

Assignment - I

Q. Explain classification of big data

→ Big data can be broadly classified into structured, semi-structured and unstructured data.

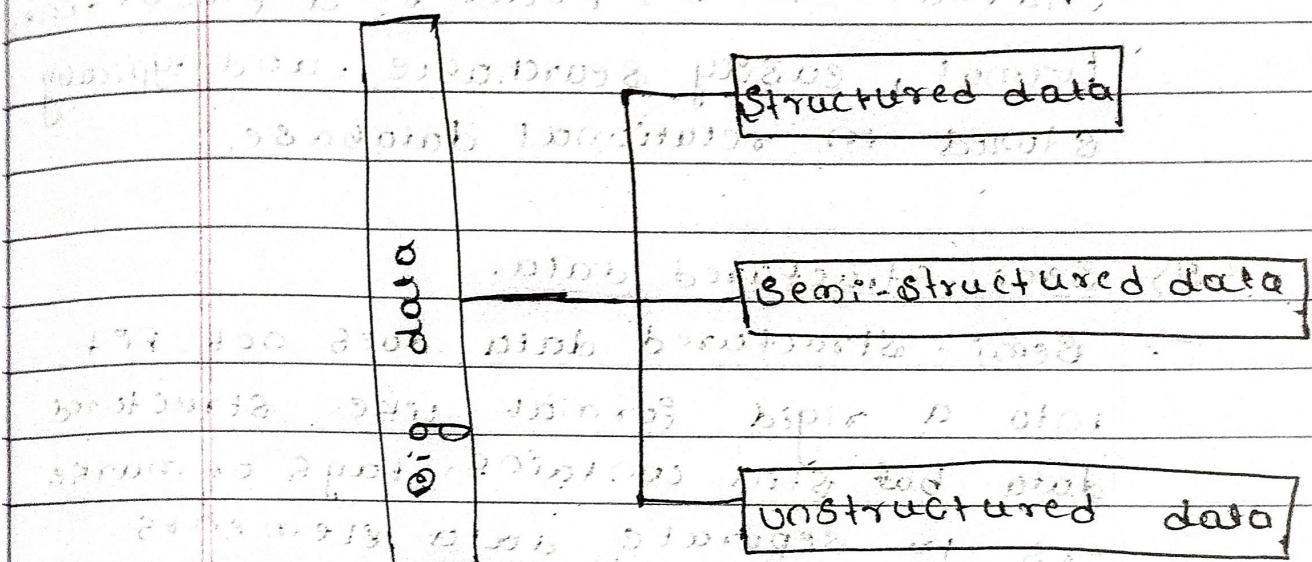


Fig. 1. Classification of Big data

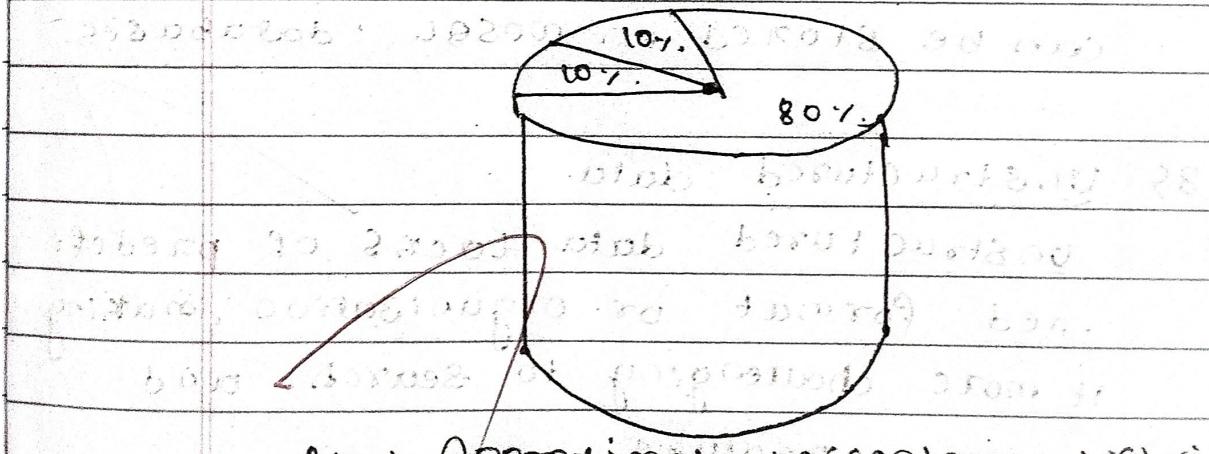


Fig: Approximate percentage distribution

Glaciers & snow fields, State of big Patagonia

structured data

→ Structured data is highly organized

Date _____ / _____
Page _____

and in a way that makes it easily searchable databases such as

ex:- financial transactions, customer information, product details.

Characteristics : follows a predefined format, easily searchable, and typically stored in relational database.

Q3) Semi-structured data.

→ Semi-structured data does not fit into a rigid format like structured data, but still contains tags or markers to separate data elements

Ex:- JSON, data, XML files, CSV files

Characteristics :- more flexible than structured data, contains metadata and can be stored in NoSQL databases.

Q4) Unstructured data.

Unstructured data lacks of predefined format or organization, making it more challenging to search and analyze.

Ex:- Social media posts, images, video's, audio files, emails.

Characteristics :- Does not fit into relational databases, requires special

needed analytics tools and often contains rich, contextual information.

Q5) List and explain characteristics of big data illustrate consider by an example of e-commerce, how big data is used.

The characteristic of Big data

composition

Condition

Data

Context

I) Composition : The composition of data deals with the structure of data, that is, the source of data, the granularity,

The types, and the nature of data as it relates to whether it's static or real time, or streaming.

II) Condition : - The condition of data deals with the state of data, that is "

can one use this data set is for analysis?" Does it require

Clearing for further enhancement and enrichment?"

3) Context :- The context of data deals with "where has this data generated?" "how sensitive is this data?" "what are the events associated with this data?" and so on.

Let's consider on e-commerce platform like Amazon or Flipkart.

How Big data is used

Is personalization : Big data helps in personalizing the shopping experience for customers by analyzing their browsing history, purchase patterns and search queries.

2) Inventory Management : E-commerce platforms use big data analytics to predict demand, manage inventory levels, and optimize supply chain operations.

3) Customer Segmentation : Big data enables segmentation of customers

based on behavior, preferences and demographics to target marketing campaigns effectively.

4) Fraud Detection :- Analyzing Big data helps in detecting and preventing fraudulent transactions by identifying causal patterns.

Characteristics of Big data :-

i) Volume : The amount of data.

ii) Velocity : The speed at which data is generated and processed.

iii) Variety : The type of data (structured, semi-structured, unstructured).

iv) Veracity : The quality and accuracy of data.

Big data often includes incomplete or inconsistent data.

Ensuring data quality and filtering out noise is essential for reliable analytics.

Value :-

The usefulness of data for decision making.

Q 83

Discuss about importance of big data Analytics.

→ Big data analytics play a crucial role in today's data driven world. Here are some important of big data analytics.

1) Reactive - Business Intelligence :- what does Business intelligence (BI) help us with? It allows the business to make faster and better decisions by providing the right information to right person at the right time in the right format.

It is about analysis of the past or historical data and then displaying the findings of the analysis or reports in the form of enterprise dashboards, alerts notifications, etc. It has support for both pre-defined specified reports as well as ad hoc querying.

2)

Relative - Big data Analytics :- Here the analysis is done on huge data sets but the approach is still reactive as it is still based on

Static data such as

No 83 proactive analytics : The analytics is done on huge datasets. It supports futuristic decision making by the use of data mining, predictive modeling, text mining, and statistical analysis. This analysis is not on big data as it still uses the traditional database management practices on big data and therefore has several limitations on the storage capacity and the processing capability.

4) proactive - big data analytics :-

This is achieved through terabytes, petabytes, exabytes of information to filter out the relevant data to analyze. This also includes high performance analytics to gain rapid insights from big data and the ability to solve complex problems using more data.

Q 84 Explain about big data and data warehouse co-existence.

Big data and data warehouses serve different purposes but can complement each other within an organization's data architecture.

Q1) Data warehouses :-

- * store structured data
- * optimized for SQL queries and business intelligence
- * used for reporting, dashboards and historical analysis

Q2) Big data platforms :-

- * handle large volume of structured, semistructured and unstructured data
- * support advanced analytics, machine learning and real-time processing
- * often use distributed computing frameworks like MapReduce or Apache Spark

How they co-existence is :-

Q3) Different purposes :-

Data warehouses are designed for structured data and support designed for

structured data and support business intelligence, reporting and analytics. Big data solutions handle large volumes of unstructured and semi-structured data for advanced analytics, machine learning and real-time processing.

Q4) Data integration :-

Organizations can integrate data from big data platforms into their data warehouse for comprehensive analysis and reporting.

Q5) Data lake architecture :-

A data lake can store big data, which can then be processed and loaded into a data warehouse for structured analytics.

Q6) Explain about typical data warehouse and hadoop environment.

A typical data warehouse environment is operational or transactional or day-to-day business data gathered from enterprise Resource planning systems, customer relationship management (CRM)

Date / /
Page

legancy systems, and several third party applications.

* The data from these sources may differ in format (data could have been housed in any RDBMS such as Oracle, MS SQL Server, DB2, MySQL and Teradata, and so on or in spreadsheet -- (.xls, .xlsx, etc) or .csv (.txt).)

* Data may come from data sources located in the same geography or different geographies.

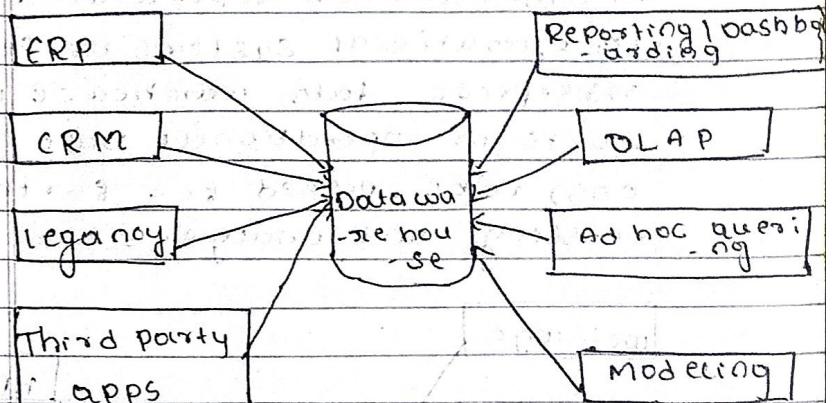
* This data is then integrated, cleaned up, transformed, and standardized through the process of Extraction, Transformation and Loading (ETL).

The transformed data is then loaded into the enterprise data warehouse (available at the enterprise level) or data marts (available at the business unit / functional level).

* A host of market leading business intelligence and analytics tools are then used to enable decision making from the use of ad-hoc queries, SQL enterprise dashboards

Date / /
Page

Data mining etc. is part of



e.g.: A typical data warehouse environment

Hadoop environment :-

* It is very different from the data warehouse environment and what exactly is this difference.

* The data source here includes disparate data from web logs to images, audios and videos to social media data to the various docs, pdfs, etc. Here the data is focus is not just the data within the company's firewall but also data residing outside the company's fire wall.

- The data is placed in Hadoop Distributed File System (HDFS). If need be this can be repopulated back to operational systems or fed to enterprise data warehouse or data marts or operational data store (ODS) to be picked for further processing and analysis.

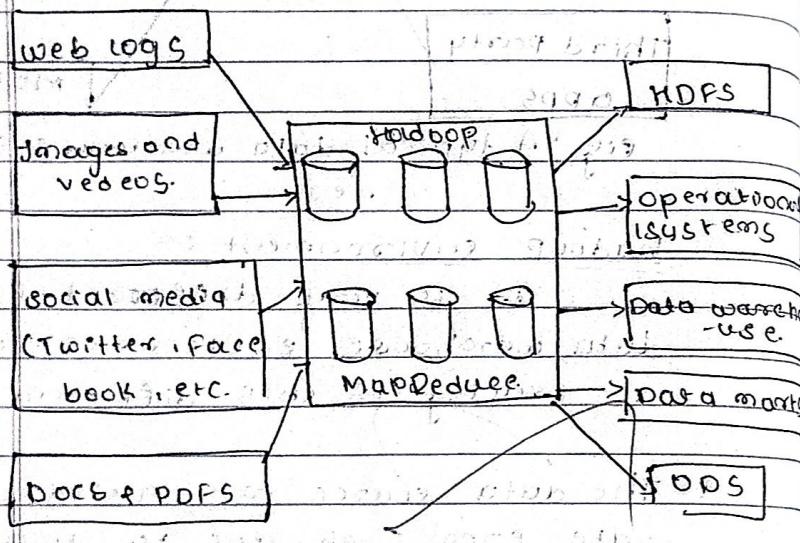


Fig : A typical Hadoop environment

Q6) Explain the technologies used in big data environment brief about few top analytical tools

⇒ Technologies used in Big data environment Big data environments utilize a variety of technologies to store, process and analyze large volumes of data. Some key technologies include:

Storage Technologies:-

Hadoop Distributed file System (HDFS)

A distributed file system for storing large amounts of data.

NoSQL Databases: Databases like HBase, Cassandra, MongoDB are designed for handling large volumes of unstructured and semi-structured data.

processing Technologies:-

Map Reduce :- A programming model for processing data in parallel across a cluster of nodes.

Apache Spark: An in-memory data processing engine for fast data processing and analytics.

Analytic technologies:-

→ Machine learning : - libraries like TensorFlow, PyTorch, and scikit-learn are used for building machine learning models.

→ Data mining : Technologies like clustering, decision trees, and association rule mining are used for discovering patterns in large datasets.

Top analytical tools

1) Description : A data visualization tool for creating interactive dashboards and reports.

2) Use case : Data visualization, reporting and business intelligence.

3) Power BI :
Description : A business analytics service by Microsoft for data visualization and reporting.

Use case : Data visualization, reporting and business intelligence.

3) Apache Spark.

Description : An in-memory data processing engine for fast data

processing and analytics.
use case : Real time data processing, machine learning and data analytic.

a) Google Analytics.

Description : A web analytics service for tracking website traffic and behaviour.

Use case : web analytics, user behaviour analysis and marketing optimization.

b) Explain about NoSQL and its types with an example.

→ NoSQL (NOT ONLY SQL).

NoSQL databases are designed to handle large volumes of unstructured and semi-structured data. They offer flexible schema designs, high scalability, and high performance.

few features of NoSQL databases are as follows.

i) They are open source.

ii) They are non-relational.

iii) They are distributed.

iv) They are schema-less.

- 5) They are cluster friendly.

Types of NoSQL Databases:-

We have already stated that NoSQL databases are non-relational.

They can be broadly classified into the following:

1) Key Value : If maintains a big hash table of keys and values.

For example: Dynamo, Redis, Riak etc

2) Simple key-value pair in key-value database

Key	Value
First name	Jessimonds
Last name	David

3) Document : It maintains data in collections, constituted of documents.

For example: MongoDB, Apache couchDB, couchbase, marklogic etc.

4) Column : Each storage block has data from only one column.

For example: cassandra, HBase etc

5) Graph : They are also called

network database ? A graph stores data in nodes. For example: Neo4j, HyperGraphDB etc.

Module - 2: Big Data

Q 1) Explain about Hadoop? Why Hadoop and history of Hadoop

→ Hadoop is an open source, distributed computing framework that enables the processing and storage of large volumes of data across a cluster of nodes. It was designed to handle massive amounts of data known as big data, provides a scalable, fault tolerant, and flexible solution for data processing and analysis.

Hadoop

Apache open-source software framework

Inspired by

- Google mapReduce.

- Google file system.

Hadoop distributed file system

- MapReduce

Fog :- Hadoop

Why Hadoop ?

Hadoop was created to address the limitations of traditional data processing systems, which were unable to handle the rapidly growing volumes of data. Hadoop's distributed architecture allows it to scale horizontally by adding more nodes to the cluster, as data volumes grow.

- Handle unstructured data, process and store data in various formats such as text, images, and videos.

- Provide fault tolerance by automatically replicating data and tasks across nodes to ensure data availability and processing consistency.

History of Hadoop

Hadoop was founded by Doug Cutting and Mike Cafarella in 2004, based on the Google File System (GFS) and MapReduce papers. The project was initially developed at Yahoo!

and later became an Apache software foundation project.

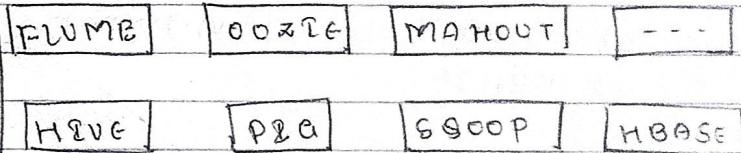
- 2004: Doug Cutting and Mike Cafarella started working on Hadoop at Yahoo!

- 2006: Hadoop graduated from the Apache incubator and became a top-level project.

- 2011: Apache Hadoop 1.0 was released, marking a significant milestone in the project's development.

Q2) List Hadoop core components and explain with appropriate diagram

Hadoop Ecosystem



core components

MapReduce programming

Hadoop distributed file system(HDFS)

Fig 2: Hadoop components

Hadoop core components

1) HDFS

→ storage component

→ distributed data across several nodes

→ natively redundant

2) MapReduce

→ computational framework

→ splits a task across multiple nodes

→ process data in parallel

→ Hadoop Ecosystem :- Hadoop Ecosystem

- system wide support projects to enhance the functionality of hadoop

- hadoop core components + these Ecosystem projects are as follows:

1) HIVE

5) PLUMBER

2) PIG

6) OOZIE

3) SQOOP

7) MAHOUT

4) HBASE

Q3) Differentiate between RDBMS & Hadoop

RDBMS

System

RDBMS

Hadoop

Relational Database management system

- flat structure

- element system

- test

Aka

Data

suitable for
structured data

suitable for
structured,
unstructured
data, support

- variety of
data formats

in real time
such as XML,
JSON, text

based flat
file format
etc

Analytical
big data
processing

Choice

when the data
needs consistent
relationship

Big data pro-
cessing which
does not
require any
consistent
relationship
b/w data

Choice

when the data
needs consistent
relationship

Big data pro-
cessing which
does not need
any relationship

		any consiste -nt relations -hip blw data
processor	needs expensi -ve hardware. or right -end processors to store huge volumes of data	In a Hadoop cluster , a node requires only q processor , a net -work card , and few hard drive -s & cost around \$ 4000 per ter objects -abytes of storage

Q4) Explain HDFS and its features

Hadoop Distributed File System (HDFS)

HDFS is a distributed file system designed to store and manage large volume of data across a cluster of nodes. It provides a scalable, fault-tolerant and flexible solution for data storage and processing.

- Strong

Key Features of (HDFS) :-

→ HDFS is a distributed file system designed to store and manage large volume of data across a cluster of nodes. It provides a scalable, fault-tolerant and flexible solution for data storage and processing.

Key Features of HDFS

* Distributed Storage : HDFS stores data across multiple nodes in a cluster, allowing for horizontal scaling and increased storage capacity.

* Fault Tolerance : - HDFS replicates data blocks across nodes to ensure data availability and durability in case of node failure.

* High throughput : - HDFS is optimized for high throughput data access, making it suitable for large-scale data processing and analytics.

* Flexible data model : - HDFS supports

Date _____

Page _____

stores a variety of data formats, including structured, semi-structured and unstructured data.

HDFS Features.

- * **Data Replication** :- HDFS replicates data blocks across nodes to ensure data availability and durability.
- * **BLOCK SIZE** :- HDFS stores data in fixed size blocks (typically 128 MB or 64 MB) which allows for efficient data processing and storage.
- * **Data integrity** :- HDFS uses checksums to ensure data integrity and detect corruption.
- * **Scalability** :- HDFS is designed to scale horizontally allowing for the addition of new nodes as storage needs grow.

