

14/2/23

Unit-11

Cassandra

* Cassandra is a distributed database, highly scalable, designed to manage large volume of structured data.

* provides high availability, no single point of failure.

* It is a NoSQL database support.

— provides schema free

→ easy replication

→ simple API

→ consistent.

* simplicity

* horizontal scaling

* it is scalable, fault tolerant

* column oriented database

* created by Facebook.

* used by fb, twitter, Cisco, Rackspace

Features

* Elastic scalability

* Always on architecture

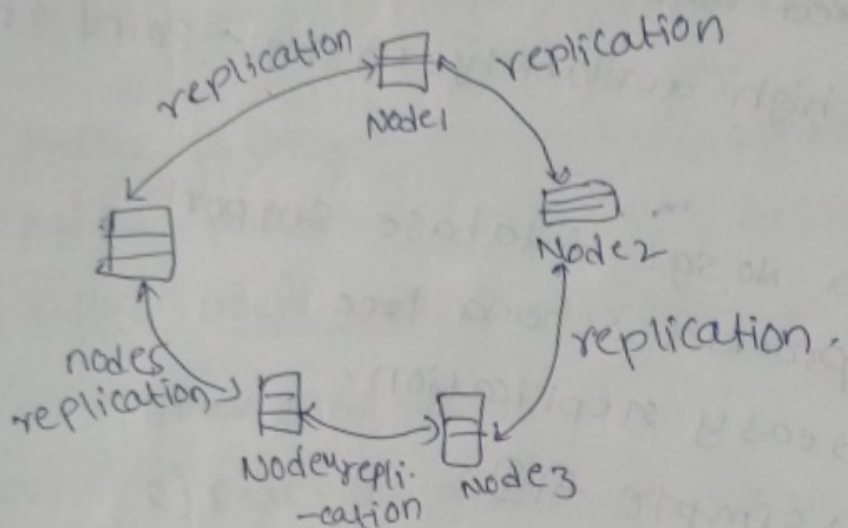
* flexible data storage

* easy data distribution

* + transaction support.

Goal:- To handle bigdata work load across multiple nodes

- All nodes in the cluster plays same role
- Each node in the cluster perform read and write request



* if the node reported with a ~~out of~~ data value, cassandra, will return the most recent value to the client.

* constant, using Gossip protocol for detect error nodes.

* After returning the most ~~resend~~ value cassandra performs a need repair in the background to update the state value.

~~create~~ create current

* Components of cassandra:-

1. Nodes
2. data centre

3. cluster

4. commit log

5. memtable

6. SS table

7. Bloom filter.

Components of Cassandra.

1. Nodes - It is the place where the data is stored.

2. Data center - collection of related nodes.

3. cluster - contains one or more data center.

4. commit log - is a crash recovery mechanism

5. memtable - is a memory resident data structure

6. SS table - it is a disk file to which the data is flushed from mem table.

7. Bloom filter - is a algorithm for testing whether an element is a member of a set or not

Cassandra Query Language (CQL)

* CQL treats the database as a container of tables

* programmer works with cql as a prompt to work.

write operation

- * Every write activity is captured through the commit log.
- * Later the data will be captured and stored in the memtable.
- * Whenever the memtable is full the data is written into sstable datafile.

Read operation

- * During read consistant gets value from memtable and checks the bloom filter to find appropriate sstable that hold the actual data.

Supports the following data.

- * Time series data.
- * Marketing data.
- * Financial data
- * IoT data
- * Graph data.

Statistics about facebook

mysql > 50 GB data.

write average ~ 300ms

Read average ~ 350ms

20/2/23

Mongo DB:

- * Mongo DB is a cross platform document oriented db that provides high performance, high availability and easy synchronization and scaling which works on collection of documents.

Features

- * It is an intelligent operational db platform
- * It belongs to a aggregation framework.
- * It efficiently rep reduces the data.
- * Converts table into Tson document, then further converts it into Bson doc.
- * Allows Sharding
- * Schema less
- * Capped Collection.
- * Stores data in cloud efficiently.

Advantages

- * schema less
- * No complex joins
- * Easy to scale out
- * Deep Querying ability.

Where Mongo DB is used?

- * Big data
- * Content Management and delivery.

* Mobile and social infrastructure

* Data hub

* Data Models in Mongo DB:

1) Embedded data model :-

2) Normal data model.

Data types in Mongo DB:

* String

* Integer

* Boolean

* null

* Symbol,

* Binary

* Regular expression,

15/12/23

Spark

* Spark is an open source cluster computing framework, purpose is to handle real time generated data.

* Started in 2009 in US Berkeley R&D Lab
2010 - open source, 2013 - Apache Foundation adopted.

Features

* Used to perform batch processing, stream processing - using Apache Storm/ky.

* Used for interactive processing (use of Neo4j, Apache Graph).

* Fast - use DAG scheduler & query optimizer
Physical execution engine.

* Easy to use - write application in Java, Scala, Python, R, and SQL provides high level operations

* Generality - provides collection of libraries.

* Light weight - it is a unified analytics engine used for large scale data processing

* Runs everywhere - run on hadoop Apache Mesos
Kubernetes, standalone and cloud.

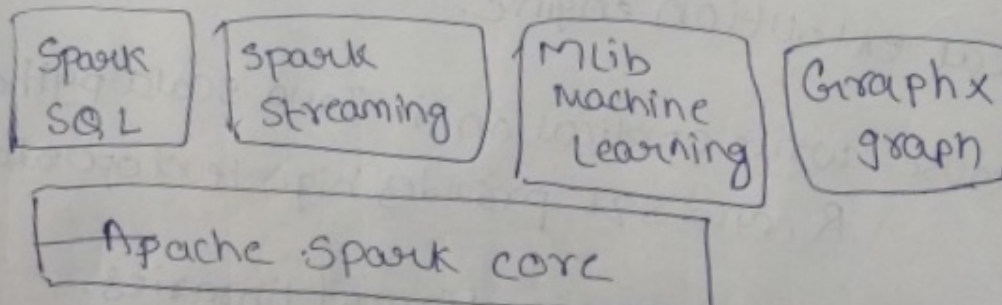
Usage of spark

1. Data Integration - to fetch consistent data from System, use ETL (Extract, Transform, Load)
2. Stream processing - work with real time generated data such as log files.
3. Machine learning.
4. Interactive analysis → handle the data interactively.

Why spark

Support multiple languages, Realtime processing
Fault tolerance in ^{memory} ~~memory~~ processing, Reliability,
cost efficient, Graph x.

Components of spark



1. Spark core:- provides execution platform for all Spark applications.
- provides in-memory processing.

2. spark SQL: provides new data abstraction called schema RDD which provides support for structured and semi structured data.

16/12/23

* Components of Spark

1. Spark core
2. Spark SQL
3. Spark Streaming
4. MLlib
5. Graph.

Resilient Distributed Dataset (RDD)

* It is a fundamental unit of data in Spark to perform parallel operation

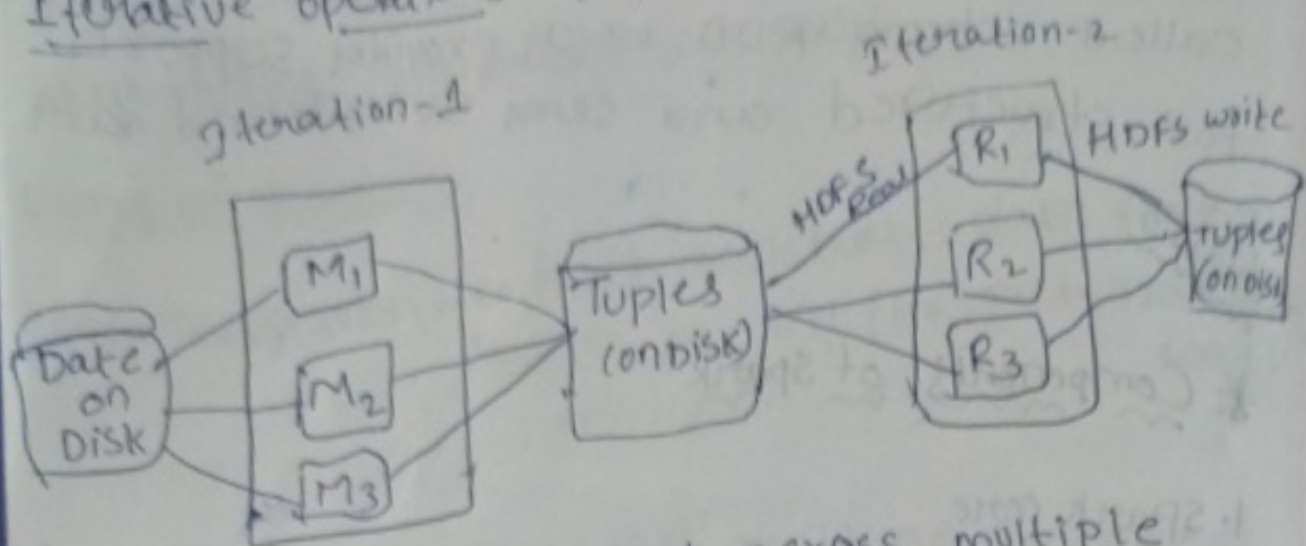
* It is immutable, Readonly, partitioned collection of records.

Two way to create RDD

1. parallelising an existing collection in your driver program
2. Referencing dataset in an external storage system such as HDFS, Hbase.

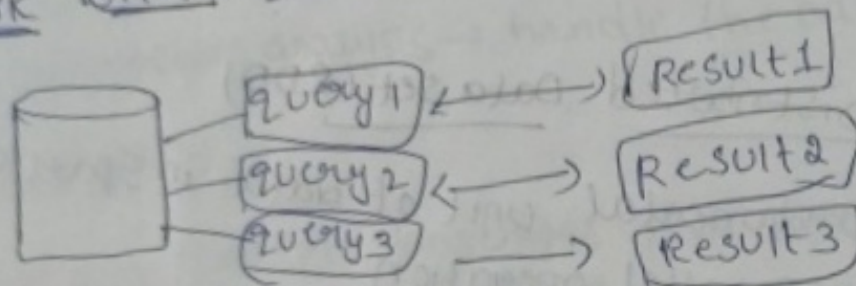
* Both Iterative and Interactive application requires faster data sharing across parallel jobs.

Iterative operations on map reduce:



* Reuse intermediate result across multiple computation in multistage application.

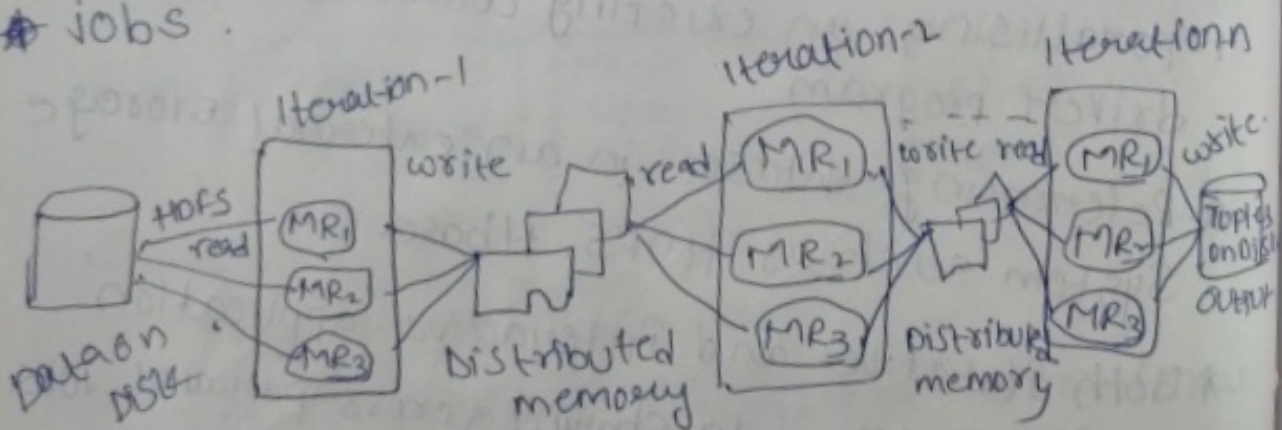
Work with interactive queries



Input from stable storage.

* In RDD supports in memory processing computation

* It stores the state of memory as an object across the jobs and object is shorable b/w jobs.



~~Q6~~

Q6. Data Visualization, It's importance, Tableau tools

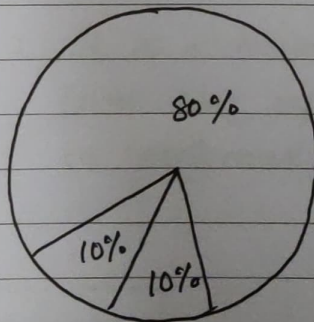
Data Visualization is the process of creating a graphical representation of data using various methods like pie charts, bar charts etc.

It is used to reveal the patterns and trends hiding under complex data

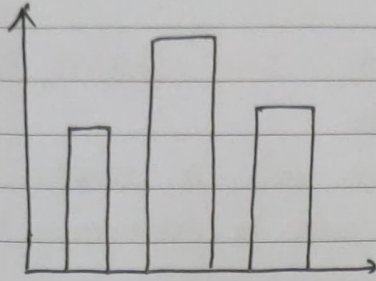
It leads to better and faster decision making

Different visualization Techniques include the following-

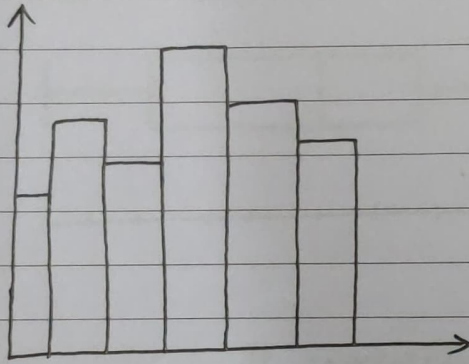
- **Pie chart** : It is divided into sections to represent the proportions of the total amount.



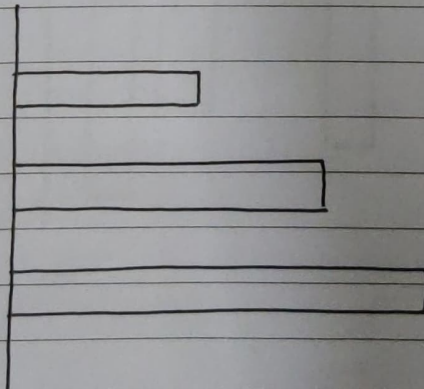
- **Bar chart** : It is used to represent different categories with the help of vertical bars



- **Histogram** : It is also divided into vertical bars representing how often something happens within a range.



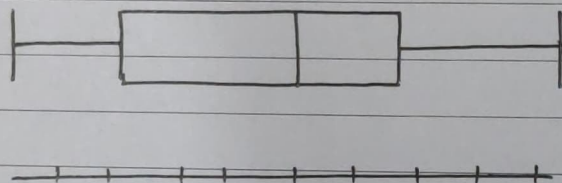
- **Gantt chart** : It is a horizontal bar chart. Each bar represents a task that is performed in a certain time.



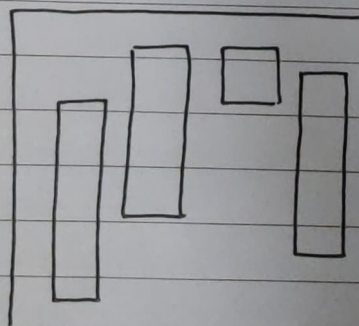
- **Heat map** : It uses different colours to show the contribution of each category.

- **Box and whisker plot** : It is in the form of distribution of data. It shows a five number summary that comprises of the following

- minimum
- first quartile (Q_1)
- median
- third quartile (Q_3)
- maximum



- **Waterfall chart** : It represents the contribution or decrease in value with time in the form of vertical bars.



Importance of data visualization →

- Using graphs and charts to visualize large amounts of data is more comfortable than studying spreadsheets and reports
- It helps to identify patterns and trends that are hidden in complex data
- It helps to make faster and better decisions
- It helps to understand the story in a single glance
- It helps us explore business insights
- It helps us identify errors in the data.

Tableau →

Tableau is a powerful visualization tool used to analyze raw data by presenting it in a visual manner.

No technical or programming skills is required to work with tableau

Tableau tools →

1. Tableau Desktop → establishes connectivity ^{btwn} with the data warehouse and other kinds of files

based on the connectivity, it is of 2 types

→ Tableau desktop personal

→ Tableau desktop professional

2. Tableau Public → Work created cannot be stored locally.

It is kept on cloud where it can be accessed and viewed by all.

3. Tableau Online → Data is stored in servers hosted on the cloud, which is maintained by the tableau group

4. Tableau Server → It can share the workbook visualization which is created by tableau desktop.

5. Tableau Reader → It allows us to view the visualisations, but we cannot edit or write to the data.