.95 \quad (2.16)$$

then we say that  $(a, b)$  is a 95% confidence interval for  $\theta$ . It is important to realize that the  $\theta$  is fixed and the intervals themselves vary.

Some conditions exist under which the distribution of the sample mean in a simple random sampling tends to normal distribution. If the **sample size** is not too small and the distribution of the population from which the sample is drawn is not different from the **normal**, then in srswor, the sample mean  $\bar{y}$  is approximately normal with mean  $\bar{Y}$  and deviation  $\frac{\sqrt{N-n}}{\sqrt{Nn}} S$  i.e.

$$\bar{y} \sim N\left(\bar{Y}, \frac{N-n}{Nn} S^2\right) \quad (2.17)$$

$$z = \frac{\bar{y} - \bar{Y}}{\sqrt{\frac{N-n}{Nn}} S} \sim N(0, 1)$$

$$\begin{aligned} \text{Hence } P\left(-z_{\frac{\alpha}{2}} \leq \frac{\bar{y} - \bar{Y}}{\sqrt{\frac{N-n}{Nn}} S} \leq z_{\frac{\alpha}{2}}\right) &= 1 - \alpha \\ \Rightarrow P\left(\bar{y} - z_{\frac{\alpha}{2}} \sqrt{\frac{N-n}{Nn}} S \leq \bar{Y} \leq \bar{y} + z_{\frac{\alpha}{2}} \sqrt{\frac{N-n}{Nn}} S\right) &= 1 - \alpha \end{aligned}$$

where  $z_{\frac{\alpha}{2}}$  is the 100  $\left[1 - \frac{\alpha}{2}\right]$  % point of normal distribution. Therefore;

$$\left(\bar{y} - z_{\frac{\alpha}{2}} \sqrt{\frac{N-n}{Nn}} S, \bar{y} + z_{\frac{\alpha}{2}} \sqrt{\frac{N-n}{Nn}} S\right) \quad (2.18)$$

is the 100  $\left[1 - \frac{\alpha}{2}\right]$  % confidence interval for  $\bar{Y}$ . For  $\alpha = 0.05, 0.025, 0.01$  values for  $z_{\frac{\alpha}{2}}$  are

1.96, 2.24 and 2.58 respectively.

**Example 2.1.** In a private library, the books are kept on 130 shelves of similar size. The numbers of books on 15 shelves picked at random were found to be 28,23,25,33,31,18,22,29,30,22,26,20,21,28 and 25. Estimate the total number  $Y$ , of books in the library and calculate an approximate 95% confidence interval for  $Y$ .

**Solution 1.**  $N = 130, n = 15, \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i = \frac{1}{15} (28 + 23 + \dots + 25) = 25.4$

$Y = N\bar{y} = 130 \times 25.4 = 3302$ . The 95% confidence interval is given by;

$$Y \pm N z_{0.05} \sqrt{\text{var}(\bar{y})}$$

$$\text{but } S^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2 = \frac{1}{N-1} \left( \sum_{i=1}^N Y_i^2 - N\bar{Y}^2 \right)$$

$$\text{which is estimated by } s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \left( \sum_{i=1}^n y_i^2 - n\bar{y}^2 \right)$$

$$\sum_{i=1}^{15} y_i^2 = (28^2 + 23^2 + \dots + 25^2) = 9947$$

$$n\bar{y}^2 = 15 \times (25.4)^2. \text{ The 95\% confidence interval for } Y \text{ at } \alpha = 0.05 \text{ will be;}$$

$$\hat{Y} = Y \pm N z_{0.05} \sqrt{\text{var}(\bar{y})}$$

$$= 3302 \pm 130 (1.96 \times \sqrt{1.14})$$

$$= 3302 \pm 272.05$$

$$\Rightarrow 3029.05 \leq Y \leq 3574.05$$

## 2.6 Sampling for proportions and percentages

In many situations, the characteristic under study on which the observations are collected are **qualitative in nature**. For example, the responses of customers in many marketing surveys are based on replies like 'yes' or 'no', 'agree' or 'disagree'. Sometimes the respondents are asked to arrange several options in the order like first choice, second choice etc. Sometimes the objective of the survey is to estimate the **proportion or the percentage** of brown eyed persons, unemployed persons, graduate persons or persons favoring a proposal, etc.

In such situations, the first question arises how to do the **sampling** and secondly how to estimate the population parameters like **population mean, population variance**, etc.

The same sampling procedures that are used for drawing a sample in case of quantitative characteristics can also be used for drawing a sample for qualitative characteristic. So, the sampling procedures **remain same irrespective of the nature of characteristic under study - either qualitative or quantitative**. For example, the SRSWOR and SRSWR procedures for drawing the samples remain the same for qualitative and quantitative characteristics. Similarly, other sampling schemes like stratified sampling, two stage sampling etc. also remain same.

### 2.6.1 Estimation of population proportion

The population proportion in case of qualitative characteristic can be estimated in a similar way as the estimation of population mean in case of quantitative characteristic. Consider a qualitative characteristic based on which the population can be divided into two mutually exclusive classes, say  $C$  and  $C^*$ .

For example, if  $C$  is the part of population of persons saying 'yes' or 'agreeing' with the proposal then  $C^*$  is the part of population of persons saying 'no' or 'disagreeing' with the proposal. Let  $A$  be the number of units in  $C$  and  $(N - A)$  units in  $C^*$  be in a population of size  $N$ . Then the proportion of units in  $C$  is;

$$P = \frac{A}{N} \quad (2.19)$$

and the proportion of units in  $C^*$  is

$$Q = \frac{N - A}{N} = 1 - P \quad (2.20)$$

An indicator variable  $Y$  can be associated with the characteristics under study and then for

$i = 1, 2, \dots, N$ .  $Y_i = 1$  if the  $i^{th}$  unit belongs to  $C$  and 0 if the  $i^{th}$  unit belongs to  $C^*$ . Now the population total is;

$$Y_{TOTAL} = \sum_{i=1}^N Y_i = A \quad (2.21)$$

and the population mean is;

$$\bar{Y} = \frac{\sum_{i=1}^N Y_i}{N} = \frac{A}{N} = P \quad (2.22)$$

Suppose a sample of size  $n$  is drawn from a population of size  $N$  by simple random sampling. Let  $a$  be the number of units in the sample which fall into class  $C$  and  $(n - a)$  units fall in class  $C^*$ , then the sample proportion of units in  $C$  is;

$$p = \frac{a}{n} \quad (2.23)$$

which can be written as  $p = \frac{a}{n} = \frac{\sum_{i=1}^n y_i}{n} = \bar{y}$ .

Since,  $\sum_{i=1}^N Y_i = A = NP$  so we can write  $S^2$  and  $s^2$  in terms of  $Q$  and  $P$  as follows;

$$\begin{aligned} S^2 &= \frac{1}{N-1} \sum_{i=1}^N (Y_i - \bar{Y})^2 = \frac{1}{N-1} \sum_{i=1}^N (Y_i^2 - N\bar{Y}^2) \\ &= \frac{1}{N-1} \sum_{i=1}^N (NP - NP^2) = \frac{N}{N-1} PQ \end{aligned}$$

$$\begin{aligned} \text{Similarly; } \sum_{i=1}^n y_i^2 &= a = np \text{ and } s^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i^2 - n\bar{y}^2) \\ &= \frac{1}{n-1} \sum_{i=1}^n (np - np^2) \end{aligned}$$

$$= \frac{n}{n-1}pq$$

Note that the quantities  $\bar{y}$ ,  $\bar{Y}$ ,  $s^2$  and  $S^2$ , have been expressed as functions of sample and population proportions. Since the sample has been drawn by simple random sampling and sample proportion is same as the sample mean, so the properties of sample proportion in SRSWOR and SRSWR can be derived using the properties of sample mean directly.

SRSWOR Since the sample mean  $\bar{y}$  is an unbiased estimator of the population mean  $\bar{Y}$  i.e.

$E(\bar{y}) = \bar{Y}$  in the case of SRSWOR, so;

$E(p) = E(\bar{y}) = \bar{Y} = P$  and  $p$  is an unbiased estimator of  $P$ . Using the expression of  $Var(\bar{y})$  the variance of  $p$  can be derived as  $Var(p) = Var(\bar{y}) = \frac{N-n}{Nn} S^2$

Similarly, using the estimate of variance can be derived as

$$\begin{aligned}\widehat{Var}(p) &= \widehat{Var}(\bar{y}) = \frac{N-n}{Nn} S^2 \\ &= \frac{N-n}{Nn} \cdot \frac{n}{n-1} pq\end{aligned}$$

$$= \frac{N-n}{N(n-1)} pq \quad (2.24)$$

SRSWR Since the sample mean  $\bar{y}$  is an unbiased estimator of population mean  $\bar{Y}$  in case of SRSWR, so the sample proportion

$E(p) = E(\bar{y}) = \bar{Y} = P$  i.e.,  $p$  is an unbiased estimator of  $P$ .

Using the expression of variance of  $\bar{y}$  and its estimate in case of SRSWR, the variance of  $p$  and its estimate can be derived as follows:  $Var(p) = Var(\bar{y}) = \frac{N-1}{Nn} S^2$

$$\begin{aligned}&= \frac{N-1}{Nn} \frac{N}{N-1} PQ \\ &= \frac{PQ}{n} \\ \Rightarrow \widehat{Var}(p) &= \frac{n}{n-1} \cdot \frac{pq}{n}\end{aligned}$$

$$= \frac{pq}{n-1} \quad (2.25)$$

## 2.6.2 Estimation of population total or total number of count

It is easy to see that an estimate of population total  $A$  (or total number of count) is  $\hat{A} = NP = \frac{Na}{n}$  its variance is  $Var(\hat{A}) = N^2 Var(p)$  and the estimate of variance is  $\widehat{Var}(\hat{A}) = N^2 \widehat{Var}(p)$

## 2.6.3 Confidence Interval estimation for P

If  $N$  and  $n$  are large, then  $\frac{p-P}{\sqrt{Var(p)}}$  approximately follows  $N(0, 1)$ . With this approximation we can write  $P\left(-z_{\frac{\alpha}{2}} \leq \frac{p-P}{\sqrt{Var(p)}} \leq z_{\frac{\alpha}{2}}\right) = 1 - \alpha$ , and the  $100(1 - \alpha)$  confidence interval of  $P$  is

$$p - z_{\frac{\alpha}{2}} \sqrt{Var(p)}, p + z_{\frac{\alpha}{2}} \sqrt{Var(p)} \quad (2.26)$$

It may be noted that in this case, a discrete random variable is being approximated by a continuous random variable, so a continuity correction  $\frac{n}{2}$  can be introduced in the confidence limits and the limits become;

$$p - z_{\frac{\alpha}{2}} \sqrt{Var(p)} + \frac{n}{2}, p + z_{\frac{\alpha}{2}} \sqrt{Var(p)} + \frac{n}{2} \quad (2.27)$$

## 2.7 Determination of sample sizes

In a field survey the statisticians would like to have a sample size that will give a desired level of precision of estimator. We note that the required precision is the difference between the estimator and the true value. This difference is denoted by  $d$ . Suppose that it is desired to find a sample size  $n$  such that the estimated value i.e. sample mean  $\bar{y}$  differs from the true value (Population mean,  $\bar{Y}$ ) by a quantity not exceeding  $d$  with a very high probability, say greater than  $1 - \alpha$ . Hence the problem is to find  $n$  such that;

$$P(|\bar{y} - \bar{Y}| \leq d) \geq 1 - \alpha \quad (2.28)$$

From srswor  $\bar{y} \sim N(\bar{Y}, \frac{N-n}{Nn} S^2)$ . Hence,

$$P\left(|\bar{y} - \bar{Y}| \leq S \sqrt{\frac{N-n}{Nn}} t\right) = 1 - \alpha \quad (2.29)$$

where  $t = z_{\frac{\alpha}{2}}$  is the 100  $(1 - \frac{\alpha}{2})$  point of normal distribution. From equation 2.28 and 2.29,  $tS \sqrt{\frac{N-n}{Nn}} = d$  where  $\frac{1}{n} = \frac{1}{N} + \frac{d^2}{t^2 S^2}$ . Hence,  $n = \frac{(\frac{tS}{d})^2}{1 + \frac{1}{N}(\frac{tS}{d})^2}$ . As a first approximation we may take  $n_o = (\frac{tS}{d})^2$ . If  $\frac{n_o}{N}$  is negligibly small, this may be taken as the satisfactory value of  $n$ . If not, one should calculate  $n = \frac{n_o}{1 + (\frac{n_o}{N})} = n_o \left(1 + \frac{n_o}{N}\right)^{-1}$ . In practice one has to replace  $S$  by an advance estimate  $s'$  (say). In case the problem is that of estimating a population proportion one may require to find  $n$  such that

$$P(|p - P| \leq d) \geq 1 - \alpha \quad (2.30)$$

For large samples in srswor  $\frac{p-P}{\sqrt{\frac{N-n}{n(N-1)} PQ}}$  is approximately a normal variable. Hence;

$$P\left(|p - P| \leq t \sqrt{\frac{N-n}{n(N-1)} PQ}\right) = 1 - \alpha \quad (2.31)$$

Equating 2.28 and 2.29 we get  $t \sqrt{\frac{N-n}{n(N-1)} PQ} = d$ . This gives;

$$n = \frac{\left(\frac{t^2 PQ}{d^2}\right)}{1 + \left(\frac{1}{N}\right) \left[\left(\frac{t^2 PQ}{d^2}\right) - 1\right]} \quad (2.32)$$



For practical purposes,  $P$  is to be replaced by some suitable estimate  $p$  of the same. For large  $N$  a first approximation of  $n$  is  $n_0 = \frac{t^2 PQ}{d^2}$ . If  $\frac{n_0}{N}$  is negligible,  $n_0$  is a satisfactory approximation to  $n$ . If not, one should calculate  $n$  as;

$$n = \frac{n_0}{1 + \left[\left(\frac{n_0-1}{N}\right)\right]} \approx \frac{n_0}{1 + \left(\frac{n_0}{N}\right)} \quad (2.33)$$

**Example 2.2.** Suppose it is required to estimate the average value of output of a group of 5000 factories in a region so that the sample estimate lies within 10 of the true value with a confidence coefficient of 95%. Determine the minimum sample size required. The population coefficient of variation is known to be 60%.

**Solution 2.** We require  $n$  such that  $P(|\bar{y} - \bar{Y}| \leq 0.1\bar{Y}) = 0.95$ . Now under normal approximation,

$$(|\bar{y} - \bar{Y}| \leq 0.1\bar{Y}) = 0.95. \text{ Hence, } 1.96S\sqrt{\frac{N-n}{Nn}} = 0.1\bar{Y}$$

$$\text{or } (1.96)^2 \left(\frac{1}{n} - \frac{1}{N}\right) = 0.01 \left[\frac{\bar{Y}}{S}\right]^2 = \frac{0.01}{0.36}$$

Solving the above equation, we get  $n = 136$  (rounded off to the next integer)

**Example 2.3.**

Consider the population consisting of 430 units. By complete enumeration of the population it was found that  $\bar{Y} = 19$ ,  $S^2 = 85.6$ . These being true population values with simple random samples, how many units must be taken to estimate  $\bar{y}$  with 10% of  $\bar{Y}$  a part from a chance of 1 in 20.

**Solution 3.**  $\bar{Y} = 19$ ,  $S^2 = 85.6 \Rightarrow S = \sqrt{85.6}$   $N = 430$ ,  $d = \frac{1}{20} = 0.05$ . 10%

$$\text{of } \bar{Y} \Rightarrow d = 0.1\bar{Y} = 0.1(19) = 1.9.$$

$$n_0 = \left(\frac{ts}{d}\right)^2$$

$$\text{but } t = z_{\frac{\alpha}{2}} = z_{\frac{0.05}{2}} = z_{0.025} = 1.96.$$

$$\Rightarrow n_0 = \frac{(1.96)^2(85.6)}{1.9^2} = 91.09167.$$

$$n = n_0 \left(1 + \frac{n_0}{N}\right)^{-1}, = 91.09 \left[\left(1 + \frac{91.09}{430}\right)\right]^{-1} = 75.166 \simeq 75.$$

## 2.8 Exercises

1. In a population with  $N = 6$ , the values of  $y_i$  are 8, 3, 1, 11, 4, and 7. Calculate the sample mean  $\bar{y}$  for all possible simple random samples of size 2. Verify that  $\bar{y}$  is an unbiased estimate of  $\bar{Y}$ .
2. For the same population in 1 above, calculate  $s^2$  for all simple random samples of size 3, and verify that  $E(s^2) = S^2$ .
3. If random samples of size 2 are drawn with replacement (from this population, show by finding all possible samples that  $Var(\bar{y})$  satisfies the equation  $Var(\bar{y}) = \frac{\sigma^2}{n} = \frac{S^2(N-1)}{nN}$ . Give a general proof of this result.

4. A simple random sample of 30 households was drawn from a city area containing 14,848 households. The numbers of persons per household in the sample were as follows: 5,6,3,3,2,3,3,3,4,4,3,2,7,4,3,5, 4,4,3,3,4,3,3, 1,2,4,3,4,2,4 Estimate the total number of people in the area and compute the probability that this estimate is within  $\hat{A} \pm 10$  per cent of the true value.
5. Consider a population consisting of 6 villages, the areas (in hectares) of which are given below;

**Table 1:** Population of 6 villages

Village	A	B	C	D	E	F
Area	760	343	657	550	480	935

- (a) Enumerate all possible WR samples of size 3. Also, write the values of the study variable for the sampled units.
- (b) List all the WOR samples of size 4 along with their area values.
6. Among the 100 computer corporations in a region, average of the employee sizes for the largest 10 and smallest 10 corporations were known to be 300 and 100, respectively. For a sample of 20 from the remaining 80 corporations, the mean and standard deviation were 250 and 110, respectively. For the total employee size of the 80 corporations, find;
- (a) the estimate
- (b) the S.E. of the estimate
- (c) the 95% confidence limits
7. Continuing with Exercise 2, for the average and total of the 100 corporations, find;
- (a) the estimate
- (b) the S.E. of the estimate
- (c) the 95% confidence limits.
8. The height (in cm) of 6 students of M.Sc., majoring in statistics, from Punjab Agricultural University, Ludhiana was recorded during 1985. The data, so obtained, are given below:
- Calculate;
- (a) Calculate the population mean  $\bar{Y}$  and population variance  $\sigma^2$ .
- (b) Enumerate all possible SRS with replacement samples of size  $n = 2$ . Obtain sampling distribution of mean, and hence show that:

i.  $E[\bar{y}]$

**Table 2:** Heights of M.Sc. students

Student	Name	Height
1	A	168
2	B	175
3	C	185
4	D	173
5	E	171
6	F	172

ii.  $V[\bar{y}] = \frac{\sigma^2}{n}$

iii.  $E[s^2] = \sigma^2$

iv.  $E[v(\bar{y})] = V(\bar{y})$

(c) Enumerate all possible SRS without replacement samples of size  $n = 2$ . Obtain sampling distribution of mean, and hence show that:

i.  $E[\bar{y}]$

ii.  $V[\bar{y}] = \frac{\sigma^2}{n}$

iii.  $E[s^2] = \sigma^2$

iv.  $E[v(\bar{y})] = V(\bar{y})$

9. Punjab Agricultural University, Ludhiana, is interested in estimating the proportion  $P$  of teachers who consider semester system to be more suitable as compared to the trimester system of education. A with replacement simple random sample of  $n=120$  teachers is taken from a total of  $N=1200$  teachers. The response is denoted by 0 if the teacher does not think the semester system suitable, and 1 if he/she does.

**Table 3:** Punjab Agricultural University

Teacher	1	2	3	4	5	6	...	119	120	Total
Response	1	0	1	1	0	1	...	0	1	72

- (a) From the sample observations given below, estimate the proportion  $P$  along with the standard error of your estimate. Also, work out the confidence interval for  $P$ .
- (b) While estimating  $P$ , the investigator feels that the tolerable error could be taken as 0.08. Do you think the sample size 120 is sufficient? If not, how many more units should be included in the sample?