

## Unit 3

### (H) Diff b/w Structured, Unstructured and semi-structured data

Properties	Structured data	Semi-structured data	Unstructured data
Technology	It is based on Relational database table	It is based on xml / Resource Description framework (RDF)	It is based on character & binary data.
Version management	versioning over tuples, rows, tables	versioning over tuples or graph is possible	versioned as a whole
Flexibility	It is schema dependent & less flexible	It is more flexible than structured data but less flexible than unstructured data.	It is more flexible & there is absence of schema.
Scalability	It is very difficult to scale in schema	Its scaling is simpler than structured data	It is more scalable.
Robustness	very robust	new technology not very strong	—
Query performance	structured query allow complex joining	Queries over anonymous nodes are possible	only textual queries are possible

### Unstructured

#### # Analytical of data Analytics.

Unstructured data is data that doesn't have a fixed form or structure. Images, videos, audio files, text files, social media data, geospatial data, data from IoT devices, & surveillance data are example of unstructured data. About 80% - 90% of data is unstructured.

Unstructured data analysis is complex & requires specialized techniques, unlike structured data, which is straightforward to store & analyze.

Here is a quick glance at all of the unstructured & analysis techniques I like discussed.

1. Keep the business objective(s) in mind.
2. Define metadata for faster data access.
3. Choose the right analytic techniques.
  - a. Exploratory data analysis techniques.
  - b. Qualitative data analysis techniques
  - c. Artificial Intelligence (AI) & Machine Learning (ML) techniques.
4. Identify the right data sources
5. Evaluate the technologies you'd want to use
6. Get real time data access
7. Store & integrate data using data lake
8. Mangle the data to get the desired features.

### Data Analytics Use Cases.

To achieve the strategic objectives of the company a better data strategy will help how to use that data in an efficient manner. These data uses which are identified in this process are called use cases.

These use cases are the key data elements for the data projects ahead.

While developing a data strategy for a business, a few optimal numbers of use cases determined. Over and over use cases result in an unrealistic & cluttered data strategy which may be risky. These data use cases are different for every business which is actually driven by the business strategy.

The main elements involved in designing & developing better data use cases are:

- ① Data evaluation.
- ② Identification of suitable analytical techniques & tools
- ③ Hierarchy of outcome metrics



Here are top five use cases for Data Analytics

### ① Security Intelligence :-

Data analytics are also deployed for improving security operations against Hackers & Cyber criminals. An IT department handles a large amount of data. Security is an important concern; many companies now use these analytics to help in obtaining better solutions to detect & prevent such attacks. User and Entity Behaviour Analytics (UEBA), security information & event management (SIEM) tools, & machine learning can be used to detect abnormalities & unusual user activities.

### ② Customer Relationship Analytics :-

One of the difficult tasks in marketing is to identify customers who are going to spend money consistently for a long time. This business insight will help the companies to gain such customers which will be a lifetime value for company. The data analytics examples in business include Telecommunications, Banking, utilities & retail. Customer segmentation also helps in establishing potential marketing strategies.

### ③ Recommendation Engines :-

You may have noticed "recommendations for you" on Youtube, Spotify or other media services. These personalized recommendations are one way of help in improving the overall user experience. This can be a winning factor as there is a lot of competition in the entertainment & media sector, you can also see such recommendations while shopping online.

### ④ E-commerce :-

A system wide ~~information~~ infrastructure with data analytics is a superior way to enhance the efficiency and performance of the business. The system metrics are used to track the performance of IT modules & the user logs. These will identify user behaviour in e-commerce sites. By using this data, retailers can gain insights that help in developing agile techniques, better business performance & profits.

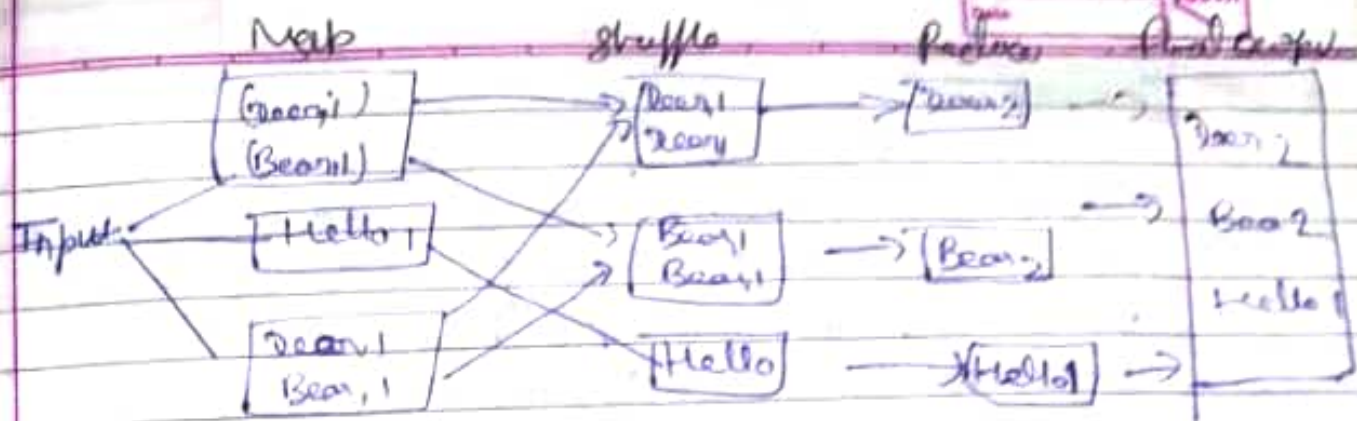
### ⑤ Forecast of future :-

There are endless ways of applying data analytics to IoT solutions for eg:- almost every industry uses sensor data to have actionable insights as a measure of predictive maintenance. A customer's movement may be tracked as a part of the security option of a device. Another example can be logistic tracking, where the vehicles are tracked from time to time & re-directed in the case of bad weather or unexpected circumstances that affect delivery time.









## At benefits of map reduce

1) Parallel processing ⇒ In map reduce job is divided among multiple nodes and each node works with a part of job simultaneously, so map reduce is based on divide and conquer paradigm which help us to process the data using different machines as the data is processed by multiple machines instead of a single machine in parallel. the time taken to process the data gets reduced by tremendous amount.

2) Data locality ⇒ In the traditional system data is sent to the processing unit & then it is processed but as the data grows & becomes very huge moving this amount of data to the processing unit posed the following issues.

a) moving huge data to processing unit is costly & degrades the network performance.

b) The master node get over burdened & may fail map reduce allows us to overcome the above issues by bringing the processing unit to the data & is distributed among multiple nodes where each node processes the part of the data residing on it this allows us to have the following advantage.

① To move processing unit to the data.

② The processing time is reduced as all the nodes are working with their part of the data in parallel.

③ Every node gets the part of data to process & therefore there is no chance of a node getting over burdened.



## # Hadoop

Hadoop is an Open Source Java based software platform that manages data processing & storage for big data applications. Hadoop works by distributing large data sets & analysis jobs across nodes in computing clusters, breaking them down into small workloads that can be run in parallel. Hadoop process structure & unstructured data & scale up reliably from a single server to thousands of machines.

It is highly scalable.

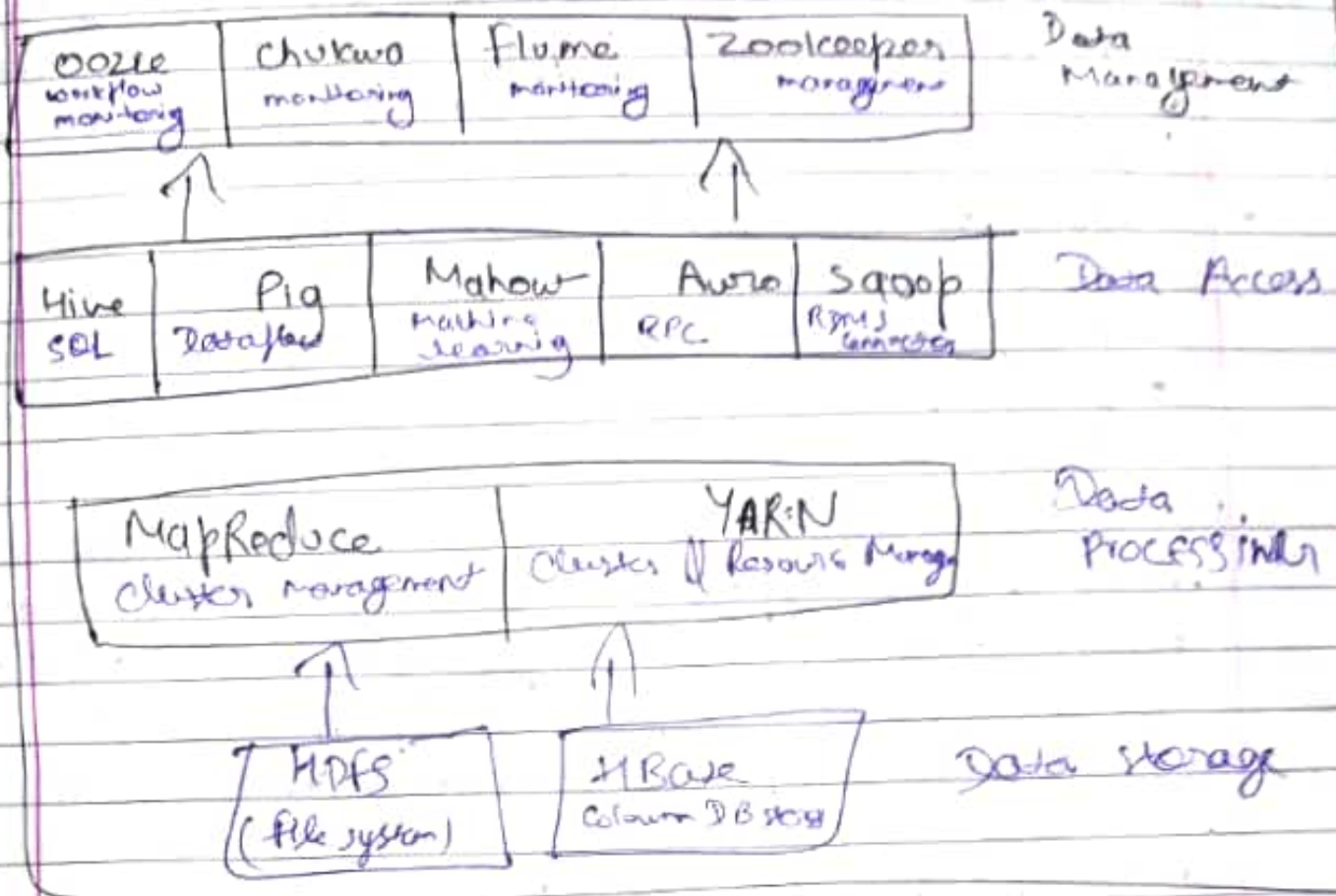
It consists of 3 components

- ① HDFS → Reliable storage system with half of the world data stored in it
- ② MapReduce → layer consists of distributed processes
- ③ Yarn → layer consists of resource manager.

## # Hadoop Ecosystem

The Hadoop ecosystem is a framework that helps in solving big data problems. The core component of the Hadoop ecosystem is a Hadoop distributed file system. HDFS is a distributed file system that has the capability to store a large stack of data sets. With the help of shell commands, HADOOP interacts with HDFS. Hadoop breaks up unstructured data & distributes it to different sections for data analysis. The ecosystem provides many components & technologies that have the capability to solve complex business tasks.

# Hadoop Ecosystem



## # HDFS

- HDFS is the primary or major component of Hadoop ecosystem and it's responsible for storing large data sets of structured or unstructured data across various nodes and thereby maintaining the metadata in the form of log files.
- HDFS consist of 2 Core Component i.e.
  - ① Name node
  - ② Data node
- Name node is the prime node which contains metadata (data about data) requiring comparatively fewer resources than the data nodes are commodity hardware in the distributed environment, undoubtedly, making Hadoop cost effective.
- HDFS maintains all the coordination between the cluster & hardware thus working as the heart of the system.



## • # Yarn

- yet another Resource Negotiator with the same implies YARN is the one who helps to manage the resources across the clusters. In short, it performs Scheduling of resource allocation for the Hadoop system.
- Consists of 3 major components i.e.
  - ① Resource manager
  - ② Node manager
  - ③ Application Manager.
- Resource manager has the privilege of allocating resources for the applications. In a system where Node managers work on the allocating of resources such as CPU, memory bandwidth per machine & data on acknowledgment resource manager. Application manager works as an interface between the resource manager & node manager & performs negotiations as per the requirement of the two.

## # MapReduce :-

- By making the use of distributed & parallel algorithms MapReduce makes it possible to carry over the processing logic & helps to write applications which transform big data sets into a manageable one.
- MapReduce makes the use of 2 functions i.e.
  - ① MapReduce makes the use of 2 functions i.e. map() & Reduce() whose task is:
    - ① Map() performs sorting & filtering of data & thereby organizing them in the form of groups. Map generates a key value pair based results which is later on processed by the Reduce() method.
    - ② Reduce(), as the name suggests does the summarization by aggregating the mapped data.



In simple Reduce() - takes the output generated by map() as input // Combines these tuples into smaller sets of tuples.

## # PIG

- Pig was basically developed by Yahoo which works on a pig latin language, which is Query based language similar to SQL.
- It is a platform for structuring the data flow, processing and analyzing huge data sets.
- Pig does the work of executing commands // in the background all the activities of MapReduce are taken care of. After the processing pig stores the result in HDFS.
- Pig latin language is specially designed for this framework which runs on Pig runtime. Just the way Java runs on the JVM.
- Pig helps to achieve ease of programming // optimization // hence is a major segment of the Hadoop Ecosystem.

## # HIVE

- With the help of SQL methodology // interface, HIVE performs reading // writing of large data sets. However, its query language is called as HQL (Hive Query Language).
- It is highly scalable as it allows real-time processing // batch processing both. Also, all the SQL datatypes are supported by Hives thus making the query processing easier.
- Similar to the Query processing frameworks HIVE too comes with two components: JDBC Drivers // HIVE Command Line.
- JDBC along with ODBC drivers work on establishing the data storage permissions // connection whereas HIVE Command line helps in the processing of queries.



## # Mahout :-

- Mahout allows machine learnability to a system application. Machine learning as the name suggests helps the system to develop itself based on some patterns, user/environmental interaction or on the basis of algorithms.
- It provides various libraries or functionalities such as Collaborative filtering, Clustering & Classification which are nothing but concepts of machine learning. It also involves algorithms as per our need with the help of its own libraries.

## # Apache Spark :-

- It is a platform that handles all the process consumption like batch processing, interactive or iterative & real-time processing, Graph Conversions, & Visualizations etc.
- It consumes less memory resource hence, it is being faster than the prior in terms of optimization.
- Spark is best suited for real-time data whereas Hadoop is best suited for structured data or batch processing hence both are used in most of the companies interchangeably.

## # Apache HBase :-

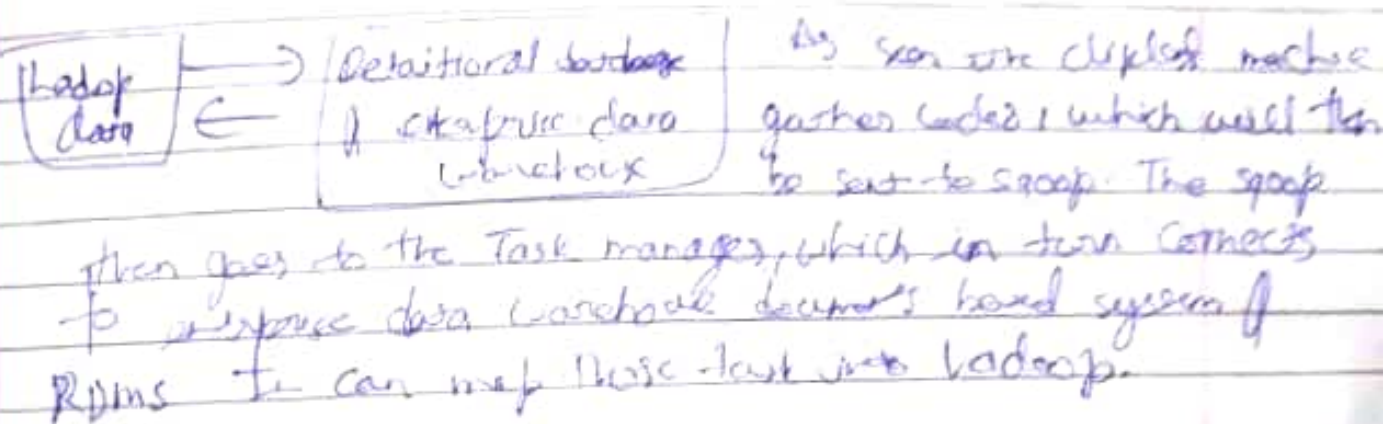
- It is a NoSQL database which supports all kinds of data & thus capable of handling anything of Hadoop Database. It provides capabilities of Google BigTable thus able to work on Big data sets effectively.
- At times when we need to search or retrieve the occurrences of something small in a huge database, the request must be processed within a short



quick span of time. At such times, HBase comes handy as it gives us a relevant way of storing limited data.

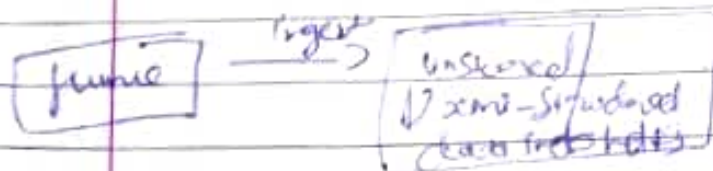
## # Sqoop

Sqoop is used to transfer data between Hadoop and external databases. Such a relational database and enterprise data warehouses. It imports data from external data stores into HDFS, Hive & HBase.



## # Flume

Flume is another data collection & ingestion tool. A distributed service for collecting aggregating & moving large amounts of log data. It ingests online streaming data from social media, log files, web server into HDFS.



data is taken from various sources, depending on your organization's needs. It then goes through the source, channel & sink. The sink feature ensure that everything is in sync with the requirements. Finally the data is dumped into HDFS.

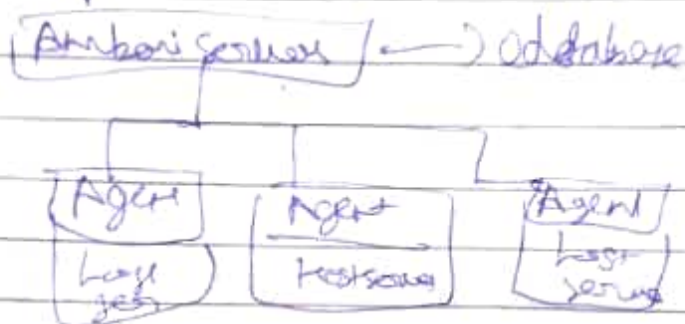


## Ambari

It is an open source tool responsible for keeping track of running applications & their statuses.

Ambari manages, monitors, & provisions Hadoop clusters, also it also provides a central management service to start stop & configure Hadoop services. ~~start stop & configure~~

### Ambari web



## Storm

The Storm is an engine that processes real-time streaming data at a very high speed. It is written in Clojure. A Storm can handle over 1 million jobs on a node in a fraction of a second. It is integrated with Hadoop to harness higher throughput.

Now that we have looked at the various data processing tools & streaming services let us take a look at the security framework in the Hadoop ecosystem.

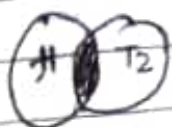


## SQL Joins

A Join clause is used to combine rows from 2 or more tables based on selected columns between them.

### Types of SQL join

(Inner) Join  $\Rightarrow$  Returns records that have matching values in both tables.



Left (Outer) Join  $\Rightarrow$  Returns all records from the left table & the matched records from the right table.



Right (Outer) Join  $\Rightarrow$  Return all records from the right table & the matches from the left table.



Full (Outer) join  $\Rightarrow$  Returns all records when there is a match in either left or right table.



### SQL Set Operations

The SQL set Operations is used to combine the two or more SQL SELECT statements.



## Types of Set Operations

### ① Union

- The SQL union operation is used to combine the results of two or more SQL SELECT queries.
- In the union operation, all the no. of datatypes & columns must be same in both tables on which UNION operation is being applied.
- The Union operation eliminates the duplicate rows from its results.

eg SELECT Column-name FROM table1  
UNION  
SELECT Column-name FROM table2;

### ② Union All

Union all operation is equal to the Union operation. It returns the set without removing duplicates & sorting the data.

eg SELECT Column-name FROM table1  
UNION ALL  
SELECT Column-name FROM table2;

### ③ Intersect

- It is used to combine two SELECT statements. The intersect operation returns the common rows from both the SELECT statements.
- In the intersect operation, the no. of datatypes & columns must be the same.
- It has no duplicates & it arrange the data in ascending order by default.

eg SELECT Column-name FROM table1  
INTERSECT  
SELECT Column-name FROM table2;



#### ④ Minus

- It combines the result of 2 Select Statements.  
minus operator is used to display the rows which are present in the first query but absent in the 2nd query.
- It has no duplicates & data arranged in ascending order by default.

eg `SELECT Column-name FROM table1,`  
`minus`

`SELECT Column-name FROM table2;`