

# 1. Introduction to Big Data

## 1. What is Big Data

- ① Big data is a term that is used to describe data that is high volume, high velocity, and high variety, requires new technologies and techniques to capture, store and analyze it.
- ② Big data refers to data that is so large, fast or complex that it's difficult or impossible to process using traditional methods.

## 2. What is Unique about big data?

- ① Companies have searched for make the best use of information to improve their business capabilities.
- ② Big data is also special because it represents both significant information which can open new doors to improve business.

## 3. Enlist tools Used to Big data

- 1. NOSQL - Database MongoDB, CouchDB, Cassandra, Bigtable,
2. Mapreduce - Hadoop, Hive, Pig, Sqoop, Kafka.
3. Storage - S3, Hadoop distributed file system.
4. Server - ElastiCLOUD, Google app engine, Elastic.
5. Processing - R, Yahoo! Pipes, Bigsheets

Above are the some tools used for big data.

→ 4. Write Advantages and disadvantage of big data

+ Advantages:

1. Big data helps in improving Science and research;
2. It improves healthcare
3. Every second additions are made.
4. One platform carries Unlimited information
5. It helps in financial trading's, sports, polling etc;
6. It helps in Optimizing business process

+ Disadvantages:

1. lot of Big data is Unstructured
2. it may increase Social stratification
3. it can be used for manipulation of public records
4. Big data is not useful in short run.
5. Traditional storage can cost lot of money to store big data;

5. Explain 5 'V's of Big Data.

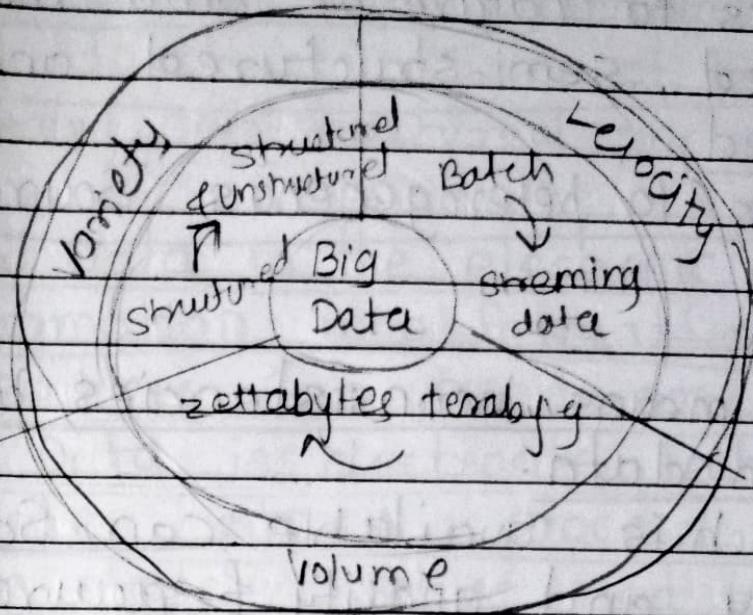
→ I. Volume      Volume

II. Variety      Variety

III. Velocity      Velocity

IV. Veracity      Veracity

V. Value      these are 5 V's of big data  
                    value



### I. Volume

- ① • The name 'Big Data' itself is related to a size.
- ② • Volume is huge amount of data.
- ③ • To determine the value of data, size of data plays a very crucial role.
- ④ • If the volume of data is very large then it is actually considered as a 'Big data'.

### II. Velocity

- ① • Velocity refers to the high speed of accumulation of data.
- ② • In Big data Velocity data flows in forms Source like machine, networks, Social media, phones etc.
- ③ • Example - There are more than 3.5 billion searches per day are made on google.

### III. Variety

- ① • It refers to Diversity of data.
- ② • Variety means different formats of data from various sources.

- ① • It refers to nature of data that is Structured, semi-structured and Un-structured,
- ④ • It refers to heterogeneous Sources.

#### IV. Veracity

- ① • Veracity means inconsistencies and Uncertainty in data.
- ② • Data which is available can sometimes get messy and quality of accuracy are difficult to control.
- ③ • Example - Data in bulk could create confusion.

#### V. Value

- ① • After having 4 V's into there comes one more v which is stands for Value,
- ② • The bulk of data having no value unless you turn it into something useful.
- ③ • Data needs to converted into something valuable to extract information.
- ④ • Value is the most imp v of all 5 v's.
- Value - Extract useful data.

1. Volume

5. Value

2. Velocity

4. Variety

3. Veracity

2. Variety

6. What is Digital Data? OR Explain types of digital data OR Types of big data or diff? Structured, semi-structured and unstructured

- Digital Data is the electronic representation of information in a format or language that machines can read or understand.
- Digital Data is electronic technology that generates, stores, process data in terms of two states: positive and non-positive!

\* Types of Digital data and Big data.

### I] Structured data -

- ① This data is basically an Organized data.
- ② It generally has defined the length and format of data.
- ③ Structured data has standurdized format, has a well defined structure.
- ④ Structured data is easily Understood by machine language.
- ⑤ example - Relational database, name, dates, address, and more.

### II] Unstructured data -

- ① This data is not-organized data.
- ② It refers to data that generally doesn't fit neatly into the traditional row and column structure of the relational database.
- ③ Unstructured data doesn't store in structured database.

④ Word, PDF, text, media logs, picture etc. are the examples of Unstructured data which can't be stored in the form of rows and columns.

### III] Semi-structured data

- ① This data basically semi-organized data
- ② It generally refers to data that has defined the length and format of data do not confirm to the formal structure of data.
  - \* Log Files, XML data examples of semi-structured data
- ③ Semi-structured data does not follow the Format of tabular data model because it does not have fixed schema.

7. What is big data analytics & explain types of big analytics.

- 
- 7 Big data analysis is the process of extracting useful information by analysing different types of big data sets;
  - 7 Big Data analytics is used to discover hidden patterns, market trends, and consumer preferences for the benefits of organizational decision making.
  - 7 There are four types of big data analysis:
    - ① Descriptive analysis
    - ② Predictive analysis
    - ③ Exploratory & discovery
    - ④ Prescriptive analysis

## I] Descriptive Analytics

- Descriptive analytics is a statistical method that is used to search and summarize historical data in order to identify patterns or meaning.
- The simplest way to define descriptive analytics is that it answers the question "What has happened?"
- Example - Company reports provide historic review

## II] Predictive analytics.

- Predictive analytics is a statistical method that utilizes algorithms and machine learning to identify trends in data and predict future behaviour.
- It has ability to "Predict" what might happen.
- This analytics are about Understanding future.
- Example - Training targets, talent management

## III] Prescriptive analysis.

- Prescriptive analysis is a statistical method used to generate recommendation and make decisions based on computational findings of algorithmic model.
- It answers the question "what should we do?"
- It focuses on advise on possible outcomes.
- Example - A training manager use prescriptive analysis.

## 8. Diff' Between structured, Unstructured and Semistructured:

Properties	structured data	unstructured data	Semi-structured data
organization	well organized	not organized at all	Partially organized
Transaction	Matured transaction	No Transaction	Transaction is adapted from DBMS, not matured.
Flexibility	less Flexible	more Flexible	more flexible than Structured data but less flexible than Unstructured data.
Scalability	Very difficult to scale	Very Scalable	simpler to scale than structured data.
Technology	It is based on relational database	It is based on character and binary data	It is based on XML / RDF.
Example	Financial data, Bar codes, Relational database.	Media logs, Video, Audio, Files	Tweets organized by hashtags, folder organized by topics.

9. What is meant by Petabyte and zettabyte?

• Petabyte :-

A petabyte is a measure of memory or data storage capacity that is equal to  $2^{50}$  bytes.

• Zettabyte :-

A zettabyte is a unit of measurement used by technology professionals and the general public to describe a computer or other device's storage capacity.

10. Inlist Big data Applications.

- 1. Healthcare
- 2. Telecommunication
- 3. Financial Firms
- 4. Retail
- 5. Energy & Utilities
- 6. Education
- 7. Media & Entertainment
- 8. Law Enforcement
- 9. Marketing
- 10. New Product development
- 11. Banking
- 12. Insurance
- 13. Agriculture

11. What are the sources of Big data?

- Big data has many sources. includes:
- Social media sources such as Twitter, Facebook generate tremendous amounts of comments and tweets.
- Machines such as Smart phones, meters, generate data.
- Image, Voice and audio data.

## 2. Introduction to Data Science

1. What is Data Science? Write down Components of data Science:

- Data Science analyzes large amount of data
- ③ Data Science is the study of data to extract meaningful insights for business;
- ④ Data science means extracting useful knowledge for data to solve business problems.
- ⑤ R-language is most important for data science.
- ⑥ Data Science is the domain of study that deals with vast volume of data using modern tools and techniques to find meaningful information.

\* Components of data science

1. Statistics - It is most important component. it is way to collect and analyze the numerical data in large amount.
2. Visualization - It represents data in visual context so that people can understand. it is easy to access the huge amount of data in visuals.
3. Data engineering - It is part of data science which involves acquiring, storing, retrieving and transforming data.

4. Advanced Computing = Heavy lifting of data science is advanced computing. it involves designing, writing, and maintaining the source code of computer programs.

5. Machine learning - It all about to provide training to machine so that it can act as a human brain.

2. Write Down Advantages and disadvantages of data Science.

• Advantages :-

1. Data Science helps organizations to reduce costs, get into new markets.
2. The domain is in high demand.
3. it offers abundance of job positions.
4. it helps in securing highly paid career.
5. it helps in achieve high objectives of business.
6. It is versatile branch of data.

• Disadvantages :-

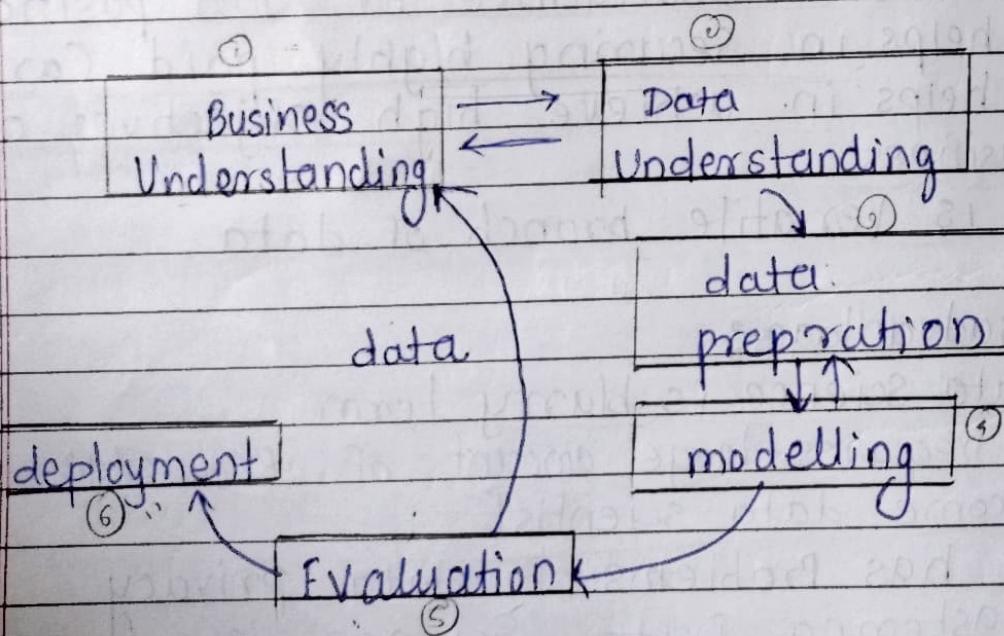
1. Data Science is blurry term.
2. it requires large amount of knowledge to become data scientist.
3. It has problems of data privacy.
4. Mastering Data Science is nearly impossible.
5. Tools used for data Science are more expensive.

3. Write down applications of Data Science

- 1. Fraud and Risk Detection
- 2. Healthcare
- 3. Virtual assistance for patients and Customer Support
- 4. Internet Search
- 5. Targeted advertising
- 6. Gaming
- 7. Website recommendations.

4. Explain Data Science Process.

- 1. Following are the process phases and stages involved in data Science process.



Data Science Process Stages.

- 3. Data Science process is systematic approach to solving a data problem. it provides a structured framework.

## 1. Business Understanding -

- The first step of this process is Setting a research goal.
- This stage has main purpose that making sure all the users understand the what, how and why of the project.

## 2. Data Understanding:-

- The second phase is Data Understanding that is data retrieval.
- You want to have data available for analysis.
- This step includes Finding suitable data.

## 3. Data Preparation -

- This phase or stage involves Data preparation and proceeds with activities in order to get familiar with data.
- In this phase data is designed, stored, modeling, formating involved.

## 4. Modeling -

- In this phase various modeling techniques are selected and applied.
- This phase comes after data preparation phase.

## 5. Evaluation -

- At this stage in the project, you have built a model that appears to have high quality.
- It is important to evaluate the model.

## 6. Deployment -

- The last step or stage of data Science process model is presenting your results and automating the analysis.
- Deployment is last phase of Data Science process;

## 5. What is Data Analytics?

- ① Data analytics is the Science of examining raw data with the purpose of drawing conclusions about that information.
- ② Data Analytics is a lifeline for the IT industry.
- ③ Data analysis is a process of inspecting, cleansing, transforming, and modeling data with the goal of discovering useful information and supporting decision making.

(4)

## 6. Need of Big Data analytics.

- Big data analytics gives future plan for organization & organization to take decision based on Big. data Analysis

That's why Big data analytics is needed and more important.

## 7. Write down advantages & disadvantages of big data analytics.

→ Advantages -

- 1) Fraud detection
- 2) Delivering Relevant Products
- 3) Optimizing & improving customer Experience.

Disadvantage:

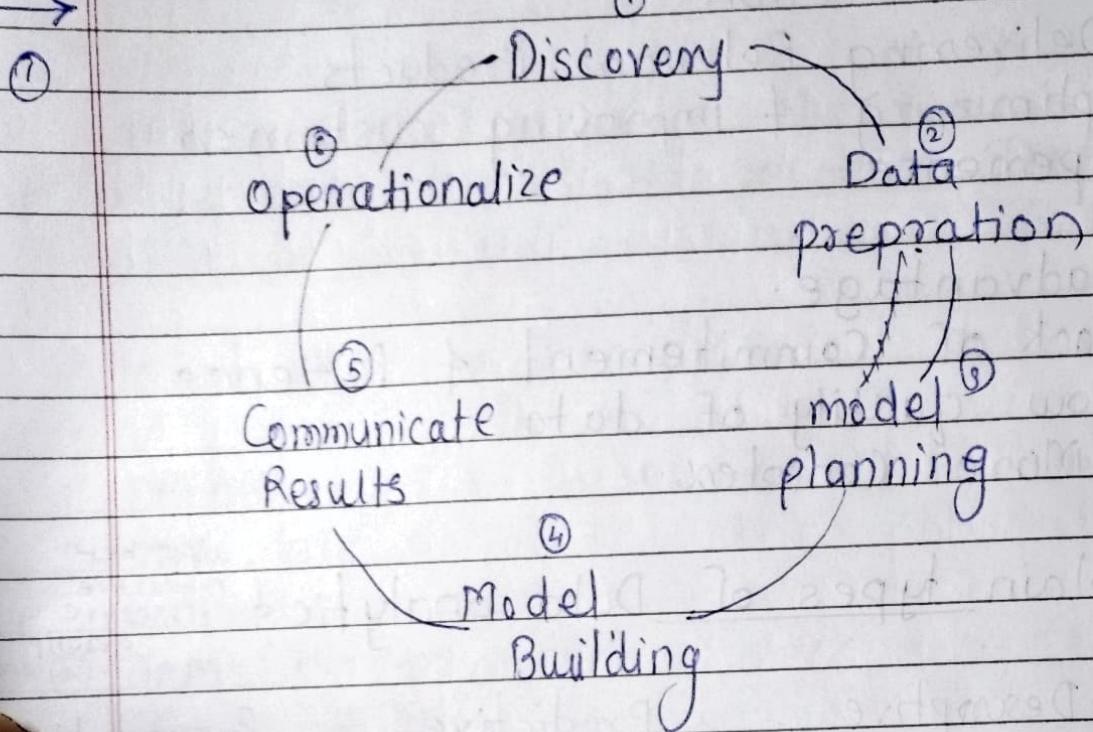
- 1) Lack of Commitment & Patience
- 2) Low quality of data
- 3) More Complex.

### 8. Explain types of Data analytics

<sup>diff' betn</sup>  
<sup>predictive</sup>  
<sup>prespective</sup>  
<sup>descriptive</sup>  
analys

Properties	Descriptive	Predictive	Prespective
Defn.	Summarizes historical data	It is the use of data to analyzes data and to identify hidden patterns.	Predict the future events provides possible outcome
Focus on	Insights into the past	Understanding the Future	Advise on possible outcomes
Answers the question	What has happened?	What might happen?	What should we do?
Tools Used	Data aggregation, Data mining	Statistics, modeling simulation	Business rule machine learning optimization
Example	Demand trends Survey results	Training targets, talent management	Training manager

g. Explain Phases of Data analytics lifecycle.



### ② Phase-1 - Discovery

- Team learns business domain.
- data Science team is trained and researches the issue
- The team comes with an initial hypothesis, which can be later confirmed with evidence.

### ③ Phase-2 Data Preparation

- This phase involves analysing, preparing data and modeling
- It requires to have an analytic sandbox
- Data preparation tasks can be repeated and not in predefined sequence
- Some tools used commonly for this process include - Hadoop, Apache Flink, Open Refine etc.

#### ④ Phase-3 Model Planning -

- In this phase, team determines the methods, techniques, and workflow of data.
- In this phase data science teams create data sets that can be used for training for testing, production and training goals.
- Some tools used in this stage are MATLAB and STASTICA.

#### ⑤ Phase-4 Model building -

- In this phase, team creates datasets for training, testing as well as production use.
- In this phase the team builds and executes model based on work done in the model planning phase.
- Tools used are free and open-source octave, R and PL/R etc.

#### ⑥ Phase-5 Communication Results-

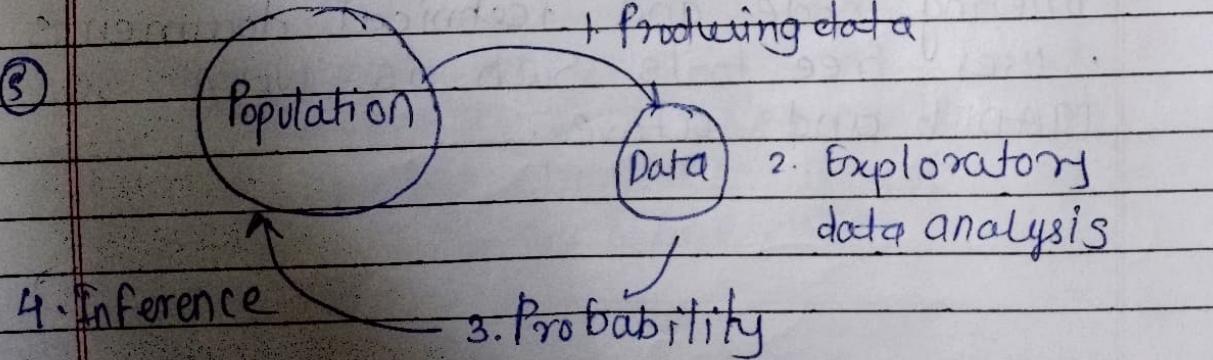
- In this phase, stakeholders determine if the results of project are success or failure based on the criteria.

#### ⑦ Phase-6 Operationalize

- In this phase team delivers final reports, briefing, code and technical documents
- Uses free tools such as WEKA, SQL MADlib and octave.

10. What is statistical Inference? Explain with example +

- ① Statistical Inference is the process of Using data analysis to deduce properties of an Underlying distribution of probability
- ② Statistical inference is method for drawing and measuring the reliability of conclusion about population based on info obtained from a Sample of the population.
- ③ Statistical Inference involves four-step process that is:
  - i. Data production
  - ii. Exploratory Data analysis
  - iii. Probability
  - iv. Inference
- ④ Statistical Inference includes two broad areas:
  - i. Statistical Estimation
    - Point Estimation
    - interval estimation
  - ii. Statistical Hypothesis Testing



## 11. What is Population Explain

- 
- ① Population refers to the collection of all elements that ~~have~~ possessing common characteristics that Comprises Universe.
  - ② Population is defined as a total of the items Under Consideration.
  - ③ Population focuses on Identifying characteristics.
  - ④ Example - automobiles with 4 wheels, people who ~~consume~~ olive oil.
  - ⑤ Types of population :-
    - i. Finite Population
    - ii. Infinite Population
    - iii. Existent population
    - iv. Hypothetical population

## 12. What is Sample.

- 
- ① Sample is a part of population chosen at random.
  - ② Sample is that part of the population from which information is collected.
  - ③ Sample is defined as proportion of population selected.
  - ④ Sample focus on Making inference about population.
  - ⑤ Example - 20 cars from each Make.

## 13. Explain Statistical & Statistical Modeling

- 
- ① Statistical modeling is an approach to statistical data analysis that helps researchers discovers something about a

## Steps of Statistical modeling:-

- ④ ① model selection ② Model fitting ③ model visualization.

Phenomenon that is assumed to exists

- ② This statistical modeling approach helps explain the variability found in the dataset
- ③ Example - Public Health data, Social media data

### 14. Define Probability

- ① Probability theory developed from the study of games of chances like dice and cards
- ② Process like flipping a coin, rolling a die or drawing a card from a deck are called probability experiments
- ③ Probability theory is the foundation for statistical inference
- ④ Probability = The no. of ways of achieving success / The total no. of possible outcomes

### 15. Probability distribution

- ① The probability distribution of a random variable  $x$  is the system of numbers.
- ② There are two types of probability distribution
  - a. discrete Probability distribution
  - b. Continuous Probability distribution

### 16. What is Correlation & types of Correlation

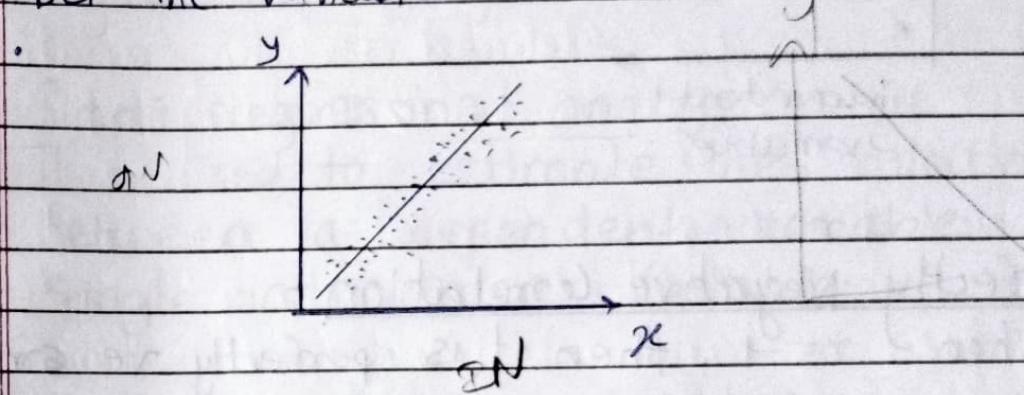
- ① Correlation is a statistical measure which determines association of two variables.

① Correlation is used to represent linear relationship between two variables.

② Types of correlation

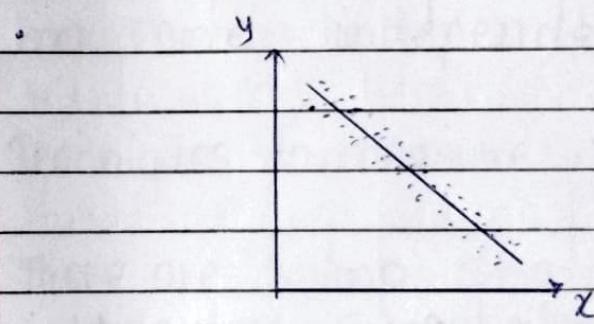
a) Positive Correlation -

- positive correlation indicates positive association bet<sup>n</sup> the variables.



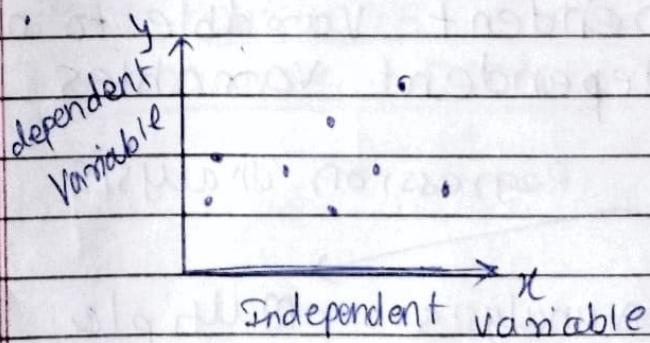
b) Negative Correlation -

- It indicates -ve association bet<sup>n</sup> the variables.



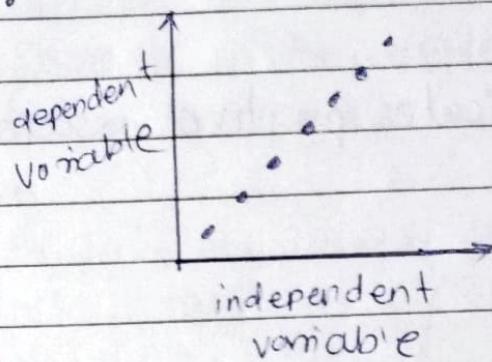
c) No Correlation

- In the no relation bet<sup>n</sup> variables.

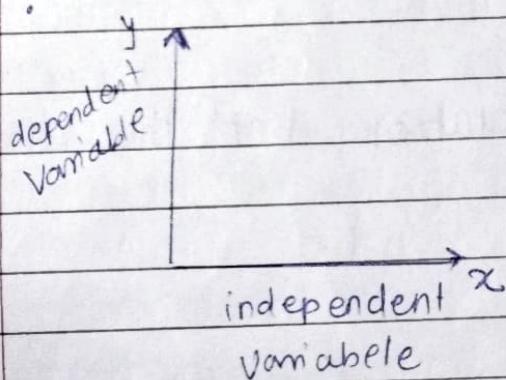


d) Perfectly positive Correlation-

- when  $r=1$ , then it is perfectly +ve correlation

e) Perfectly Negative Correlation-

- When  $r=-1$ , then it is perfectly -ve correlation



## 17. What is Regression &amp; Regression analysis.

- ① Regression is a statistical technique that relates a dependent variable to one or more independent variables.

## ② Type of Regression analysis

Simple Regression analysis  
example. relationship  
bet<sup>n</sup> crop yields &  
rainfall

multiple  
example. relationship  
bet<sup>n</sup> salaries of  
employee & their  
experience of education

③ The statistical technique that expresses a functional relationship between two or more variables in the form of an equation to estimate the value of a variable, based on the given value of another variable is called "regression analysis".

④ Simple regression analysis -

→ it is used to estimate the relationship between a dependent variable and a single independent variable.

⑤ Multiple regression analysis

→ it is used to estimate the relationship between a dependent variable and one or more independent variables.

\* Techniques to measure Correlation.

→ There are 3 imp statistical tool used to measure Correlation.

a) Scatters Diagrams

b) Karl Pearson's Coefficient of Correlation

c) Spearman's rank correlations;

# 3. Introduction to Machine Learning

1. What is machine learning & write down need of machine learning.  
→
  1. Machine learning is defined as the capability of a machine to imitate intelligent human behavior;
  2. Machine learning is subset of artificial intelligence, which focuses on using statistical techniques to build intelligent Computer System to learn from available database;
  3. Machine learning is a scientific discipline that is concerned with the design and development of algorithms that allows computers to evolve behaviors based on available data from database.
4. Need of Machine learning-
  - Diversity of data - data is being generated from diff' channels and its natures are diff?
  - Capacity and dimensions - There is large number of data available led to growth of data
  - Speed - As data volume increases, speed at which data is captured and transformed.
  - Complexity - With increase in Complexity of data, high data quality and security is required.
  - Applicability - It means applicability of data to business and performance improvement.

2. Write down advantages and disadvantages of machine learning.

• Advantages-

1. ML is used in variety of applications such as banking and financial sector, healthcare etc.
2. ML has capability to handle multi-dimensional and multi-variety data.
3. ML allows time cycles reduction and efficient utilization of resources.
4. The process of automation of tasks is easily possible.
5. Due to ML, there are tools available to provide quality improvement in complex data.

• Disadvantages-

1. Acquisition of relevant data is the major challenge.
  2. It is impossible to make immediate accurate predictions with a machine learning system.
  3. Machine learning needs a lots of training data for future prediction.
  4. Interpretation of results is also a major challenge to determine effectiveness of machine learning algorithms.
  5. Possibility of High error.
- Advantage- ML is used for data analysis and knowledge discovery in data base, data mining, etc.

### 3. Introduction to R language & Write features of R language

- ① R is a programming language and used software environment for statistical analysis, graphic representation and reporting.
- ② R was created by Ross Ihaka and Robert Gentleman at the University of Auckland, New Zealand and currently developed by the R Development Core Team.
- ③ Features of R :-
  - R is well developed, simple and effective programming language
  - R has an effective data handling and storage facility
  - R provides graphical facilities for data analysis

④ R language includes various components such as :-

- 1. Attributes
- 2. Data types
  - a. Vectors
  - b. List
- 3. Data frame
- 4. Control structures i.e if . else,  
For, while etc.

### 3. Introduction to R language & write features of R language

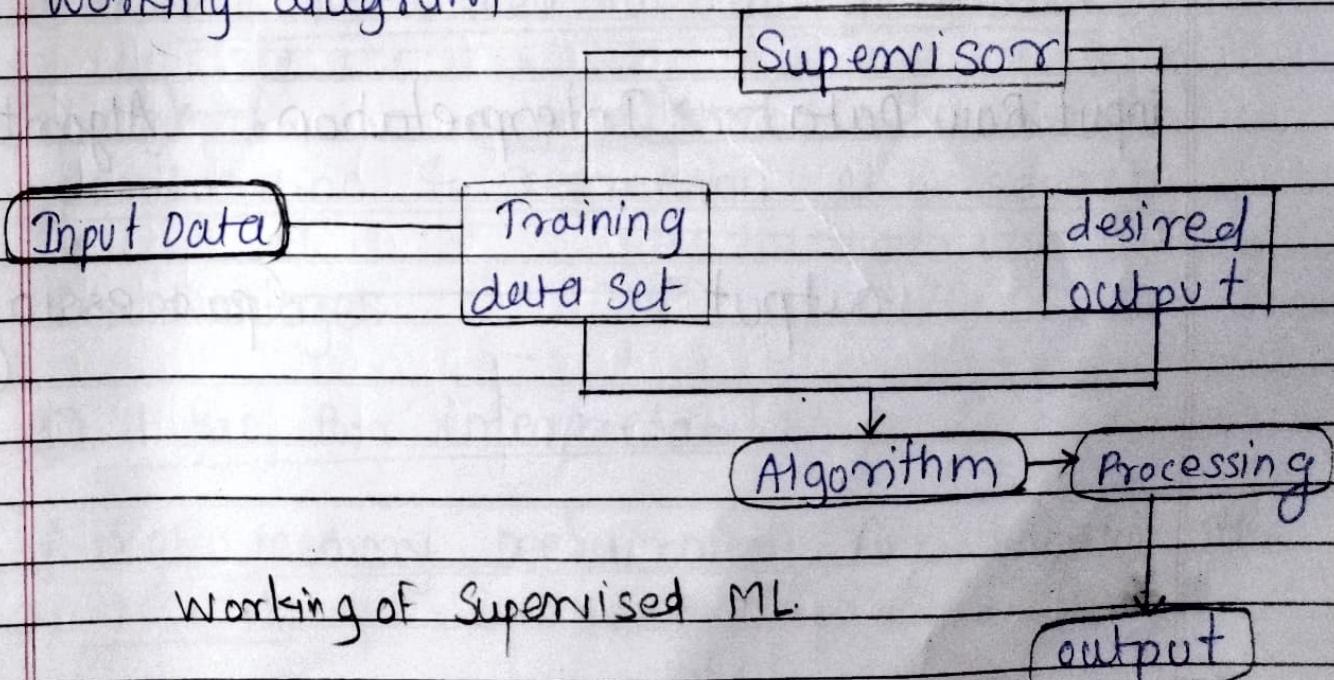
- ① R is a programming language and software environment for statistical analysis, graphic representation and reporting.
- ② R was created by Ross Ihaka and Robert Gentleman at the University of Auckland, New Zealand and currently developed by the R Development Core Team.
- ③ Features of R :-
  - R is well developed, simple and effective programming language.
  - R has an effective data handling and storage facility.
  - R provide graphical facilities for data analysis.
- ④ R language includes various components such as :-

1. Attributes
2. Data types
  - a. Vectors
  - b. List
3. Data frame
4. Control structures i.e if...else,  
For, while etc.

4. Explain two main type of Machine learning OR

\* I. Supervised Machine Learning

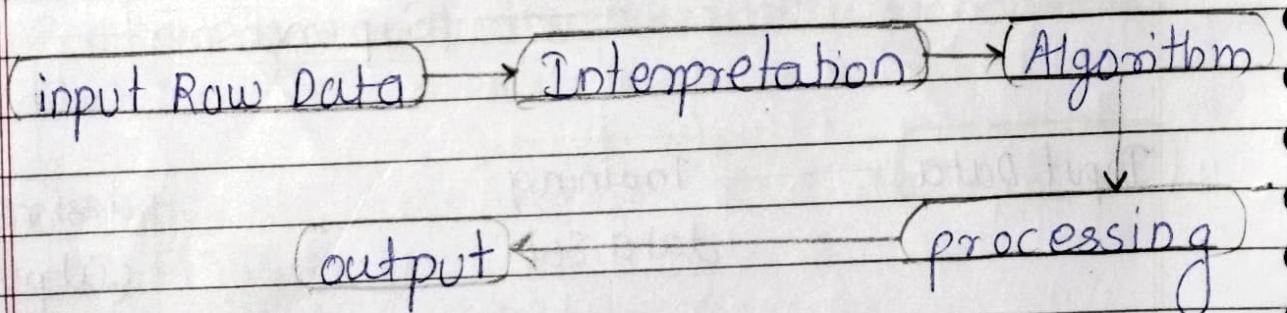
- ① Also known as Predictive analysis.
- ② In Supervised Machine learning, the training is controlled by external Supervisor.
- ③ Common use of Supervised Machine learning is to use historical data to predict statistical future events;
- ④ In Supervised learning, Computer is provided with example input that are labeled with their desired output;
- ⑤ Supervised learning is the type of machine learning in which machines are trained using well "labelled" training data, and on basis of that data, machines predict the output;
- ⑥ Supervised learning Used to train machines so Can develop predictive data model;
- ⑦ Working diagram



- ⑧ SML described as "learn from the past and predict the future!"

## \* II. Unsupervised Machine learning

- ① In unsupervised learning there is no supervisor to observe the process.
- ② In unsupervised learning, the network is provided with inputs but not with desired output.
- ③ In unsupervised learning data is Unlabeled.
- ④ Unsupervised learning is used to create recommendation system.
- ⑤ Following are the some examples of unsupervised machine learning applications
  - a) Social network analysis
  - b) Image Segmentation
  - c) Google News
- ⑥ Working diagram of Unsupervised learning



## 5. short Note on K-Nearest Neighbors (KNN)

- ① K-NN stands for K-Nearest Neighbors.
- ② KNN is one of the Simplest ML algorithm.
- ③ KNN based on Supervised learning algorithm.
- ④ KNN does not learn anything in training period hence KNN is called as lazy learner.
- ⑤ KNN model is used for classification as well as regression:
  - a. In Case of classification, new data points get classified in a particular class.
  - b. In Case of regression, new data gets labeled based on the avg. value of KNN.
- ⑥ Advantages:
  - ① Algorithm is simple and easy to implement.
  - ② There is no need to build a model.
  - ③ algorithm is versatile, that is if used for classification & regression also search.
- ⑦ Disadvantages:

- ① It has Poor interpretability.
- ② High memory requirement for storing the results.

## 6. Naive Bayes Algorithm (short note)

- ① Naive Bayes algorithm is a Supervised learning algorithm;
- ② Naive Bayes is a classification technique based on Baye's Theorem.
- ③ Naive Bayes algorithm model is easy to build & particularly useful for all very large data sets;
- ④ The principle behind Naive Bayes theorem is Bayes theorem or Bayes Rule.
- ⑤ Application of Naive Bayes Algorithm
  - a. it is used for credit scoring
  - b. it is used in medical data classification
  - c. it is used in text classification such as spam filtering etc.
- ⑥ Baye's Theorem
  - a. also known as Bayes' Rule or Bayes' law.
  - b. Bayes' law is used to calculate the conditional probability:
  - c. Bayes theorem represented as :

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} = \frac{P(B|A) P(A)}{P(B) P(A) + P(B) P(\bar{A})}$$

### ⑦ Advantages

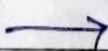
- .. it is easy and fast to predict class of test data sets,

2. it also performs well in multi-class prediction
3. it is effortlessly trainable, even with a small available data set;

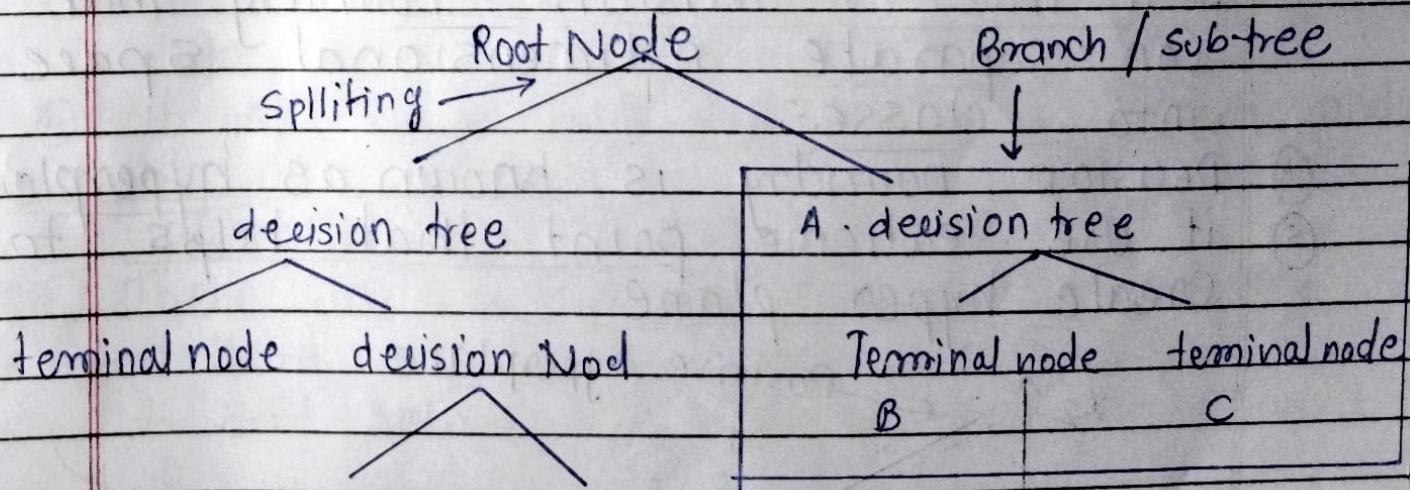
### ⑧ Disadvantages

1. Partially dependent
2. assumes independence of features

### 7. short note on decision tree.



- ① A decision tree also called as prediction tree.
- ② A decision tree uses a tree structure to specify sequences of decisions (conditions)
- ③ It is called decision tree because, structure is similar as tree.
- ④ decision tree starts with root node & then expands its branches.
- ⑤ decision tree simply asks a question based on answer is (Yes / No).
- ⑥ structure of decision tree :-



Terminal node terminal node.  
A is parent node of B and C

## 6. Naive Bayes Algorithm (short note)

- ① Naive Bayes algorithm is a Supervised learning algorithm;
- ② Naive Bayes is a classification technique based on Baye's Theorem.
- ③ Naive Bayes algorithm model is easy to build & particularly useful for all very large data set.
- ④ The principle behind Naive Bayes theorem is Bayes theorem or Bayes Rule.
- ⑤ Application of Naive Bayes Algorithm
  - a. it is used for credit Scoring
  - b. it is used in medical data classification
  - c. it is used in text classification such as spam filtering etc.
- ⑥ Baye's Theorem
  - a. also known as Bayes' Rule or Bayes' law.
  - b. Bayes' law is used to calculate 'the conditional probability'
  - c. Bayes theorem represented as :

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

### ⑦ Advantages

- .. it is easy and fast to predict class of test data set.

2. it also performs well in multi-class prediction
3. it is effortlessly trainable, even with a small available data set.

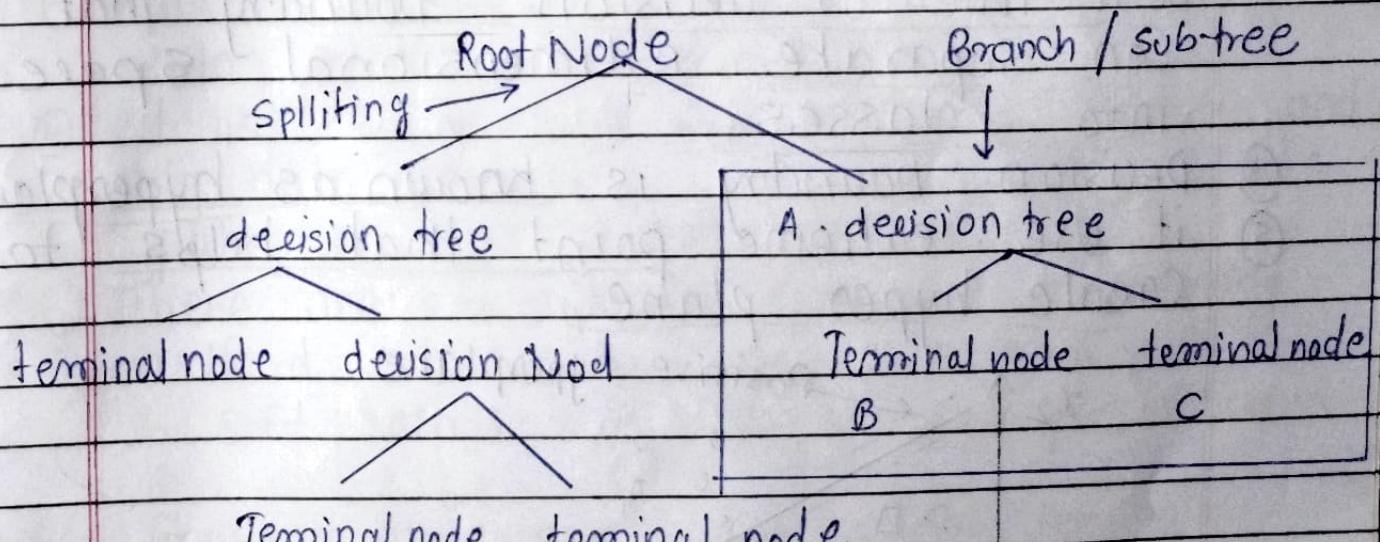
### ⑧ Disadvantages.

1. Partially dependent
2. assumes independence of features

### 7. short note on decision tree.

→

- ① A decision tree also called as prediction tree.
- ② A decision tree uses a tree structure to specify sequences of decisions (conditions)
- ③ It is called decision tree because, structure is similar as tree.
- ④ decision tree starts with root node & then expands its branches.
- ⑤ decision tree simply asks a question based on answer is (Yes / No).
- ⑥ Structure of decision tree :-



A is parent node of B and C

### ⑦ Advantages-

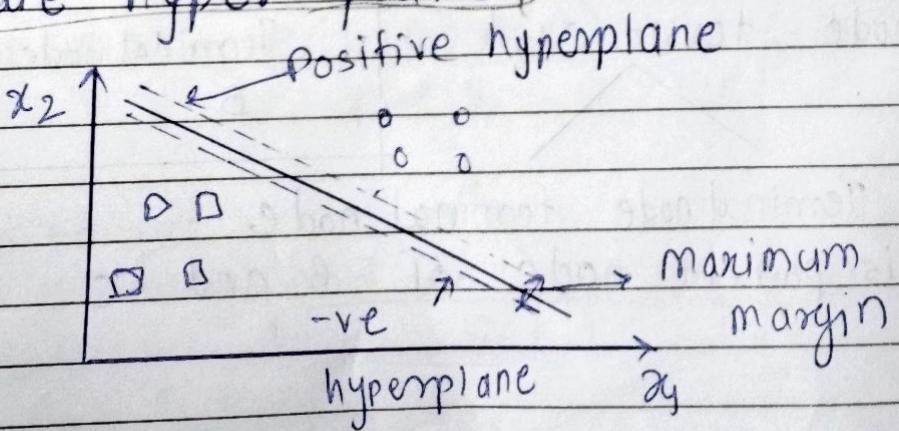
- a. Simple to Understand and interpret
- b. Able to handle both numerical and Categorical data.
- c. requires little data preparation.
- d. performs well with large data set
- e. fast and accurate.

### ⑧ Disadvantages

- a. Trees can be very non-robust.
- b. Poor Performance.
- c. Expensive.

### 8. Write short Note on SVM.

- ① SVM stands for Support Vector Machine
- ② SVM is one of the most popular supervised learning algorithm.
- ③ The goal of SVM algorithm is to create best line or decision boundary that can separate n-dimensional space into classes.
- ④ Decision boundary is known as hyperplane
- ⑤ it uses extreme point that helps to create hyper plane.



### ⑥ Advantages:

1. it works well with clear margin separation.
2. it is effective in high dimensional spaces.
3. it uses subset of training points.

### ⑦ disadvantages

1. SVM does not directly provide probability.
2. When data set has more noise it does not perform very well.
3. When we have large data set then SVM have training time higher.

## g. Cluster Analysis

- ① cluster analysis is a data one of the best data analysis technique,
- ② cluster analysis is a method of Unsupervised machine learning,
- ③ cluster analysis is a task of dividing data points into a number of groups,
- ④ it is very useful for data mining and big data because it automatically finds patterns in data,
- ⑤ There are two main type of clustering
  - a. Hard clustering eg. k-means.
  - b. Soft clustering. eg. Fuzzy C-Means

Grouping the data items such that each item only assigned one cluster

Method of grouping data items such data item can exist in multiple cluster

10. Write short note on K-Means.

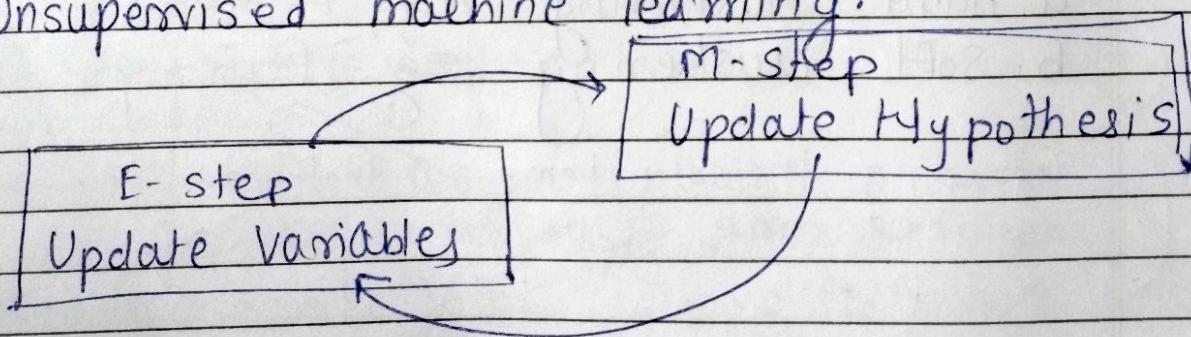
- 
- ① K-means clustering is one of the Simplest and popular Unsupervised machine learning algorithm
  - ② This algorithm involves telling the algorithms how many possible cluster there are in the dataset.
  - ③ K-means is a partitional clustering algorithm.
  - ④ Advantage
    - 1. simple easy to understand and to implement
    - 2. it provide easily interpretable clustering results.
    - 3. Fast and efficient to use

⑤ disadvantage

- 1. algorithm is only applicable if the mean is defined
- 2. do not work properly for large amount of data points.

11. short note on EM-Algorithm.

- 
- ① EM stands for Expectation-Maximization algorithm
  - ② it is defined as the Combination of Various Unsupervised machine learning.
  - ③



## 12. Association Rule Mining

- ① Association Rule mining is a procedure which is Meant to Find Frequent patterns, correlations, associations from datasets
- ② Association rule mining discovers strong association or correlation relationship among data
- ③ 1. Basket Data analysis  
2. Cross marketing  
3. Cohort design are the main applications of Association rule mining.

## 13. Apriori Algorithm

- ① Apriori algorithm was proposed by Agrawal and shrikant in 1994.
- ② Apriori is designed to operate on database Containing transactions
- ③ Apriori algorithm uses "bottom up" approach
- ④ Advantages:
  - 1. Efficient and effective with large item sets.
  - 2. Easily implementable
  - 3. provides real insight
- ⑤ Disadvantage:
  - 1. Assumes translational database in memory
  - 2. Requires many database scans
  - 3. Slow
  - 4. Higher runtime for execution.

## 14. What is Regression Analysis:

- ① Regression analysis is a set of statistical processes of estimating the relationship among variables.
- ② Regression analysis is a very useful and powerful tool technique.
- ③ Regression analysis has two main types

### i] Linear Regression analysis:

- ① linear regression analysis is one of the type of regression analysis.
- ② linear regression is used to predict the value of variable based on the value of another variable.
- ③ The variable you want to predict is called the dependent variable.
- ④ The variable you are using to predict the other variable's value is called independent variable.
- ⑤ General mathematical eq<sup>n</sup> linear regression

$$y = ax + b$$

### ii] Non-linear Regression analysis:

- ① it is form regression analysis in which observational data are modeled by a function which is a non-linear combination of model parameters and depends on one or more independent variables.

## 4. Data analytics with R/WEKA machine learning

### 1] What is WEKA?

- ① WEKA is an open source software provides tools for data preprocessing, implementation of machine learning algorithms.
- ② A WEKA is a collection of machine learning algorithms.
- ③ A WEKA is used for data mining task.
- ④ WEKA tool contains data preprocessing, classification, regression, clustering, visualization etc.
- ⑤ it is suited for developing new machine learning schemes.

### 2] What is Data manipulation?

- ① Data manipulation is the modification of information to make it easier to read and to be more organized.
- ② When data collection process is done by machine so many errors are in the data but data manipulation removes those errors & arranges data in structured order.
- ③ Example:  
Data log is sorted in alphabetical order to easy findout.

### 3] What is data Visualization?

- ① Data visualization is a technique Used for the graphical representation of data.
- ② Data visualization Uses charts, histograms, graphs.
- ③ Data Visualization make it easy to Understand data or info.
- ④ Data visualization make it easy to recognize patterns, trends and exceptions in our data.

### ⑤ Advantage of data visualization in R :-

- ① R offers broad collection of visualization libraries.
- ② R also offer data visualization in the form of 3d model.

### ⑥ Disadvantage of data visualization in R :-

- ① R is only preferred for data visualization
- ② R Data visualization Using R is slow for large amount of data;

### 4] what is Pipe Operator?

- ① Pipe operator is a special operational Function.
- ② Pipe operator wrap multiple functions together.
- ③ Pipe operator it is denoted as %>%.
- ④ Pipe operator make it easy when we need to perform various operations on a dataset to derive the results.

# Data Analysis

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

## 5) What is data analysis.

- ① Data analysis is defined as a process of cleaning transforming and modeling data to discover useful information for business decision making.
- ② The main purpose of data analysis is to extract useful information from data and making decision of business.
- ③ following are the some time types of data analysis:
  - ① Text analysis
  - ② statistical analysis
  - ③ Diagnostic analysis
  - ④ Predictive & Prescriptive
  - ⑤ Inferential
- ④ Process of data analysis / steps of data analysis process

Identify business problems ①

  └ Data acquisition ②

    └ Process/clean data ③

      └ Exploratory Analysis ④  
      └ Step

      └ Generate

Test

① Identify business model ⑤

statistical

problem

└ validate

Diagnostic

② data acquisition model ⑥

Predictive

③ Process / clean data

└ visualize

descriptive

④ Exploratory analysis model ⑦

⑤ model Generation &

Validation

# \* Remaining Important Questions \*

classmate

Date \_\_\_\_\_  
Page \_\_\_\_\_

## \* Remaining important Questions.

Q. What is statistical Modeling?

- ① statistical modeling is the process of applying a statistical analysis to a dataset
- ② statistical modeling is a mathematical representation of observed data
- ③ example:
  - 1) Public health data
  - 2) Social media data
- ④ steps of statistical modeling process
  - 1) model selection
  - 2) model fitting
  - 3) model validations

Q. What is Correlation & Write techniques for measuring correlation.

→ 1. Correlation analysis is method of statistical evolution used to study the strength of a relationship bet<sup>n</sup> two numerically measured, continuous variables.

2. Correlation analysis when we use we want to establish possible connection bet<sup>n</sup> variables.

3. For techniques for measuring Correlation see chapter no. 2 Notes back.

## Q: DIFF' Bet' Population and Sample

### Population

### Sample

- |   |  |
|---|--|
| 1. Population refers to the collection of all elements possessing common characteristics. | 1. Sample means a subgroup of members of population chosen for participation in the study. |
| 2. it includes each and every unit of group.  | 2. it includes only a handful of units of population.                                      |
| 3. Parameter is characteristic.   | 3. Statistic.  |
| 4. it focuses on identifying characteristics.   | 4. it focuses on making inferences about population.                                       |
| 5. it is defined as total of the items under Considerations.                              | 5. it is defined as a proportion of the population selected.                               |
| 6. example - Automobiles with four wheels, people who consume olive oil.                  | 6. example - 20 cars from each make.   |

## Q. What is Regression.

1. A regression is a statistical technique that relates a dependent Variable to one or more independent Variable.
2. Regression Significance is used to determine whether the relationship exists or not
3. example:- we can say age & height is described Using linear regression model.
4. Type of Regression model

Regression analysis

Simple

Multiple

### 5. I. Simple Regression:-

- It is used to estimate/ identify a relationship betn dependent Variable to a Single independent variable
- example - A relationship betn crop & rainfall.

### 6. II. Multiple Regression

- It is used to identify a relationship betn dependent Variable to a one or more independent Variable

- example - A relationship betn salary & employee & their experience education.

## 7. Regression analysis steps

1. Identify the statement of problem.
2. choice of relevant variables
3. collection of data on relevant variables
4. specification of model
5. choice of method for fitting the data
6. fitting of model
7. validation of model
8. Using the chosen model.

## 8. Diff' bet' Correlation and Regression



Points	Correlation	Regression
Meaning	Correlation is a relationship bet' two variables	Regression is a relationship bet' dependent to independent variable
Usage	To represent linear relationship bet' two variable	to fit a best line bet' dependent to independent Variable
Dependent & independent var	No difference	Both are different variables.
Objective	To find the numerical value bet' variables	To estimate the value of fixed Variables.

Q. describe Hadoop?

- ① Hadoop is one of the best data analysis tool used in Big data.
- ② Apache Hadoop is an open source framework.
- ③ Hadoop is used to store process and analyze data which are very huge in volume.
- ④ Hadoop is used to efficiently store and process large datasets ranging in size from gigabytes to petabytes of data.
- ⑤ Hadoop allows clustering multiple computers to analyze massive datasets in parallel more quickly.
- ⑥ Hadoop is written in Java.

Q. Write down tools used for big data

- ① NoSQL - Database, MongoDB, Cassandra.
- ② MapReduce - Hadoop, Hive, Pig, Kafka.
- ③ Storage - Hadoop distributed file system, S3
- ④ Server - El2, Google app engine
- ⑤ Processing - R, Yahoo! pipes etc.

Above all tools used for big data.

Q. Why hadoop is used for big data analysis?

- ① Hadoop makes it easier to use all the storage and processing capacity in cluster servers, and to execute distributed process against huge amount of data.
- ② Hadoop provides building blocks on which other services and applications can be built.
- ③ Hence Hadoop is used in big data analysis.

\* What are the features of hadoop.

- ① Hadoop is Open Source.
- ② Hadoop is easy to use.
- ③ Hadoop is very cost-effective.
- ④ Hadoop is faster in Data processing.
- ⑤ Hadoop provides feasibility.

Q. Explain need for modern Corporate Sector to go for Big data strategy.

- ① Companies use big data strategy in their systems to improve operations, provide better customer service, and take other more actions that ultimately can increase revenue and profits.
- ② Big data processes large amount of data and help to grow organisation and organisation wealth. hence mordern Corporate Sector

need to go for Big data strategy.

Q. What is training, testing and cross-validation of machine learning models.

→ ① Training a model means learning good values for all the weight and bias from labeled.

② The testing statistical based test for machine learning model.

③ The team would train, test and validate their model on real world data.

④ Cross-validation is a technique for validating model efficiency by training it on the subset of input data.

Q. Explain tools used in Big data.

I] Apache Hadoop -

it is open source framework from apache and runs commodity hardware. It is used to store process and analyze Big data.

II] Apache Spark -

It is considered as a successor of Hadoop as it overcomes the drawbacks of its. Spark unlike Hadoop supports both real time as well as batch processing.

### III] Apache storm -

1. Apache storm is big data analysis tool.
2. Apache storm is a free and open source distributed real-time computation system.
3. Apache storm is efficiently process unbounded storm of data.
4. Storm is designed to process vast amount of data in a fault-tolerant.

### Q. What is Overfitting and Underfitting.

→ ① Overfitting means that your model makes not accurate predictions. In case of overfitting train error is very small and val/test error is large.  
Overfitting can happen due to low bias to high variance.

② Underfitting means that your model makes accurate but initially incorrect prediction. In case of underfitting, train error is large and val/test error is large too.  
Underfitting can happen due to high bias to high variance.

### Q. Explain Deep learning.

→ ① Deep learning is a machine learning technique that teaches computers to do what comes naturally to humans.

② Deep learning is a key technology behind driverless cars.

③ Deep learning is subset of machine learning.

④ Applications.

1. Self driving car
2. Natural lang processing
3. Virtual assistant
4. Visual recognition
5. Fraud detection.
6. Healthcare.

Q. Explain Role of machine learning algorithm.

- 1. Machine learning algorithms uses historical data as input to predict new output values.
- 2. A machine learning is a system by which AI system conduct its task.
- 3. The two main process of machine learning algorithm are classified classification and regression.

Q. Short note on Random forest.

- ① it is supervised machine learning algorithm
- ② it used widely in classification and regression problems,
- ③ it builds decision tree on different samples and make decision tree.
- ④ it provides better results for classification problems.

## Q. Understood machine learning

1. The Broad - ML Predicting things based on what they done in the past
2. the Practical - ML finds the relationship bet<sup>n</sup> in your data.
3. The Technical - ML USE statistical methods to predict value.
4. The mathematical - ML attempts to predict the value v given an input.