## MODULE-3 SAMPLING TECHNIQUES

**Sampling** may be defined as the selection of some part of an aggregate or totality on the basis of which a judgement or inference about the aggregate or totality is made. In other words, it is the process of obtaining information about an entire population by examining only a part of it.

Sampling is used in practice for a variety of reasons such as:

- 1. Sampling can save time and money. A sample study is usually less expensive than a census study and produces results at a relatively faster speed.
- 2. Sampling may enable more accurate measurements for a sample study is generally conducted by trained and experienced investigators.
- 3. Sampling remains the only way when the population contains infinitely many members.
- 4. Sampling remains the only choice when a test involves the destruction of the item under study.
- 5. Sampling usually enables us to estimate the sampling errors and, thus, assists in obtaining information concerning some characteristics of the population.

**Universe/Population:** From a statistical point of view, the term 'Universe 'refers to the total of the items or units in any field of inquiry, whereas the term 'population' refers to the total of items about which information is desired.

The attributes that are the object of study are referred to as **characteristics** and the units possessing them are called **elementary units**. The aggregate of such units is generally described as population. Thus, all units in any field of inquiry constitute the universe and all elementary units (based on one characteristic or more) constitute population.

The population or universe can be finite or infinite. The population is said to be **finite** if it consists of a fixed number of elements so that it is possible to enumerate it in its totality. An **infinite** population is that population in which it is theoretically impossible to observe all the elements. Thus, in an infinite population the number of items is infinite i.e., we cannot have any idea about the total number of items.

**Sampling frame:** The elementary units or the group or cluster of such units may form the basis of sampling process in which case they are called sampling units. A list containing all such sampling units is known as sampling frame. Thus, the sampling

frame consists of a list of items from which the sample is to be drawn. For instance, one can use telephone directory as a frame for conducting opinion survey in a city. Whatever the frame may be, it should be a good representative of the population.

**Sampling design:** A sample design is a definite plan for obtaining a sample from the sampling frame. It refers to the technique or the procedure the researcher would adopt in selecting some sampling units from which inferences about the population is drawn. Sampling design is determined before any data is collected.

\*Statistic(s) and parameter(s): A statistic is a characteristic of a sample, whereas a parameter is a characteristic of a population.

Thus, when we work out certain measures such as mean, median, mode or the like ones from samples, then they are called statistic(s) for they describe the characteristics of a sample. But when such measures describe the characteristics of a population, they are known as parameter(s). For instance, the population mean  $\mu$  is a parameter, whereas the sample mean is a statistic.

\*Precision: Precision is the range within which the population average (or other parameter) will lie in accordance with the reliability specified in the confidence level as a percentage of the estimate ± or as a numerical quantity. For instance, if the estimate is Rs 4000 and the precision desired is ± 4%, then the true value will be no less than Rs 3840 and no more than Rs 4160. This is the range (Rs 3840 to Rs 4160) within which the true answer should lie. But if we desire that the estimate should not deviate from the actual value by more than Rs 200 in either direction, in that case the range would be Rs 3800 to Rs 4200.

\*Confidence level and significance level: The confidence level or reliability is the expected percentage of times that the actual value will fall within the stated precision limits. Thus, if we take a confidence level of 95%, then we mean that there are 95 chances in 100 (or .95 in 1) that the sample results represent the true condition of the population within a specified precision range against 5 chances in 100 (or .05 in 1) that it does not. Precision is the range within which the answer may vary and still be acceptable; confidence level indicates the likelihood that the answer will fall within that range, and the significance level indicates the likelihood that the answer will fall outside that range. We can always remember that if the confidence level is 95%, then the significance level will be (100 - 95) i.e., 5%; if the confidence level is 99%, the significance level is (100 - 99) i.e., 1%, and so on.

**Sampling error:** Sample surveys do imply the study of a small portion of the population and as such there would naturally be a certain amount of inaccuracy in the information collected. This inaccuracy may be termed as sampling error or error variance. In other words, sampling errors are those errors which arise on account of sampling, and they generally happen to be random variations (in case of random sampling) in the sample estimates around the true population values. The magnitude of the sampling error depends upon the nature of the universe, the more homogeneous the universe, the

smaller the sampling error. Sampling error is inversely related to the size of the sample i.e., sampling error decreases as the sample size increases and vice-versa. A measure of the random sampling error can be calculated for a given sample design and size and this measure is often called the precision of the sampling plan.

Sampling error = Frame error + Chance error + Response error

\*Standard Error: The standard deviation of sampling distribution of a statistic is called S.E.

The utility of the concept of standard error in statistical induction arises on account of the following reasons:

1)The standard error helps in testing whether the difference between observed and expected frequencies could arise due to chance.

[FOR FURTHER READING ONLY: The criterion usually adopted is that if a difference is less than 3 times the S.E., the difference is supposed to exist as a matter of chance and if the difference is equal to or more than 3 times the S.E., chance fails to account for it, and we conclude the difference as significant difference. This criterion is based on the fact that at X ± 3 (S.E.) the normal curve covers an area of 99.73 per cent. Sometimes the criterion of 2 S.E. is also used in place of 3 S.E. Thus, the standard error is an important measure in significance tests or in examining hypotheses. If the estimated parameter differs from the calculated statistic by more than 1.96 times the S.E., the difference is taken as significant at 5 per cent level of significance. This, in other words, means that the difference is outside the limits i.e., it lies in the 5 per cent area (2.5 per cent on both sides) outside the 95 per cent area of the sampling distribution. Hence, we can say with 95 per cent confidence that the said difference is not due to fluctuations of sampling. In such a situation our hypothesis that there is no difference is rejected at 5 per cent level of significance. But if the difference is less than 1.96 times the S.E., then it is considered not significant at 5 per cent level and we can say with 95 per cent confidence that it is because of the fluctuations of sampling. In such a situation our null hypothesis stands true. 1.96 is the critical value at 5 per cent level. The product of the critical value at a certain level of significance and the S.E. is often described as 'Sampling Error' at that particular level of significance. We can test the difference at certain other levels of significance as well depending upon our requirement.]

2) The standard error gives an idea about the reliability and precision of a sample. The smaller the S.E., the greater the uniformity of sampling distribution and hence, greater is the reliability of sample. Conversely, the greater the S.E., the greater the difference between observed and expected frequencies. In such a situation the unreliability of the sample is greater. The size of S.E. depends upon the sample size to a great extent, and it varies inversely with the size of the sample. If double reliability is required i.e., reducing S.E. to 1/2 of its existing magnitude, the sample size should be increased fourfold.

3) The standard error enables us to specify the limits within which the parameters of the population are expected to lie with a specified degree of confidence. Such an interval is usually known as confidence interval.

## **SAMPLING THEORY**

Sampling theory is a study of relationships existing between a population and samples drawn from the population. Sampling theory is applicable only to random samples. The theory of sampling studies the relationships that exist between the universe and the sample or samples drawn from it. The main problem of sampling theory is the problem of relationship between a parameter and a statistic. The theory of sampling is concerned with estimating the properties of the population from those of the sample and also with gauging the precision of the estimate. This sort of movement from particular (sample) towards general (universe) is what is known as statistical induction or statistical inference.

Sampling theory is designed to attain one or more of the following objectives:

- (i) Statistical estimation: Sampling theory helps in estimating unknown population parameters from a knowledge of statistical measures based on sample studies. In other words, to obtain an estimate of parameter from statistic is the main objective of the sampling theory. The estimate can either be a point estimate or it may be an interval estimate. Point estimate is a single estimate expressed in the form of a single figure, but interval estimate has two limits viz., the upper limit and the lower limit within which the parameter value may lie. Interval estimates are often used in statistical induction.
- (ii) Testing of hypotheses: The second objective of sampling theory is to enable us to decide whether to accept or reject hypothesis; the sampling theory helps in determining whether observed differences are due to chance or whether they are really significant.
- (iii) Statistical inference: Sampling theory helps in making generalization about the population/ universe from the studies based on samples drawn from it. It also helps in determining the accuracy of such generalizations.

The theory of sampling can be studied under two heads viz., the sampling of attributes and the sampling of variables and that too in the context of large and small samples (By small sample is commonly understood any sample that includes 30 or fewer items, whereas a large sample is one in which the number of items is more than 30). When we study some qualitative characteristic of the items in a population, we obtain statistics of attributes in the form of two classes; one class consisting of items wherein the attribute is present and the other class consisting of items wherein the attribute is absent. The presence of an attribute may be termed as a 'success' and its absence a 'failure'.

We generally consider the following three types of problems in case of sampling of attributes:

- (i) The parameter value may be given, and it is only to be tested if an observed 'statistic' is its estimate.
- (ii) The parameter value is not known, and we have to estimate it from the sample.
- iii) Examination of the reliability of the estimate i.e., the problem of finding out how far the estimate is expected to deviate from the true value for the population.

The theory of sampling can be applied in the context of statistics of variables (i.e., data relating to some characteristic concerning population which can be measured or enumerated with the help of some well-defined statistical unit) in which case the objective happens to be:

- (i) To compare the observed and expected values and to find if the difference can be ascribed to the fluctuations of sampling.
- (ii) To estimate population parameters from the sample, and
- (iii) To find out the degree of reliability of the estimate.

**Sampling distribution**: If we take a certain number of samples and for each sample compute various statistical measures such as mean, standard deviation, etc., then we can find that each sample may give its own value for the statistic under consideration. All such values of a particular statistic, say mean, together with their relative frequencies will constitute the sampling distribution of the statistic, say mean. Accordingly, we can have sampling distribution of mean, or the sampling distribution of standard deviation or the sampling distribution of any other statistical measure. It may be noted that each item in a sampling distribution is a particular statistic of a sample. The sampling distribution tends quite closer to the normal distribution if the number of samples is large. The significance of sampling distribution follows from the fact that the mean of a sampling distribution is the same as the mean of the universe. Thus, the mean of the sampling distribution can be taken as the mean of the universe.

## IMPORTANT SAMPLING DISTRIBUTIONS

Some important sampling distributions, which are commonly used, are:

- (1) sampling distribution of mean
- (2) sampling distribution of proportion
- (3) student's 't' distribution
- (4) F distribution
- (5) Chi-square distribution.

Sampling distributions: Any statestic, being a random veriable, how a probability distribution of a statester is called Sampling distribution.

Types of Sampling distribution

(1) Sampling distribution of Mean

2) Sampling distribution of proportion

3 t-distribution X2-distribution

(Chi-square)

(5) = - distribution.

Standard error: The standard deviation of a of the Campling distribution of a statistic is called Its Standard error.

(1) Sampling distribution of mean (X) It is the probability distribution of all the possible means of random Samples of a given Size that we take from the population.

observations from a population with mean el and S.D G.

het X denote the mean of these Observations 1e, X = X1+X2+ · · Xn X is a random variable known as Sampling distribution of Sample mean.  $F(X) = F(\overline{X1 + \cdots \times Xn})$ = \frac{1}{n} \{ \mathbb{E}(X))+ - \mathbb{E}(X)\} = I { ut... ul }  $V(X) = V(X_1 + \cdots \times_n)$ = - 1= { v(x1)+ .. v(x,0)}  $=\frac{1}{m^2}\left\{\sigma_1^2+\cdots\sigma_r^2\right\}$ = 62 :. S.D(X) = 5 Note: 5x = 5 is called standard error of the Sample mean. · If X1, X2 · · Xn be a random sample from a normal population with mean u and S.D.G. 18, Xia N(U, 52) then by CLT X~N(4,52)

To reduce the sampling distribution of mean to unit normal distribution 1e, N10,1) we can write

by CLT, X is normally distributed provided the no. of Sample items

are large & n ≥ 303

· Standard error of mean when population S.D is known

· S.E of mean when population S.D

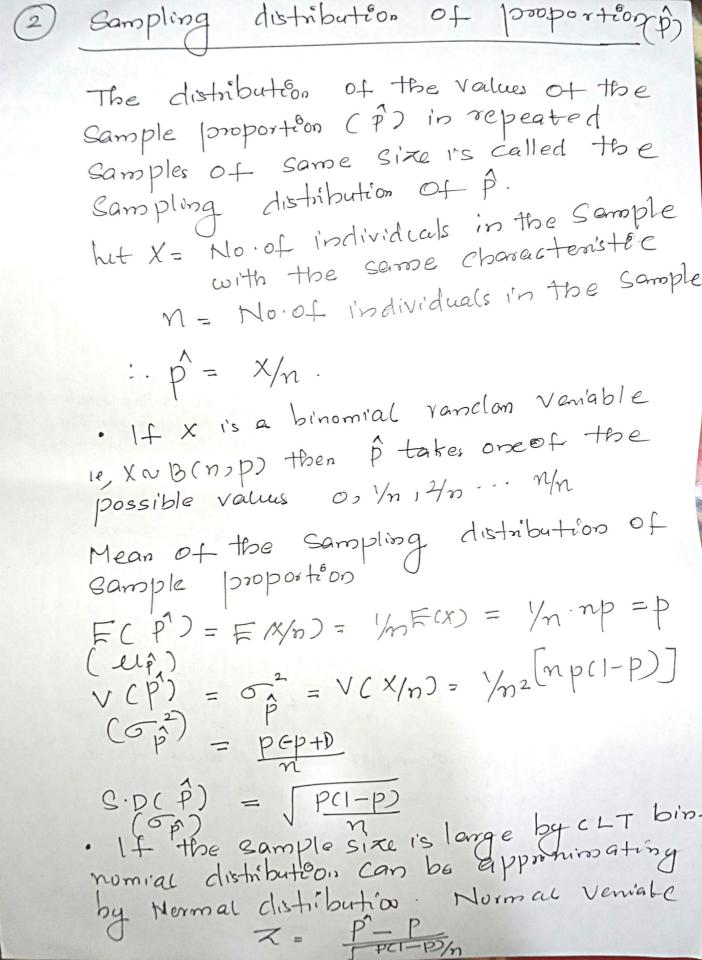
$$\frac{\sigma_{X}}{\sqrt{n}} = \frac{\sigma_{S}}{\sigma_{S} - S \cdot D \cdot of sample}$$

$$\frac{\sigma_{S} - S \cdot D \cdot of sample}{\sqrt{n-1}}$$

of two samples

(a) When two samples are drawn from the same population

2 If op is not known, sample s.D for Combined Samples denoted by ognation may be Substituted? (b) When two samples are drawn from different population  $\frac{1}{X_1-X_2} = \sqrt{\frac{p_1^2}{\gamma_1} + \frac{p_2^2}{\gamma_2}}$ Elf op, op are not known, then os, and os, respectively may be substituted. · S.F. of mean when population is finite  $\frac{\sqrt{N-n}}{\sqrt{n}} = \frac{\sqrt{N-n}}{\sqrt{N-1}}$ N-n is called Correction factor.



S.E of spoportion 
$$CP$$
)
$$\hat{p} = \sqrt{\frac{pq}{n}} \text{ or } \sqrt{\frac{pq}{n}} \sqrt{\frac{N-n}{n}}$$

n = Samplesize P = population 9=1-p Proportion N = populateon

Pq N-n Q=1-p prof Coupen population N= population Size is given)

S. E of difference blw proportions of two samples

a) 
$$\frac{1}{P_1 - P_2} = \sqrt{\frac{pq}{m_1} + \frac{pq}{m_2}}$$

where  $P = n_1 P_1 + n_2 P_2$  { It is the best estimate of proportion in the population.

 $M_1 = No \cdot of events in sample 1$   $M_2 = 1$ M2 = "

b) 
$$rac{P_1-P_2}{V_1-P_2} = \sqrt{\frac{P_1Q_1}{N_1} + \frac{P_2Q_2}{N_2}}$$

when samples are drawn from two beterogeneous population.

Distribution of the sample proportion is

· Normal If np > 5 and n (1-p) > 5

o Binomial If np25 and nc1-po25

## 3 Students to distribution

If we down a random sample of nobservations from a normally distributed as population with mean it and variance of then X = X - ut follows Std normal distribution. In storme cases or is not known (population S.D is denoted by so reflace or in (1) by sample standard deviation of so so (1) does not follow normal distribution. This new type of distribution is called to distribution

10, 
$$t = \frac{\overline{X} - u}{\overline{S}}$$
 where  $\overline{S} = \sqrt{\frac{\overline{X}(\overline{X} - \overline{X})^2}{M - 1}}$ 

bas (m-1) degrees of freedom.

Degrees of Freedom of a Statistic is a tree integer denoted by D is the no. of values in the final Calculation of a Statistic that one free to vary.

In general,  $\hat{v}=n-k$ 

M = Sample Size

K = no of independent Constrains
imposal on the Observations
in the Sample.

Note: when population S.D 6 is not known and the Sample is of a Small Sixe (ic, n 430), we use to distribution for the Sampling distribution of mean.

A F-distributeon

This test is usually used to test equality of vernances of two independent Sample of Size no and no mespectively taken from two independent normal population having the Same variance,  $\sigma_{12}^{2} = \sigma_{12}^{2}$ 

F =  $\frac{\sigma_{s_1}^2}{\sigma_{s_2}^2}$  where  $\frac{\sigma_{s_1}^2}{\sigma_{s_2}^2} = \sum \frac{(x_{1i} - x_{1})^2}{m_{1-1}}$ 

bas an F distribution with n<sub>1</sub>-1 and n<sub>2</sub>-1 degrees of freedom. Fration's Computed in a way that the larger Variance is always to the numerator.

X2 distribution with degrees of - Sreedom k is the distribution of a Sum of the squares of k independent standard normal random veniables

the following.

1) To check the relationships between Categorical variables.

2) To check the pendence of two criteria of classification of multiple qualitative verniables 3) To conduct a The Chi-Saucrotest.

The formula for the chaquare statestec used in a Cha-savare test is

 $\mathcal{X}_{\nu}^{2} = \sum_{i=1}^{k} \frac{(Q_{i} - E_{i})^{2}}{E_{i}^{2}} \xrightarrow{\mathcal{D}} \xrightarrow{\text{Degrees of freedom}} V_{i}$   $\mathcal{X}_{\nu}^{2} = \sum_{i=1}^{k} \frac{(Q_{i} - E_{i})^{2}}{E_{i}^{2}} \xrightarrow{\text{Observed value}} V_{i}$ 

Tables are there that geve the value of  $\chi^2$  for given def which may be used with calculated value of  $\chi^2$  for relevent def at a desired level of Significance for testing hypotheses.

Precision: It is the range within which the population parameter well lie in accordance with the reliability in accordance with the reliability expectived in the Considence level specified in the Considence level as a percentage of the estimate to as a numerical quantity.

Precision elesized is ±4% then love value is in blu 9600 and 10400 \ 2 10000 - tox10000, 10000 + tox10000}

Considence level: It is the enpected of of times that the adual value will of times that the adual value will be stated precision limits. Fall within the stated precision limits. Fall within the accordance level of Eq. If we take a Confidence level of their of their are are as that there are as chances in 100 that the sample as chances in 100 that the force Condition of the population with in a specific of the population with in a specific of the population with in a specific of the population range against to chances in 100 that It does not.

Bignificance level! It is the probability
that the answer well fall outside
precision.

if Confidence level is 95%. Then significance level is 5%.

Confidence interval. A probability that a parameter will fall between a Set of values. If t is the Statistic used to estimate perameter 0, then 1-x Confidence Cimits for O + + S.E(+) + x/2 where tx12 is the Significant or entical Value of t at level of significance & for a two tealed test. 1e, p(t- s.E(t)tx/2 , t+ s.E(t)tx/2) 1-00 is called Considence Westicient are a is called level of significance. Result: If the comple size is longe than the tunder laying distribution of the Standardised veriate corresponding to the Sampling distribution of the statestic t follows numace distribution 1, Z = t= (1) N(0)) · 30 1-2 Confidence Limits for 0 13

+ + ZX/2° S.F.(+)

To estimate population means sample mean I is the best estimater of population mean (el) 1-a Confidence limits - for el 1's · X ± Za/2 os popular fon size 13 Y ± Za/2 os N-n population size(N) X ± ta/2 of sample · To estimate population proportion Sample proportion is the best exturna ter for population propostion 1-a Confidence lumits often population proportion is p + z pq cample of 36 New

.