# Sampling Distribution

Sampling distribution is a statistic that determines the probability of an event based on data from a small group with in a large population.

Since the population is too large to analyze, you can select a smaller group and repeatedly sample & analyze them. The gathered data, & statistic is used to calculate the likely occurance (or) probability of an event

Each random sample selected may have a different value assigned to the statistic being studied.

For ex:- if you randomly sample data three times & determine mean (or the average of each sample, all three means are likely to be different & fall somewhere along the graph. That's variability.

bell curve →

→ The number Observed in a population (N)

→ The no. Observed in the Sample (n)

→ The method of choosing the Sample

## Types of distribution

① Sampling distribution of mean

② The most common type of Sampling distribution is the mean. It focuses on calculating the mean of every Sample group chosen from the population & plotting the data points. The graph shows a normal distribution where the center is the mean of the Sampling distribution, which represents the mean of the entire population.

② Sampling distribution of Proportion.

This Sampling distribution focuses on Proportions in a population you select Samples & calculate their Proportions. The means of the Sample proportions from each group represent the Proportion of the entire population.

## ③ T-distribution

A T-distribution is a sampling distribution that involves a small population. It is used to estimate the mean of the population & other statistics such as intervals, statistical differences & linear regression.

## Sampling & Resampling

## Sampling

Sampling is a process of selecting group of observations from the population to study the characteristics of the data to make conclusions about the population.

Probability          Non-Probability

# Probability Sampling (Random Sampling)

In this type, data is randomly selected so that every observations of population gets the equal chance to be selected for Sampling.
- Simple Random Sampling
- Cluster Sampling
- Stratified "
- Systematic "

## Non Probability Sampling

In this type, data is not randomly selected. It mainly depends upon how the Statistician wants to Select the data
- convenience
- judgmental/purposive
- Snowball/Referral
- Quota

## Sampling error :- Errors which occur during Sampling Process are known as SE(∞)

Difference b/w Obs value of a Sample Statistics & the actual value of a population parameters

## Advantage of Sampling

Reduce Cost & Time
Accuracy of Data
Less resource needed

## Resampling

→ Resampling is the method that consist of drawing repeatedly drawing samples from the population.

K-fold cross-Validation          Bootstrapping

In this method population data is divided into k equal sets in which one set is considered as the test set for the experiment while all other set will be used to train the model

## Bootstrapping

:- In bootstrapping samples are drawn with replacement (i.e One Observation can be repeated in more than One group & the remaining data which are not used in samples are used to test the model.

## Statistical Inferences (or) Inferential Statistics

Statistical Inference is a method of making decisions about the parameters of a population, based on random sampling. It helps to assess the relationship b/w. dependent & independ Var

The purpose of SI to estimate the uncertainty (or) sample to sample variation

Components ⟵ Sample size
Variability in the Sample
Size of the Observed differences

$$\sqrt{\frac{1}{(N)} \sum (Y_i - \hat{Y}_i)^2}$$

# Prediction Error

In prediction error refers to the diff b/w the predicted values made by some model & the actual values

→ LINEAR REGRESSION
→ LOGISTIC "

Linear regression — used to predict the value of some continuous response variable

We typically measure the prediction error of a linear regression model with a metric known as RMSE

$\hat{y}_i$   $y$   $RMSE = \sqrt{\sum(\hat{y}_i - y_i)^2 / n}$ ——— sample size

14   12

Predicted value    Observed value

Logistic Regression — used to predict the value of some binary response variable.

Prediction error of Logistic Reg model is with a metric known as total misclassification rate

$$\text{Total Misclassification rate} = \frac{\#\ \text{incorrect prediction}}{\#\ \text{total Pre}}$$