

INFERENTIAL STATISTICS:

It is a way of making inferences about populations based on samples.

- Given information about a subset of examples, how do we draw conclusions about the full set (including other specific examples in that full set)
- Inferential / Inferences based on principles of evidences we use sample statistics.

We will discuss,

1. Probability
2. Distributions and
3. Hypothesis testing

CENTRAL LIMIT THEOREM (CLT) :

(137)

It states that if you have a population with mean ' μ ' and standard deviation ' σ ' and take sufficiently large random samples from the population with replacement, then the distribution of the sample means will be approximately normally distributed.

STEPS :

The steps used to solve the problem of CLT that are either involving ' $>$ ', ' $<$ ' or 'between' are as follows:

1. The information about the mean, population size, standard deviation, sample size and a number that is associated with ' $>$ ', ' $<$ ' or two numbers associated with both values for range of 'between' is identified from the problem.

2. A graph with a centre as mean is drawn.
3. The formulae $Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$ is used to find the Z-score.

4. The Z-table is referred to find the 'Z' value obtained in the previous step.

5.

CASE - i.e CLT ' $>$ '

Subtract the Z-score value from 0.5

CASE - ii: CLT ' $<$ '

Add 0.5 to the Z-score value

CASE - iii: CLT 'BETWEEN'

Step 3. is executed

6. The Z-value is found along with \bar{x} .

The last step is common to all the three cases, that is to convert the decimal obtained into a percentage.

(139)



- find mean/average salary
- Entire population of India
- 1.3 Billion.
- Inferential statistics + μ_{sal}

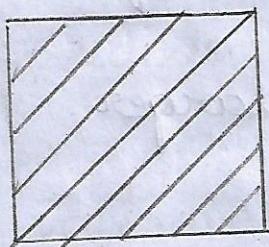


$$\mu_{\text{sal}} = ?$$

Step-1: collect the data from all 1.3 Billion.
(no need to worry about cost or time)

Step-2:
$$\frac{\sum_{i=1}^{1.3 \text{ BILLION}} \text{sal}_i}{1.3 \text{ Billion}} = \mu_{\text{sal}}$$

Step-3: Using inferential stats, cost & time
gets reduced.



$$\downarrow$$

$$\mu_{\text{sal}}$$

10K
DATA
POINTS

$$\frac{\sum_{i=1}^{10000} \text{sal}_i}{10,000}$$

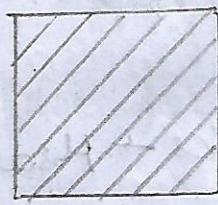


$$\text{mean of sample} = \bar{x}$$

μ = mean of population

σ = standard deviation of population

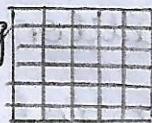
s = standard deviation of sample.



(140)

1 Lakh

Sampling



\bar{x} (sample mean)



μ_{height}

$\rightarrow \mu_{\text{height}} - \bar{x}$

$\rightarrow \mu \approx \bar{x} \rightarrow \text{point estimate}$

\therefore Sampling space is biased.

SAMPLING:

It is a process used in statistical analysis in which a predetermined number of observations are taken from a larger population.

The methodology used to sample from a larger population depends on the type of analysis being performed, but it may include simple random sampling or systematic sampling.

SAMPLING TECHNIQUES:

(21)

1. CONVENIENT SAMPLING:

It is a type of non-probability sampling in which people are sampled simply because they are "convenient" sources of data for researches.

↳ In probability sampling, each element in the population has a known non-zero chance of being selected through the use of a random selection procedure.

2. VOLUNTEER SAMPLING:

It is a sampling technique where participants self-select to become part of a study because they volunteer when asked, or respond to an advert.

3. RANDOM SAMPLING:

It is a part of the sampling technique in which each sample has an equal probability of being chosen.

- A sample chosen randomly is meant to be an unbiased representation of the total population.
- An unbiased random sample is important for drawing conclusions.

4. STRATIFIED SAMPLING:

In a stratified sample, researchers divide a population into homogeneous subpopulations called strata (the plural of stratum) based on specific characteristics like race, gender, location, etc.

- Every member of a population should be in exactly one stratum.

(142)

* Random sampling and stratified sampling are used frequently.

CONFIDENCE INTERVAL:

A confidence interval refers to the probability that a population parameter will fall between a set of values for a certain proportion of times.

* This measures the degree of uncertainty or certainty in a sampling method.

^RINTERPRETATION OF CONFIDENCE INTERVAL:

The proper interpretation of a confidence interval is probably the most challenging aspect of this statistically concept.

The most common interpretation of the concept is as follows:

There is a 95% probability that, in the future, the true value of the population parameter (e.g., mean) will fall within the interval between x [lower bound] and y [upper bound].

In addition, we may interpret the confidence interval using the statement below:

We are 95% confident that the interval between x [lower bound] and y [upper bound] contains the true value of the population parameter.

However, it would be inappropriate to state the following:

There is a 95% probability that the interval between x [lower bound] and

(145)

$\hat{Y}[\text{upper bound}]$ contains the true value of
the population parameter.