

STATISTICAL INFERENCE II

Statistical Inference could be defined as the act of drawing inferences about a population or characteristics from information contained in a sample. In this process, if a population value is known, there will be no need to make inferences about them.

SAMPLING THEORY

Sampling theory is a study of relationship of a population and samples drawn from the population. It is of great value in many connections. It is useful in estimation of unknown population quantities such as population mean, variance etc. Often called population parameter or briefly PARAMETERS from a knowledge of corresponding sample qualities such as sample mean, variance, etc. Often called sample statistics or briefly STATISTICS.

In general, a study of inferences made concerning a population by use of samples drawn from it, together with indications of the accuracy of such inferences using probability theory is called STATISTICAL INFERENCE.

SAMPLING DISTRIBUTION

If we take a sample from a given size of a population and repeat the sampling as many times as possible, calculate the mean of each sample and summarize into a frequency distribution, the resulting distribution is called a SAMPLING DIST of the sample mean.

Suppose we have a population

$$S = (1, 2, 3, 4, 5, 6)$$

A possible interpretation of the population can be that it is a list of all possible outcomes of throwing a die once. The probability distribution of this population can be as follows:

Fig 1.1 \Rightarrow probability of throwing a die once

x	Prob(x)
1	$\frac{1}{6}$
2	$\frac{1}{6}$
3	$\frac{1}{6}$
4	$\frac{1}{6}$
5	$\frac{1}{6}$
6	$\frac{1}{6}$
Total	1

Since this is a finite population, we can examine the relationship between the population and Sampling distribution. We do this by considering the arithmetic mean as the most important parameters in statistics.

Ex 1:

- Take a sample (without replacement of any two) from the above population and calculate the mean and repeat the process until no samples have the same members.
- Carry out the same process for a sample of size four.

Solution

Sample	$\bar{x}(n=2)$	Sample	$\bar{x}(n=4)$
1, 2	1.50	1, 2, 3, 4	2.50
1, 3	2.00	1, 2, 3, 5	2.75
1, 4	2.50	1, 2, 3, 6	3.00
1, 5	3.00	1, 2, 4, 5	3.25
1, 6	3.50	1, 2, 4, 6	3.50
2, 3	2.50	1, 2, 5, 6	3.25
2, 4	3.00	1, 3, 4, 5	3.50
2, 5	3.50	1, 3, 4, 6	3.75
2, 6	4.00	1, 3, 5, 6	4.00
3, 4	3.50	1, 4, 5, 6	3.50
3, 5	4.00	1, 2, 3, 4, 5	3.50
3, 6	4.50	2, 3, 4, 6	4.00
4, 5	4.50	2, 3, 5, 6	4.25
4, 6	5.00	2, 4, 5, 6	4.50
5, 6	5.50	3, 4, 5, 6	4.50

$${}^M C_n \quad {}^P C_s \quad {}^6 C_2 = 15$$

EX 2.

classify the means resulting from samples of each of ② ③ ④ of Example 1 and calculate their relative frequency.

SOLUTION

\bar{x}	f	$P(x)$	\bar{x}	f	$\bar{x} P(x)$
1.50	1	1/15	2.50	1	2.50
2.00	1	1/15	2.75	1	2.75
2.50	2	2/15	3.00	2	3.00
3.00	2	2/15	3.25	2	3.25
3.50	3	3/15	3.50	3	3.50
4.00	2	2/15	3.75	2	3.75
4.50	2	2/15	4.00	2	4.00
5.00	1	1/15	4.25	1	4.25
5.50	1	1/15	4.50	1	4.50

Ex 3.

Calculate the mean and variance of the samples in

fig 1.1.

SOLUTION

x	f	f_x	$x - \bar{x}$	$(x - \bar{x})^2$	$f(x - \bar{x})^2$
x	$P(x)$	$x P(x)$	$x - \bar{x}$	$(x - \bar{x})^2$	$f(x - \bar{x})^2$
1	1/6	1/6	-2.5	6.25	6.25/6
2	1/6	2/6	-1.5	2.25	2.25/6
3	1/6	3/6	-0.5	0.25	0.25/6
4	1/6	4/6	0.5	0.25	0.25/6
5	1/6	5/6	1.5	2.25	2.25/6
6	1/6	6/6	2.5	6.25	6.25/6
		21/6			17.50

$$\text{Mean} = \bar{x} = \frac{\sum x_i P(x)}{\sum P(x)} = \frac{21}{6} = 3.50$$

$$\text{Variance} = \sigma^2 = \frac{\sum (x_i - \bar{x})^2 P(x)}{\sum P(x)} = \frac{17.50}{6} = 2.917$$

Similarly, from the Sampling distribution showed in Ex. 2, we can calculate the corresponding Sample mean and Variance for each Sampling data.

SOLUTION

For $n=2$

\bar{x}	$P(\bar{x})$	$\bar{x} P(\bar{x})$	$\bar{x} - 3.5$	$(\bar{x} - 3.5)^2$	$(\bar{x} - 3.5)^2 P(\bar{x})$
1.5	1/15	1.5/15	-2.0	+4.00	4.0/15
2.0	1/15	2.0/15	-1.5	2.25	2.25/15
2.5	2/15	5.0/15	-1.0	1.00	2.0/15
3.0	2/15	6.0/15	-0.5	0.25	0.50/15
3.5	3/15	10.5/15	0.0	0.00	0.0
4.0	2/15	8.0/15	0.5	2.25	0.50/15
4.5	2/15	9.0/15	1.0	1.00	2.0/15
5.0	1/15	5.0/15	1.5	2.25	2.25/15
5.5	1/15	5.5/15	2.0	4.00	4.0/15
		52.5/15			17.50/15

$$\text{Mean} = \frac{52.5}{15} = 3.50$$

$$= 3.50$$

$$\text{Variance} = \frac{\sum (\bar{x} - 3.5)^2 P(\bar{x})}{15} = 17.50$$

$$= 1.167$$

Following the same method, we can calculate the Mean and Variance for $n=4$ as

$$\text{Mean} = 3.50 \text{ and Variance} = 0.2917$$

Note the equality of the parent population mean given in Fig 1.4 and the two sampling distribution means as given above. This is not by accident but a property of the method of calculating Sample mean.

This property is that the sample mean is an UNBIASED estimate of the population mean. The Sample Variance calculated above can be derived by using this formula:

$$\text{Sampling Variance} = \frac{N-n}{N-1} \frac{\sigma^2}{n}$$

PP \Rightarrow Population

where N is the pp size, σ^2 is the pp variance and n is the Sample Size.

In the above example, $n=2$, $N=6$, $\sigma^2 = 2.917$

$$\text{Sampling variance} = \frac{1}{n} \left(\frac{\sigma^2}{N-n-1} \right) = \frac{2.917}{5} = 0.5834$$

Also for $n=4$, $N=6$, $\sigma^2 = 2.917$

$$\text{Sampling variance} = \frac{2}{n} \left(\frac{\sigma^2}{N-n-1} \right) = \frac{2}{4} \left(\frac{2.917}{5} \right) = 0.2917$$

* When N is very large, the formula given for Sampling Variance above reduces to σ^2/n .

DEFINITION: The Square root of Sampling Variance is called the Standard Error.

N.B.: The mean of sampling distribution is the same as mean of the population i.e. the Sampling Mean is an UNBIASED value of the pp mean.

SAMPLING AND SAMPLING DISTRIBUTION.

Population and Samples: A major purpose of doing research is to refer or generalize from a Sample to a larger population. This process of inference is accomplished by using statistical methods.

based on probability. Population is the term used to describe a large set or collection of items that has something in common. The Universe of population consists of the total collection of items or elements that fall within a scope of statistical investigation.

The purpose of defining a statistical population is to provide very explicit limits for the data collection process and for the inferences and conclusion that may be drawn from the study.

A Selected in such a way that

A Sample on the other hand is a subject of the

PP Selected in such a way that it is representative of the larger population. The term population and sample are relative. An aggregate of element which constitute a PP for one purpose may merely be a sample for another. E.g. the average age of students in a class.

✓ REASONS FOR SAMPLING

- 1) Sample can be studied more simply than population.
- 2) A study of a sample is less expensive than PP. Since a smaller number of items are examined.
- 3) A study of an entire PP is impossible in most situations hence, sampling may represent the only possible or practicable method to obtain the desired information. For example, in the case of processing such as manufacturing where the universe is conceptually infinite including all features as well as current production. It is not possible to accomplish a complete enumeration of the PP. Also, in destructive sampling of a finite PP. After investigation it is possible to effect a complete enumeration of the PP but, it will not be practical to do so. For example the production of bulbs.
- 4) Samples can be selected to reduce heterogeneity i.e. (very difference).
- 5) Samples results are often more accurate than the result based on PP, since more time and resources can be ^{spent} based on the people who performed the observation and collection of data.
- 6) If samples are properly selected, probability method can be used to estimate an error in resulting statistics. It is this aspect of Sampling that permits investigator to make probability statements about observations in a study.

METHODS OF SAMPLING

Items can be selected from statistical universes of duration in a variety of ways. It is useful to distinguish random from non-random method of selection. The best way to ensure that a sample will yield reliable and valid inferences is to use probability samples or random samples in which the prob. of being included in the sample is known for each subject in population. In a nutshell, this is a definition of RANDOM SAMPLING METHOD.

NON RANDOM SAMPLING METHODS: These are referred to as judgment sampling i.e. selection method in which judgment is exercised in deciding which element of a universe is to be included in the sample. The basic reason random sampling is preferable to non random sampling is that in judgement selection, there is no objective method of measuring the precision or reliability of estimate made from the sample. On the other hand in random sampling, the precision with which estimate of pp values can be made obtainable from ^{the} Sampling Value. ~~is to a~~ This is a very important advantage since random sampling techniques provide an objective basis for measuring errors due to the sampling process and for stating the degree of confidence to be placed upon estimate of pp values.

SAMPLE RANDOM SAMPLING: This is one of the methods of sampling. It has been that a random sample is a sample drawn in such a way that the prob. of inclusion of every elements in the population is known. A ~~is~~ simple random sample of size ' n ' is a sample drawn in such a way that every combination of ' n ' elements has an equal chance of being the sample selected since no practical sampling

Situation involves Sampling without replacement. It is useful to think of this type of sample as one from which each of the 'N' population elements has an equal prob. $\frac{1}{N}$ of being the one selected, the first drawn, $\frac{1}{N}$ of being the one selected, the 2nd one drawn and so on. Until the nth sample has been drawn since there are N^n possible samples of 'n'. The prob. that any sample of size 'n' will be drawn $\frac{1}{N^n}$. e.g. Let

Let $N = 5, n = 2$ the possible samples will be ${}^5C_2 = 10$. Supposing the 5 nos are labeled. Let's take 2 possible samples we have.

METHOD OF SIMPLE RANDOM SAMPLING.

- 1) Drawing Chips from a box: If attention is restricted to the most straight forward situation in which the elements are easily identified and can be numbered. For example, suppose there are 100 students in class and we wish to draw a simple random sampling of 20 of these students without repeating assigning nos from 1-100 to each of the student and place these numbers on physically similar slips of paper which could be placed in a box. Shake the box to accomplish a thorough mixing of the slips then draw the sample. The slip is drawn and the number of it is recorded. Also, the 2nd is drawn and the no of it is recorded and so on. The student corresponding to these 20 numbers constitute the required simple random sampling.
- 2) Tables of Random Numbers: If the population size is very large, the above procedure can become quite unwieldy and time consuming. It may even introduce bias. If the slips are not thoroughly mixed, the tables of random digits can be used to

Select the required samples. These tables are useful.

• Random Number Tables: A table of random digits is a simple

table of digits which have been generated by a random process.

3) Stratified Random Sampling: In stratified random

Sampling, the population is classified into mutually exclusive sub-groups or strata and probability samples are drawn independently for each of these strata.

A sample from each stratum may be obtained by simple random sampling or some other forms of probability sampling system. Carefully combined to together, it covers all estimate of a parameter or they may be compared with one another to reveal between strata differences. It will concentrate on the case where the objective is to obtain an overall estimate of a parameter by combining the result of stratification as compared to simple random sampling to obtain reduction in sampling error or synonymous increase in precision. To do so you reduce your sampling error, you are indirectly increasing precision. The reduction in sampling error

elements which are more alike with respect to the characteristics under investigation. Stratification

is most effective when the elements within strata are homogeneous as possible and difference of elements among strata as great as possible to minimises differences among elements within strata and to maximise differences among strata.

4) Cluster Sampling: This is a technique in which the

Population is sub-divided into groups or clusters, then a probability sample of these clusters is

drawn and studied. The primary aim of Cluster Sampling is to a complete cost savings in sample design. Sampling error is reduced in cluster sampling when the units within clusters are as heterogeneous as possible. For instance, if a single cluster duplicates all of the heterogeneity, which exists in the population, we would only need to draw this cluster into our sample to have a good description of the population.

5.) SYSTEMATIC SAMPLING - You select k^{th} elements of the population until you have all the sample. In the population k can be any no. In many practical applications, Systematic Sampling is used in place of simple random sampling as a method of obtaining random selection. In a systematic sample, every k^{th} elements is drawn from the listed or arranged in some specified manner. The starting point is selected at random from 1st k element. E.g. If a researcher has severally arranged list of farmers and he selects every 5th name on the list where $k = 5$, then a starting point is selected at random; the 1st 5 names and then every 5th name selected after the 1st no. k . Sometimes k is determined by dividing the no. of items in sample frame by the sample size. If the means of a population occurs in random or the other result of the systematic sampling very similar to those of simple random sampling. However, systematic sampling should be avoided whenever there is some periodic or cyclical vary in the order of units in population.

DEFINITION OF TERMS

SAMPLING FRAME/FRAME: A listing of all the elements in the population is called the sampling frame or simply frame.

TYPES OF ERROR

The concept of error is a central one throughout all statistical works. Whenever we have measurement inferences or decision making, the possibility of error is present. There are 2 different types of errors which may be present in statistical measurement, namely:

— Systematic error

— Random error

Systematic errors cause a measurement to be incorrect in some systematic ways. They are of two types which persist even when the sample size is increased. They are errors involved in the procedures of a statistical investigation and may occur in the planning stages or during or after the collection process.

Systematic errors are sometimes called BIASED errors. Among the causes of biased are:

faulty design of a questionnaire such as

surveys with leading questions

refusals by respondent to provide correct information, incomplete Sampling frames, mistakes in planning and processing of data and so on.

These errors may be viewed as arising primarily from inaccuracies or deficiencies in the measuring instruments. These errors can be reduced or completely eliminated.

Random or Sampling errors arise from operation of a large number of uncontrolled factors called chance. For example, if repeated random samples

of the same size are drawn from a population without replacement. A particular statistics such as the mean will differ from sample to sample even if the same definition and procedures are used. These samples tend to distribute themselves below and above the true population parameter. The difference between the mean \bar{x} said to be a random error or a sampling error. The complete collection of factors which could be explained why the sample mean differed from the population mean is unknown and errors decreases in the average.

for this reason that a larger sample of observations is preferred to a smaller one because sample errors are smaller for larger samples. The results are more reliable and more precise.

PARAMETER:> The value of a measure such as the mean, a median or a standard deviation computed from a population is called parameter. The value of such a measure computed from a sample is called statistics. Thus, statistics is derived from a sample and parameter is from population value. For statistics, we use Greek letters & Roman letters for population & sign and not for sample.

FUNCTION:> This is a value by which every member of one set is assigned to or paired with one member of another set. Eg suppose there are 28 students in a seminar and that y is a set of projects to present. If f is a rule that assigns to every element x in the set X a unique element y in the set Y , then f is said to be a function that maps X into Y . On the other hand, a random

Variable β is a function which assigns numerical value to the different outcomes defined by a Sample Space. Don't forget that in the above example the well β "choose a random topic" which is the random variable.

Example 2: Suppose a coin β is tossed twice so that the sample space β x , let x represent the number of heads which comes up so the random variable β is how the number of x obtained in two tosses of a fair coin. Then in this case it's standard x in a sample space and discrete random sample.

SAMPLING DISTRIBUTION

Let us consider how statistic differ from samples to samples. If repeated simple random sampling of the same size are drawn from a statistical population. E.g. A given statistic such as a mean or proportion will vary from one sample to another. The distribution of means can then be examined to estimate the amount of variation that can be from one sample to another. Probability distribution of such a statistic is referred as a Sampling distribution. Thus, we have a Sampling distribution of mean, proportion and so on.

SAMPLING DISTRIBUTION OF MEAN: Let μ be the prob. distribution of some given pp from which we take a sample of size n . Then, the probability distribution of a statistic \bar{x} is called the Sampling distribution of means. We use the following theorems to explain this: