

1. List and define 5V's of Big Data. List different challenges and issues in Big Data.

→ The 5V's of Big Data are as follows -

(1) volume

→ The name 'Big Data' itself is related to a size which is enormous. Volume is a huge amount of data. To determine the volume of data, size of data plays a crucial role. If the volume of data is very large then it is actually considered as 'Big Data'.

(2) velocity

→ It refers to the high speed of accumulation of data. In Big Data velocity data flows in from sources like machines, networks, social media, mobile phones, etc.

(3) variety

→ It refers to nature of data that is structured, semi-structured and unstructured also.

(4) veracity

→ It refers to inconsistencies and uncertainty in data which is available can sometimes get messy and quality and accuracy are difficult to control.

(5) value

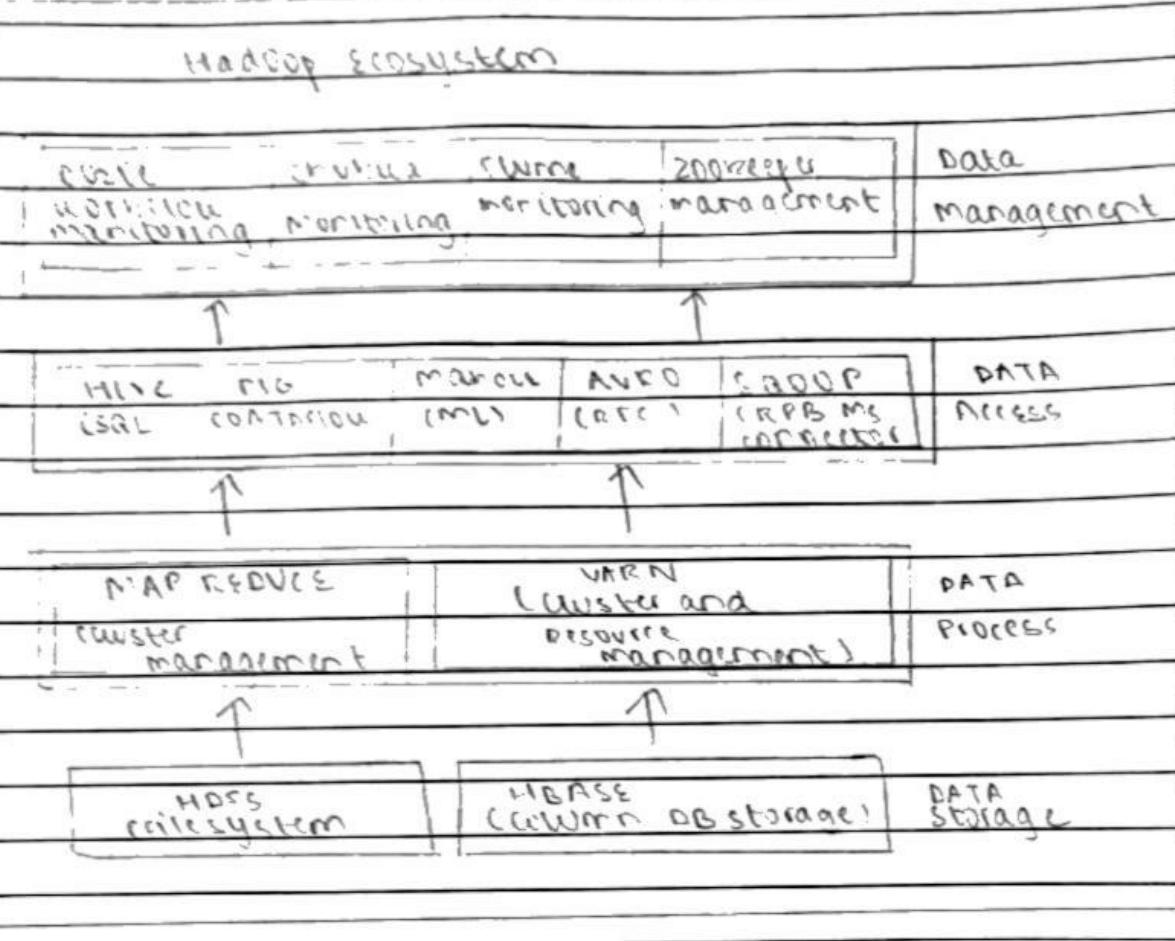
→ The bulk of data having no value is no good to the company, unless you turn it into something useful. Data in itself is of no use or importance but it needs to be converted into something valuable to extract information.

- The Galaxy Education System
- some of the Big Data challenges are -
    - (1) sharing and accessing data
    - (2) Privacy and security
    - (3) Analytical challenges
    - (4) Technical challenges
  - 2. list and explain different components of Hadoop
    - (1) HDFS
      - Hadoop distributed file system is the backbone of Hadoop which runs on java language and stores data in Hadoop applications. They act as a command interface to interact with Hadoop. The two components of HDFS are - Data Node and Name Node.
      - Name Node manages file system and operates all data nodes and maintains records of meta data updating. In case of deletion of data, they automatically record it in edit log.
      - Data Node requires vast storage space due to reading and writing operations. They work according to the instructions of the Name node.
    - (2) HBASE
      - It is an open source framework storing all types of data and doesn't support the SQL database. They run on top of HDFS and written in java language.

- The two major components of HBASE are :- HBASE master and regional server. HBASE master is responsible for load balancing in a Hadoop cluster and controls the failover. The regional servers role would be a worker node and responsible for reading, writing data in cache.
- (3) YARN  
→ It is an important component in the ecosystem and called an operating system in Hadoop which provides resource management and job scheduling task. The components are resource and Node manager, application manager and container. They also act as guards.
- (4) Sqoop  
→ It is a tool that helps in data transfer between HDFS and MySQL and gives hand on to import and export data, they have a connector for fetching and connecting data
- (5) Apache spark  
→ It is an open source cluster computing framework for data analytic and an essential data processing engine
- (6) Apache Flume  
→ It is distributed service collecting a large amount of data from the source and moves back to its origin and transferred to HDFS.

- The Galaxy Education System
- (7) Hadoop Map Reduce
    - It is responsible for data processing and acts as a core component of Hadoop
  - (8) Apache Pig
    - Data manipulation of Hadoop is performed by Apache Pig. It helps in the reuse of code and easy to read and write code.
  - (9) Hive
    - It is an open source platform for performing data warehousing concepts; it manages to query large data sets stored in HDFS
  - (10) Apache Drill
    - Open source SQL engine which process non-relational databases and file systems
  - (11) Apache Zookeeper
    - It is an API that helps in distributed coordination.
  - (12) Oozie
    - It is a java web application that maintains many workflows in a Hadoop cluster

## 3. Draw and explain the Hadoop Ecosystem



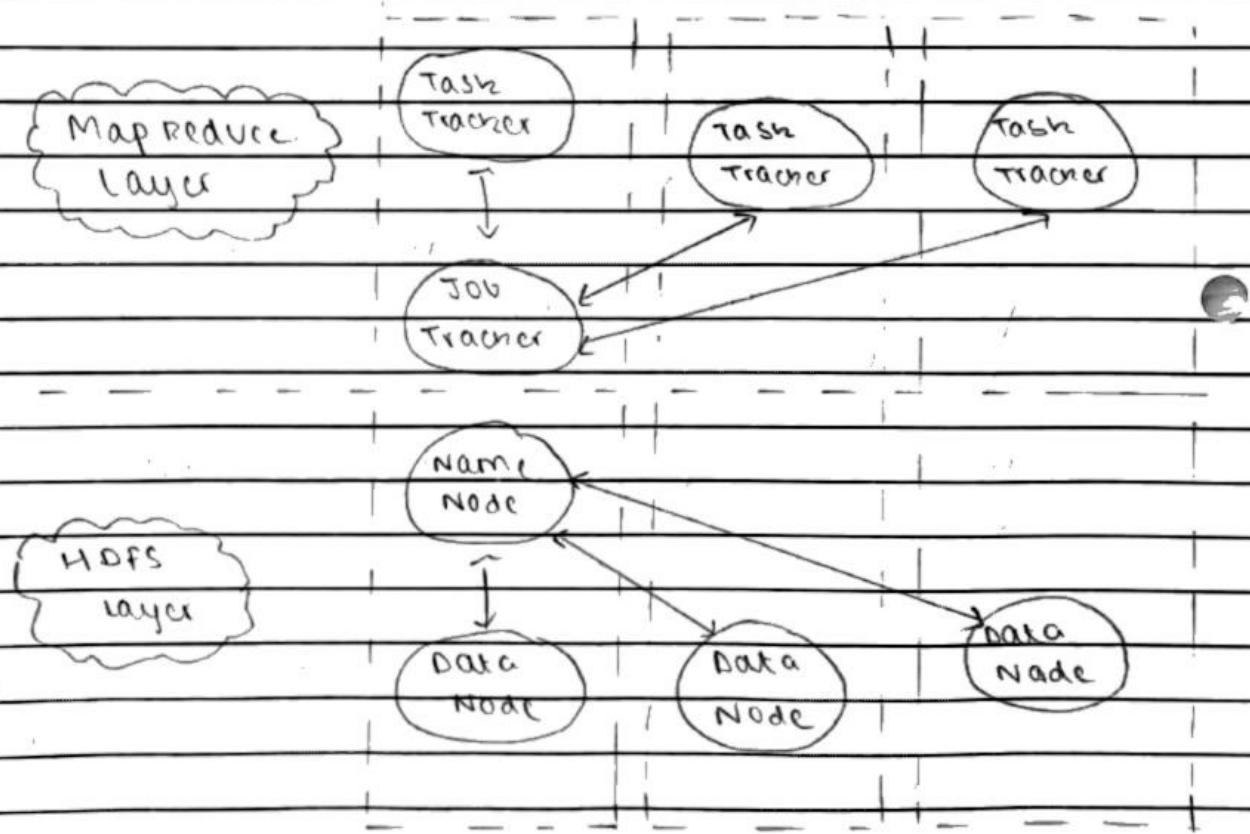
- Hadoop ecosystem is a platform or a suite which provides various services to solve the big data problems
- It includes Apache projects and various commercial tools and solutions
- There are 4 major elements of Hadoop - HDFS, Map Reduce, YARN and Hadoop common. Most of the tools work collectively to provide services such as absorption, analysis, storage and maintenance of data.

→ Components of Hadoop -

1. HDFS → Hadoop distributed file system
2. YARN → Yet Another resource negotiator
3. Map Reduce → Programming based data process
4. Spark → In memory data processing
5. PIG/HIVE → Query based processing on data service
6. HBASE → No SQL Data Base
7. Mahout → ML algorithm libraries
8. SOLAR, WIENNE → searching and indexing
9. ZOOKEEPER → Managing cluster
10. OOZIE → Job scheduling

#### The Galaxy Education System

4. Draw and explain architecture of Hadoop. List limitations of Hadoop



- The Galaxy Education System
- Name node  
It represents every files and directory which is used in the namespace
  - Data Node  
It helps you to manage the state of an HDFS node and allows you to interact with the blocks.
  - Master Node  
It allows you to conduct parallel processing of data using Hadoop map reduce
  - Slave Node  
Allows you to store data for complex calculations.
  - Limitations of Hadoop →
    - (1) Problem with small files
    - (2) Vulnerability
    - (3) Low performance in small data surroundings
    - (4) Lack of security
    - (5) High up processing
    - (6) Batch processing is only supported
  - 5. What is Map Reduce Engine? Explain its principle execution
    - It is a processing technique and a program model for distribution
    - The Map Reduce algorithm contains two main tasks, namely Map and Reduce

- Map takes a set of data and converts it into another set of data, where individual elements are broken down
- Reduce task, which takes the output from a map as an input and combines those data tuples into a smaller set of tuples. As the sequence of the name MapReduce implies, the reduce task is always performed after map job.

The Galaxy Education System

#### 6 Advantages Of NoSQL

- Big data capability
- No single point failure
- Easy replication
- Fast performance and horizontal scalability
- Can handle unstructured data as well
- Multiple vendor choices are available
- Easy and human readable query

7. List the different NOSQL data architectural patterns? Explain with diagram and give one example of each

→ (1) Key-value store database

→ This model is one of the most basic models of NOSQL database. As the name suggests, the data is stored in form of key-value pairs. The value is typically linked or related to the key. The key value pair storage database generally stores data as a hash table where each key is unique.

Eg - DynamoDB, Berkeley DB

key	value
$k_1$	$v_1$
$k_2$	$v_2$
:	:
$k_m$	$v_n$

→ (2) Column store database

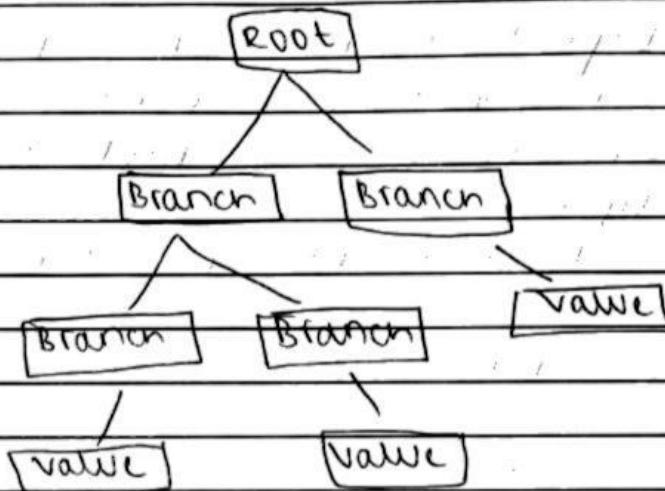
→ The data is stored in individual cells which are further grouped into columns. Column oriented databases work only on columns.

They store large amount of data

Eg - HBASE, Bigtable by google

row id	column family	column name	timestamp	value

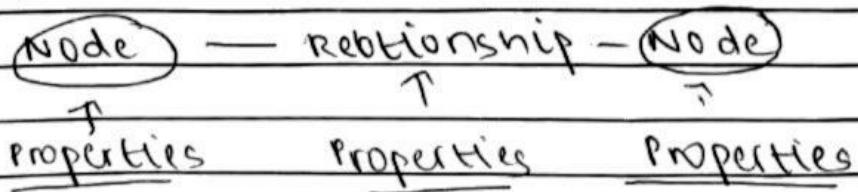
- (3) Document database
- It fetches and accumulates data in form of key value points but values are called documents. It has tree like structure and can be complex.
- Eg - MongoDB, Couch DB



- + (4) Graph database

→ It deals with storage and management of data in graphs. Objects are called as nodes and are joined together by relationship called edges

Eg - Neo4J, Flock DB



## 8. Hamming distance

- calculates distance between two binary vectors
- mostly used when dealing with one-hot-encode categorical columns

$$\text{Hamming distance} = \sum_{i=1}^N \text{abs}(x_i - y_i)$$

## Euclidean distance

- calculates distance between two real valued vectors
- mostly used for columns with numerical values like int and float

$$\begin{aligned}\text{Euclidean distance} &= \sqrt{\sum_{i=1}^N (x_i - y_i)^2} \\ &= \sqrt{\sum_{i=1}^N (x_i - y_i)^2}\end{aligned}$$

## Manhattan distance

- more useful to vectors that describe objects on uniform grid like chessboard

$$\text{Manhattan distance} = |x_1 - x_2| + |y_1 - y_2| \text{ for point } x_1, y_1, x_2, y_2$$

## Minkowski distance

- it is generalization of Euclidean and Manhattan distance

$$\text{Minkowski distance} = \sqrt[p]{\sum (x_i - y_i)^p}$$

## 9. Stream data

- stream data is data that continuously generated by different sources.

Characteristics - time sensitive

- continuous
- Heterogeneous
- Imperfect
- volatile

challenges - choosing data formats

- algorithm testing, life cycle management
- scaling and performance

- sampling methods
  - Bernoulli sampling
  - reservoir sampling
  - stratified sampling
  - based reservoir sampling

	Traditional Data Approach	Big Data Approach
→ Data integration is easy		Data integration is difficult
→ centralized and managed in centralized form		Managed in distributed form
→ volume ranges from 6B to 1TB		volume ranges from Petabytes to zettabytes

### The Galaxy Education System

*	Document store	Graph store
→	stores the data in form of document as value	stores the data in node connected by edges
→	e.g. mongo DB	e.g. Neo4J
*	Job tracker	Task tracker
→	runs on separate node	runs on all nodes
→	replaced by resource manager	replaced by node manager in MR2
→	receives requests for execution from client	assigned mapper and reducer tasks to execute

- |                              |   |
|------------------------------|---|
| * Name Node                  | Data Node   |
| → centerpiece of HDFS        | responsible for storing actual data                 |
| → known as master            | known as slave                                      |
| → single point of failure    | If its down, it doesn't affect availability of data |
|                              |   |
| * DBMS                       | NoSQL   |
| → Data stored as file system | Non-relational database system                      |
| → supports single user only  | multiple user                                       |
| → mix of open source         | Open source only                                    |