

## "UNIT-01"

SPI DATE: \_\_\_/\_\_\_/\_\_\_  
PAGE: \_\_\_

### ■ Introduction to Big Data Platform:

- ⇒ Data is nothing but facts and statistics stored or free flowing over a network, generally it's raw and unprocessed.
- ⇒ When data are processed, organised, structured or presented in a given context so as to make them useful, they are called "information".

### ■ Data:

- ⇒ The quantities, characters or symbols on which operations are performed by a computer.
- ⇒ Stored in form of magnetic signals.
- ⇒ Recorded on magnetic, optical or mechanical recording media.

### ■ Big Data:

- ⇒ Big Data is also data but with a huge size.
- ⇒ Collection of data that is huge in size and yet growing exponentially with time.
- ⇒ None of the traditionally data management tools are able to store it or process it efficiently.

### ④ Importance:

- i) Big data can bring dramatic cost reduction
- ii) Substantial improvements in time required to perform a computing task.
- iii) New product and service offerings.

### ④ Example: New York Stock Exchange Good Write

1T new data generated every day.

## Types of Big Data:

- Structured
- Unstructured
- semi-structured.

## Difference between Traditional Data and Big Data:

Traditional Data	Big Data
→ Deals with structured Data.	→ Deals with structured, unstructured or semi-structured data.
ii) Volume ranges from 1GB to 1.Terabyte).	ii) Volume ranges from Petabytes to Zettabytes.
iii) Generated per hour or per day or more.	iii) Generated per second.
iv) Centralized and managed in centralized form.	iv) Distributed and managed in distributed form.
v) Data integration is very easy.	v) Data integration is very difficult.

## ■ Structured Data:

- ⇒ Any data that can be stored, accessed or processed in form of fixed format is termed as "structured data".
- ⇒ An "Employee" table in a database is an example of Structured Data.

## ■ Unstructured Data:

- ⇒ Any data with unknown form or the structure is classified as unstructured data.
- ⇒ A heterogeneous data source containing a combination of simple text files, images, video etc.
- ⇒ Example of unstructured data : the output returned by "Google Search".

## ■ Semi- structured data:

- ⇒ Contains both forms of data.
- ⇒ Actually not defined with a table definition in relational DBMS.
- ⇒ Example = a data represented in an XML file.

Good Write

## Three Characteristics

Difference between Structured, Unstructured and Semi-structured data:

Structured	Unstructured	Semi-structured
i) It is based on Relational database table.	i) It is based on characters or binary data.	i) It is based on XML / RDF.
ii) Versioning over tuples, rows, tables whole.	ii) Versioned as a graph.	ii) Versioning over tuples or graph is possible.
iii) It is schema dependent and less flexible.	iii) It is more flexible and there is absence of schema.	iii) More flexible than unstructured and less flexible than unstructured data.
iv) Allows complex joining.	iv) Only textual queries possible.	iv) Queries over anonymous nodes are possible.
v) Matured transaction	v) No transaction management	v) Transaction is adapted from DBMS.

## Three Characteristics of Big Data:

- i) Volume → Data quantity.
- ii) Velocity → Data Speed.
- iii) Variety → Data Types.

## Challenges Faced by conventional systems:

- i) Large quantities - impossible to process.
- ii) Meaningful and collected in real time.
- iii) Multiple sources.
- iv) Collect correct data.

**Conventional System:** System consists of one or more zones each having either manually operated call points or automatic detection devices or combination of both.

### ④ Challenges:

- i) Uncertainty of Data Management Landscape.
- ii) Big Data Talent Gap.
- iii) The talent gap that exists in the industry and getting data into the big data platform.
- iv) Need for synchronization across data sources.
- v) Getting important insights through the use of big data analytics.

### ④ Other three challenges:

- i) Data
- ii) Process
- iii) Management

Good Write

## Challenges of Big Data:

- i) Meeting the need for speed.
- ii) Visualization helps organizations perform analyses.
- iii) The degree of granularity increases.
- iv) Understanding the data.
- v) Addressing data quality.
- vi) Displaying meaningful results.
- vii) Dealing with outliers.

## Intelligent Data Analysis: (IDA)

IDA: An interdisciplinary study concerned with the effective analysis of data.

- Used for extracting useful information from large quantities of online data.
- Extracting desirable knowledge or interesting patterns from existing database.

→ At a level of abstraction higher than the data, and information on which it is based and can be used to deduce new information and new knowledge.

\* Goal: Extract useful knowledge, the process demands a combination of extraction, analysis, conversion, classification, organization, reasoning, and so on.

## ① Uses of TDA:

- i) Data visualization
- ii) Data pre-processing
- iii) Data engineering
- iv) Database mining techniques, tools & applications
- v) Big data applications.
- vi) Evolutionary algorithms.
- vii) Machine Learning.
- viii) Neural Nets.
- ix) Fuzzy logic
- x) Statistical pattern recognition.
- xi) Post-processing

## ② Importance of TDA:

- i) Decision making
- ii) Data processing
- iii) Epidemiological study.
- iv) Multi dimensionality of problems is looking for methods for adequate and deep data processing and analysis.

## ■ Illustration of TDA by using Sec5:

⇒ application.names : lists the classes to which cases may belong and the attributes used to describe them.

→ application.data : provides information on the training cases from which Sec5 will extract patterns.

⇒ application-test : provides information on the test case.

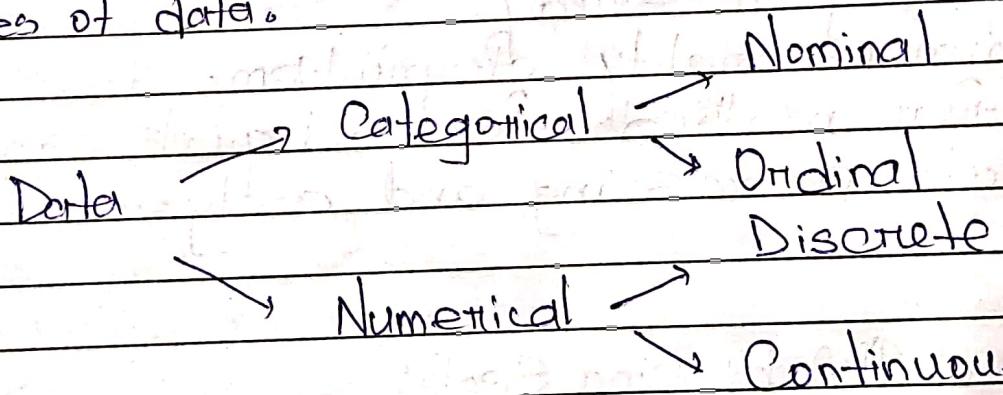
■ Nature of data:  
⇒ Set of values of qualitative or quantitative variables; in short, pieces of data are individual pieces of information.

⇒ Measured, collected and reported and analyzed

① Properties of Data ⇒

- i) Amenability of use
- ii) Clarity
- iii) Accuracy
- iv) Essence
- v) Aggregation
- vi) Compression
- vii) Refinement.

■ Types of data:



① Nominal: The characteristics of nominal are as follows

② Order: The order of the response of observation does not matter.

③ Distance: Nominal scales do not hold distance. The distance between a 1 and 2 is no the same as a 2 and 3.

④ True zero: There is no true or real zero.

In nominal, zero is uninterrupted.

Good Write

- ④ Appropriate statistical form nominal scale:  
mode, count and frequencies.
- ⑤ Display: histograms or bar chart.

- ii) Ordinal Scale: The characteristics are as follows.
- ⑥ Order: The order of the response on observation matters.
  - ⑦ Distance: Ordinal scales do not hold distance. The distance between first and second is unknown as is the distance between first and third and so on along with all observations.
  - ⑧ True zero: There is no true or real zero. An item, observation or category cannot finish zero.
  - ⑨ Appropriate statistics: mode, mean and frequencies.
  - ⑩ Displays: histogram or bar charts.

iii) Interval Scale: Classical interval scale are Likert scale (eg 1 - strongly agree and 5 - strongly disagree).

Semantic interval scale

are (1 - dark and 9 - light).

The characteristics are as follows -

- ⑪ Order: The order of response on observation does matter.
- ⑫ Distance: Interval scales do not offer distance.

⑬ True zero: There is no zero. However, data can be rescaled in a manner that contains zero.

- ④ Appropriate statistics: count, frequencies, mode, median, mean, standard deviation, skewness and kurtosis.
- ⑤ Displays: histograms or bar chart, line charts and scatter plots.
- iv) Ratio scale: Ratio scale appears as nominal scale with a true zero.  
The characteristics are as follows:
  - ① Order: The order of the response of observation matters.
  - ② Distance: Have an interpretable distance.
  - ③ True Zero: There is true zero.
  - ④ Appropriate statistics: count, frequencies, mode, median, mean, standard deviation, skewness and kurtosis.
  - ⑤ Displays: histograms or bar chart, line charts and scatter plots.

## ■ Analytic process and tools:

There are 6 analytic process:

### i) Deployment:

→ We deploy the result of analysis.

→ Monitoring and maintenance.

→ Produce a final report and review the project.

### ii) Business understanding:

→ Determine the business objective.

→ Assess the situation.

→ Determine the data mining goals.

→ Produce the project plan.

### iii) Data exploitation:

→ Need to gather initial data, describe and explore the data and verify data quality to ensure it contains the data we require.

→ Data collected from the various sources is described in terms of its application.

### iv) Data preparation:

→ Need to select data as per need, clean it, construct it to get useful information.

→ Integrate it all.

### v) Data Modeling:

- Select a modeling technique, generate test design, build a model and assess the model built.
- Analyze relationships between various selected objects in the data.
- Test cases are built for assessing the model and model is tested and implemented.

### vi) Data evaluation:

- At the end of analysis, we evaluate the quality of the measurements and results comparing them to our original design criteria.

## ■ Analysis and Reporting:

■ Analysis = The process of exploring data on reports in order to extract meaningful insights which can be used to better understand and improve business performance.

■ Reporting = The process of organizing data into informational summaries in order to monitor how different areas are performing.

Good Write

## Difference between Analysing and Reporting:

<u>Analysing</u>	<u>Reporting</u>
i) Provides what is needed	i) Provides what is asked for.
ii) Typically customized.	ii) Typically standardized.
iii) Involves a person.	iii) Does not involve a person.
iv) Extremely flexible.	iv) Fairly flexible.
v) Transforms data and information into insights.	v) Translates raw data into information.
vi) Interprets the data at a deeper level and providing actionable recommendations.	vi) Helps companies to monitor their online business and alert them to when data falls outside of expected range.

## Modern Data Analytic tools:

- i) Hadoop : helps in storing and analyzing data.
- ii) MongoDB : used on datasets that change frequently.
- iii) Talend : used for data integration & management.
- iv) Cassandra : a distributed database used to handle chunks of data.

Good Write

- v) Spatio: used for real-time processing and analysing large amount of data.
- vi) STORM: an open source real-time computational system.
- vii) Kafka: a distributed streaming problem that is used for fault-tolerant storage.

## Statistical Concepts:

- ⇒ Data consists of information coming from observations counts, measurements or response.
- ⇒ Statistics is science of collecting, organizing, analysing and interpreting data in order to make decisions.
- ⇒ Population is the collection of all outcomes, responses, measurements, or counts that are of interest.
- ⇒ Statistics is a piece of data from a portion of a population.

### ② Basic statistical operations:

- i) Mean: A measure of central tendency for quantitative data i.e. the long-term average.
- ii) Median: A measure of central tendency for quantitative data i.e. the half-way point.
- iii) Mode: The most frequently occurring (discrete) value or where the probability density function peaks.

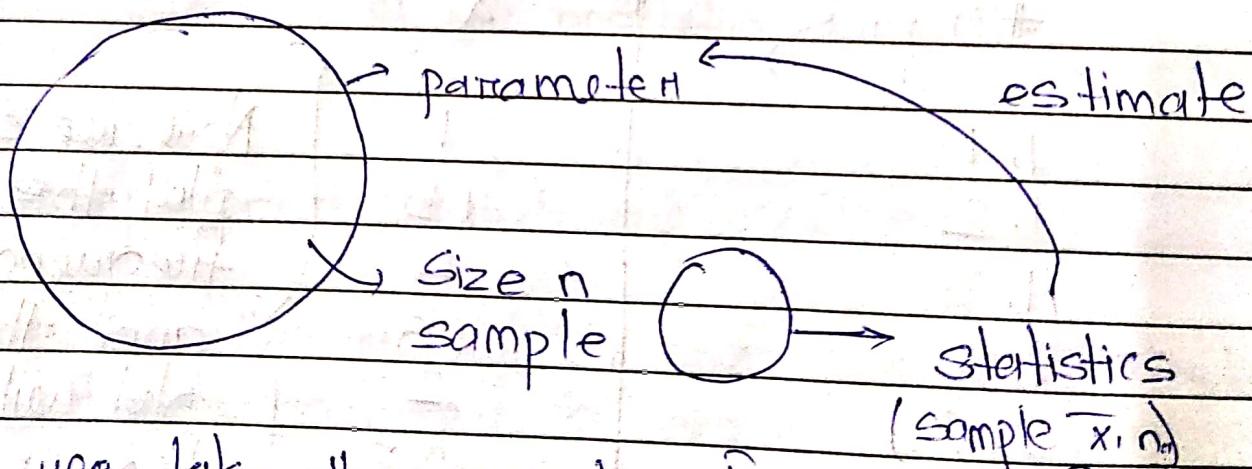
## ■ Sample Distribution:

Sample: A sample is "a smaller" (but hopefully representative) collection of units from a population used to determine truths about that population.

### ① Types of Samples $\Rightarrow$

- i) Probability Samples: Systematic random sample  
Stratified random sample  
Multistage sample  
Multiphase sample  
Cluster sample!

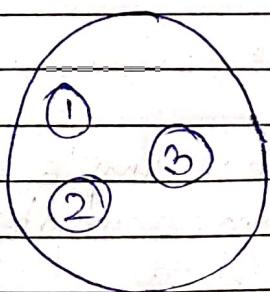
- ii) Non-probability Samples: Convenience sample  
Purposive sample  
Quota



Here, we take the sample of size  $n$  from a population; calculate the statistics and estimate the parameter of the population.

The distribution / frequency w/ which we can get different values for the statistics that is trying to estimate this parameter is called "Sampling distribution".

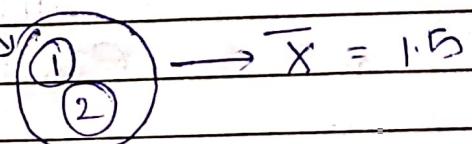
For example,



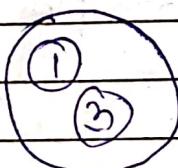
Population parameter here,  $\mu = \frac{1+2+3}{3} = 2$

Three balls population

① Taking samples like



$$\bar{x} = 1.5$$



$$\bar{x} = 2$$

likewise,

#'s pick

$\bar{x}$

1, 1

1

Now, we can

1, 2

1.5

plot the

1, 3

2

frequencies

2, 1

1.5

and the

2, 2

2

plot will be

2, 3

2.5

known as

3, 1

2

"Sampling

3, 2

2.5

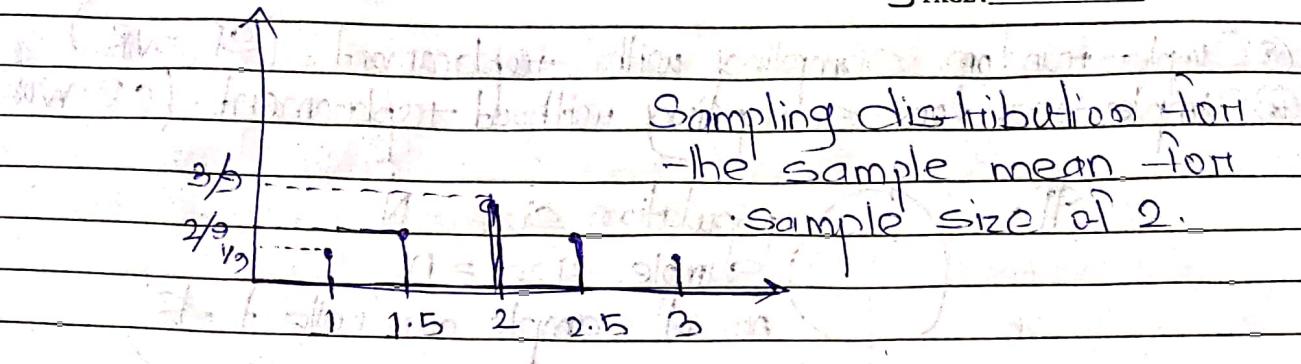
distribution"

3, 3

3

of the sample means.

Good Write



### ■ Specific types of samples:

- i) Stratified sample: Ensures that each segment from the population is represented.
- ii) Cluster sample: All members from randomly selected segments of a population.
- iii) Systematic sample: A sample in which each member of the population is assigned a number.
- iv) Convenience sample: Consists only of available members of the population.

### ■ Sampling distributions are used to:

- calculate the probability that sample statistic could have occurred by chance.
- to decide whether something that is true of sample statistic is true of the population.
- come to conclusions about a population

- ① Simple random sampling with replacement (SRSWR)  
 ② Simple random sampling without replacement. (SRSWOR)

PF, population size =  $N$

Sample size =  $n$

No. of sample can collect =  $\binom{N}{n}$

For SRSWR,  $k = N^n$

For SRSWOR,  $k = {}^N C_n$

## II Standard error:

Population size =  $N$

mean =  $\mu$

Standard deviation =  $\sigma$

Variance =  $\sigma^2$

$n$	$n$	$n$	$n$
$x_1$	$x_2$	$x_3$	$x_k$
$s_1$	$s_2$	$s_3$	$s_k$
$s_1^2$	$s_2^2$	$s_3^2$	$s_k^2$

sample      Sample 2      sample k

## Sampling distribution of mean

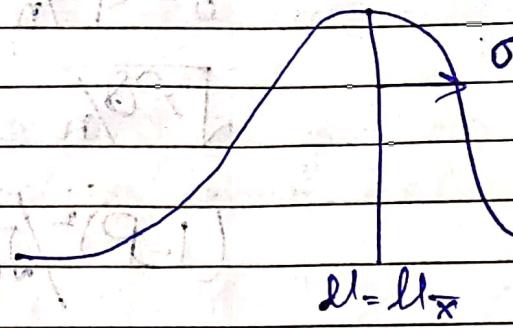
$x_1$	$P(\bar{x}_1)$
$\bar{x}_2$	$P(\bar{x}_2)$
$\bar{x}_k$	$P(\bar{x}_k)$

Good Write

→ Standard deviation of sampling distribution of mean  
 is "Standard Error" of mean.

④ mean of sampling distribution of mean,  $\mu_{\bar{x}} = \mu$  (population mean)

⑤ standard deviation of sampling distribution of mean or standard error of mean,  $\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$



$$\sigma/\sqrt{n}$$

$\sigma$  = standard deviation of population  
 $n$  = sample size

Uses:

i) Setting up "Confidence Interval"

ii) To set the "Confidence Interval"

↳ (the mean of estimate plus and minus with a factor of the variation in that estimate).

iii) Test of Significance

↳ (the procedure for comparing an observed data with a claim the truth of which being assessed).

# Statistics | areas for unit 1 in Standard Edition

- ① Sample mean,  $\bar{X}$  follows  $\bar{X} \sim N(\mu, \sigma^2/n)$
- ② Sample standard deviation,  $s$  follows  $s \sim \sqrt{\sigma^2/n}$
- ③ Sample variance,  $S^2$  follows  $S^2 \sim \frac{\sigma^2}{n-1}$
- ④ Sample proportion,  $P$  follows  $P \sim \sqrt{P(1-P)/n}$
- ⑤ Sample correction,  $r$  follows  $r \sim (1-P)^2/n$
- ⑥ Difference of two sample mean ( $\bar{X}_1 - \bar{X}_2$ ) follows  $\bar{X}_1 - \bar{X}_2 \sim \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

- The Central Limit Theorem says, "for random sampling, as the sample size  $n$  grows, the sampling distribution of  $\bar{X}$  approaches a normal distribution."
- The sampling distribution will be normal no matter what the population distribution's shape as long as  $n > 30$ .

## Re-sampling:

- The method that consists of drawing repeated samples from the original data samples.
- It is a non-parametric method of statistical inference.
- Re-sampling uses experimental methods, to generate the unique sampling distribution.
- Re-sampling is performed to get additional information about model fitness.

example: M M P P P T

→ Finding standard error and confidence interval

1	2	1	2	3	4	5	6	7	8	9	10
2	3	1	2	3	4	5	6	7	8	9	10
3	4	2	3	4	5	6	7	8	9	10	
4	5	3	4	5	6	7	8	9	10		
5	6	4	5	6	7	8	9	10			
6	7	5	6	7	8	9	10				
7	8	6	7	8	9	10					
8	9	7	8	9	10						
9	10	8	9	10							
10		9	10								

→ It checks the consistency of the model statistics.

⇒ Consistency of model parameters.

⇒ Checks the confidence of the model so that the model is not overfitted for the data.

⇒ Consistency of model performance.

⇒ Checks if there is high deviation; then

there is something wrong with the model.

Good Write

## Re-sampling methods:

~~Re-sampling~~ There are four major re-sampling methods available:

1) Permutation: "Fisher's tea taster"

- ④ 8 cups of tea are prepared
    - four with tea poured first
    - four with milk poured first
  - ⑤ The cups are presented to her in random order.
    - let's say T T T T M M M M

## Solutions

- ⇒ Make a deck of eight cards, four marked "T" and four marked "M".
  - ⇒ Deal out these 8 cards successively in all possible orderings.
  - ⇒ Record how many of those permutations show  $\geq 6$  matches.
  - ⇒ Repeat many times.

## Algorithm:

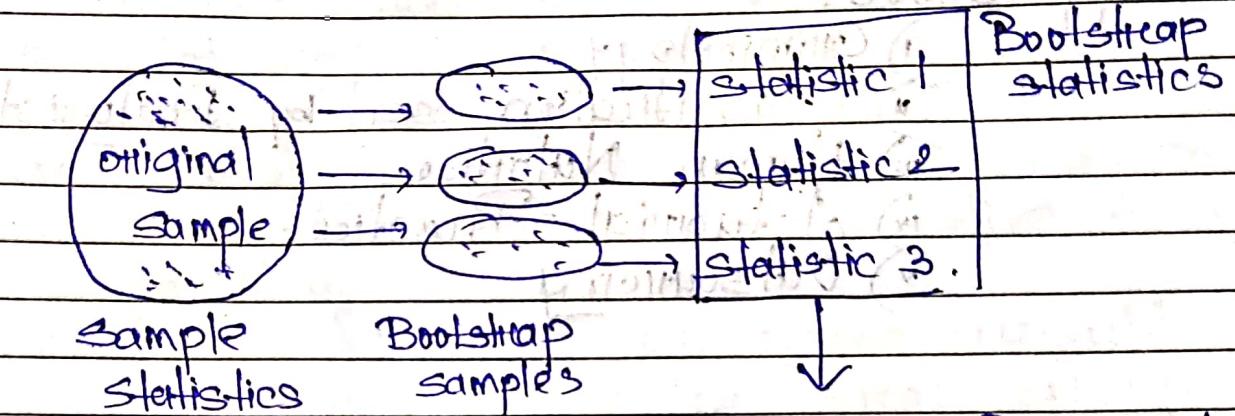
- Step 01: Collect Data from Control & Treatment group
  - Step 02: Merge sample to form pseudo population.
  - Step 03: Sample without replacement from pseudo population to simulate control Treatment
  - Step 04: Compute target statistics for each example.

~~for all players with  
affinity for him.~~

## *Good Write*

i) **Bootstrap:** Bootstrap generates distinct data sets by repeatedly sampling observations from the original dataset.

→ These generated data sets can be used to estimate variability in lieu of sampling independent data sets from the full population.



**Algorithm:** A bootstrap sample is a random sample conducted with replacement.

- Step 01: Randomly select an observation from the original data.
- Step 02: "Write it down"
- Step 03: "Put it back" (i.e. any observation can be selected more than once).
- Step 04: Repeat (step 1-3)  $N$  times;  $N$  is the number of observations in the original sample.

**R Final Result:** One "bootstrap sample" with  $N$  observation.

Good Write

- ⇒ Bootstrapping makes no assumption regarding the population.
- ⇒ No normality of error terms.
- ⇒ No equal variances.
- ⇒ Allows for accurate forecasts of intermittent demand.
- ⇒ Applications:

- i) Criminology
- ii) Classification used by Ecologists
- iii) Human Nutrition
- iv) Actuarial Practice
- v) Outsourcing

### Difference between Parametric Bootstrap and non-Parametric Bootstrap.

#### Parametric Bootstrap vs Non-Parametric

i) Generates its samples from the (assumed) distribution of the data using the estimated parameter values.

i) Generates its sample by sampling with replacement from the observed data.

→ If the assumed distribution is not true, it will give wrong results.

→ If the assumed distribution is true, it will give correct results.

→ If the assumed distribution is not true, it will give wrong results.

→ If the assumed distribution is true, it will give correct results.

## Comparison between Bootstrap and Jackknife

- Bootstrap can be viewed as closely related method of Jackknife
- Bootstrap is 10 times more computationally intensive than Jackknife
- Bootstrap is conceptually simpler than Jackknife.
- Bootstrap performs well than Jackknife.
- Jackknife is more conservative than Bootstrap.
- Jackknife produces same result each time while Bootstrap gives different result.

## iii) Cross Validation: A technique used to protect against overfitting in a predictive model.

- In cross validation, we make a fixed number of folds (or partitions) of data, run the analysis on each fold, then average the overall error estimate.
- It is not a re-sampling technique.
- Requires large amount of data.
- Useful in data mining and AI.

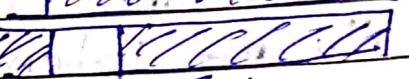
## Types of Cross Validation:

### i) Leave One Out Cross validation:

Suppose, we have 1000 records (dataset).

Test Train

For exp 1, 

For exp 2, 

Test Train

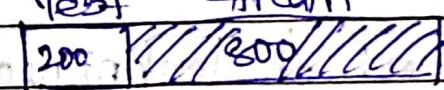
Check each data, so we have to do many iterations and it will have small error but low bias.

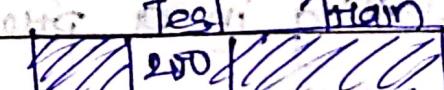
### ii) K-Fold Cross validation:

Suppose, we have 1000 records.

Consider, the random value of  $K=5$

Test Train

Acc=1  $\leftarrow$  For exp 1, 

Acc=2  $\leftarrow$  For exp 2, 

Iterations will be minimum of all the folds. Iterations will be equal to  $K$ .

The minimum accuracy getting in this model will be the minimum of all the accuracy and

The maximum accuracy getting in this model will be the maximum of all the accuracy.

## ii) Stratified Cross validation:

No of instances of  
oil of Potassium or manganous each class should  
be proportionate.  
There should be  
proportionate amount  
of "Yes" "No" / "0,1"  
data in test set or  
training set.

(x) Yes & No & (x) 0,1  
Test set

## Statistical Inference:

① **Inference:** Use a random sample to learn something about a large population.

② **Statistical inference:** The process of making guesses about the truth from a sample.

③ **Classification:**  
 Statistical inference is divided into:  
 Estimation  
 Testing  
 Hypothesis testing  
 Least squares method  
 Good Write

## ■ Types of Statistical Inference:

### i) Confidence Intervals:

Range of values =  $m$  is expected to lie within 95% confidence interval.  
 = 0.95 probability that  $m$  will fall within range  
 = probability is the level of confidence.

Finding:

$$x - z_{cv}(\sigma_x) \leq m \leq x + z_{cv}(\sigma_x)$$

Lower limit                              Upper limit

### ii) Test of Significance (Hypothesis Testing):

→ Sample data to evaluate a hypothesis about a population parameter.

→ A hypothesis is an assumption about the population parameter.

→ The technique is introduced by considering a "One-sample Z test".

#### ④ Hypothesis Testing Steps:

- Null and alternative hypothesis.
- Test statistic
- P-value and interpretation
- Significance level.

**Null Hypothesis:** A type of statistical hypothesis that proposes that no statistical significance exists in a set of given observations.

**Difference between Null Hypothesis and Alternative Hypothesis:**

Null Hypothesis	Alternative Hypothesis
-----------------	------------------------

i) It denotes there is no relationship between two measured phenomena.	i) It's a hypothesis that random cause may influence the observed data.
ii) It is represented by $H_0$ .	ii) It is represented by $H_1$ .
iii) Statement example $\Rightarrow$ Rohan will win at least Rs. 10000 in lucky draw.	iii) Statement example $\Rightarrow$ Rohan will win less than Rs. 10000 in lucky draw.

**Prediction Error:** Prediction error is the failure of some expected event to occur.

$\text{Prediction error} = \text{actual outcome} - \text{predicted outcome}$

$y_i = \hat{y}_i + e_i$

Good Writing  $\rightarrow$  good marks

④ Bias =  $\frac{\sum E_i}{\text{no. of predictions}}$  |  $E = \text{errors}$

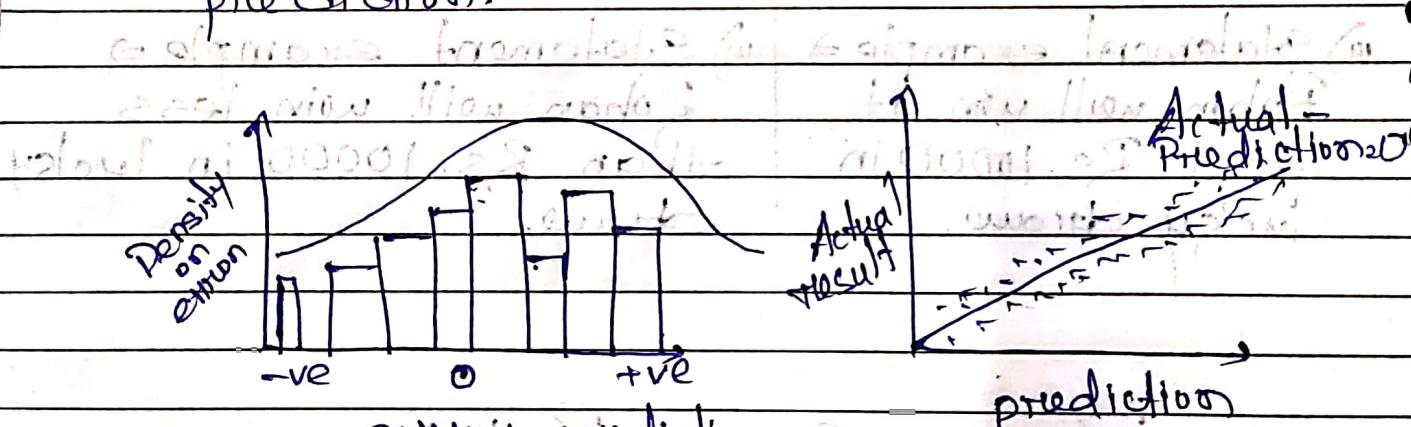
Explanation: no. of predictions

⑤ Root Mean Square Error

$$= \sqrt{\frac{\sum (\text{prediction error})^2}{\text{no. of predictions}}}$$

→ Bias tells us which direction our predictions are going. It's not good for the magnitude of the prediction.

→ RMSE gives us the magnitude of the prediction.



Regressions differing in accuracy of prediction. The standard error of the estimate is a measure of the accuracy of predictions. The regression line is the line that maximizes the sum of squared Good Write deviations of prediction.