

Big Data

Big Data - is data that exceeds the processing capacity of conventional database systems. Collection of data sets so large and complex that it becomes difficult to process using on-hand data management tools.

Big data is high-volume, high-velocity and high-variety information assets that demand cost-effective, innovative forms of information processing for enhanced insight and decision making.

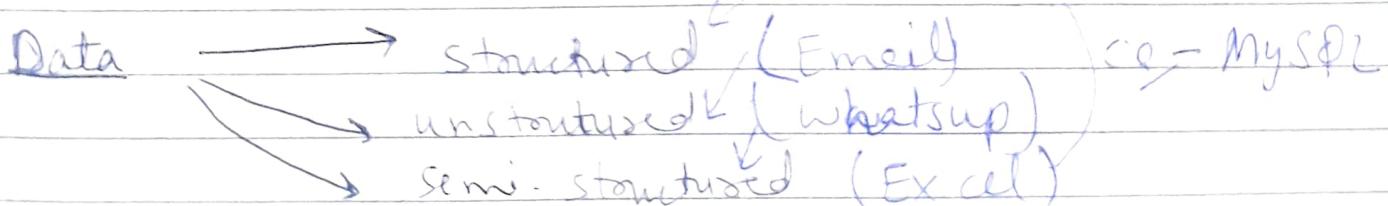
Where Is This 'Big Data' Coming From?

social data, machine data, transactional data.

Big Analytics

process of examining data - typically of a variety of source types, volumes and / or complexities - to uncover hidden pattern, unknown correlations and other useful information
 ⇒ intent is to find business insights.

- Deeper insights - Rather than looking at classifications groups, you will insights into all the individual.
- Broader insights - take account in all the data to produce more accurate insights.



Structured data - data stored in rows and columns, mostly numerical, where meaning of each data item is defined. It consists of 10% of total data and is accessible through database management systems or hadoop

Unstructured data - data of different forms like eg- text, image, video, document etc. It can also be in the form of customer complaints, contacts, or internal emails. It comprises of 90% of data.

Unstructured data cannot be stored using traditional relational databases. Storing data with such a variety and complexity requires the use of adequate storage systems, commonly referred to as NoSQL database. Eg - MongoDB

- semi-structured data and unstructured data can only be handled through Hadoop system.
- CRUD CYCLE is not supported by Hadoop.

Semi-structured data - is a form of structured data that does not obey the tabular structure of data model associated with relational databases or other forms of data tables, but nonetheless containing tags or other markers to separate semantic elements and enforce hierarchies of records and fields with data.

What happen and Why? ↪

Diagnostic data - is data that is automatically recorded by S/w for the purpose of troubleshooting problems.

Descriptive data - form of advance analytics which examines data or content to answer the question.

Predictive Analytics - data describe the basic features of the data in a study and provide summary about the sample and the measures.

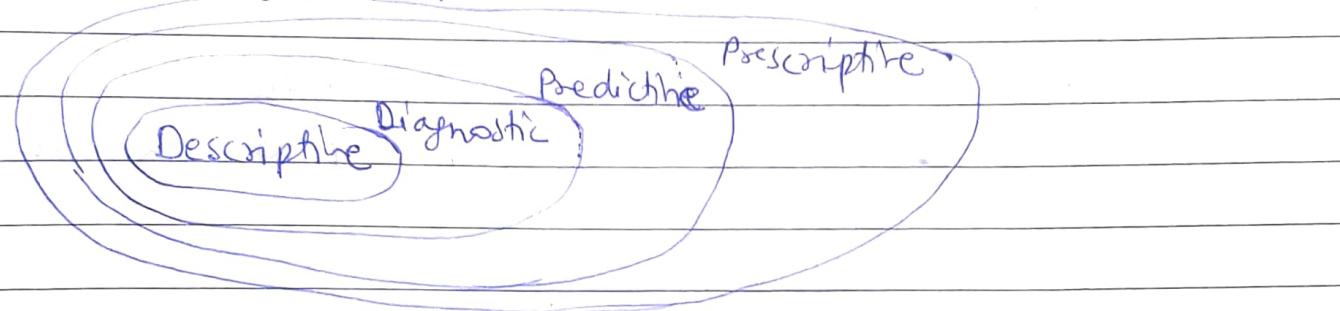
Predictive Analytics - is the use of data, statistical algorithms and machine learning techniques to identify the likelihood of future outcomes based on historical data.

Descriptive - analytics addresses the issue of what happened.

Diagnostic - analytics answers the question of why something happened.

Predictive - analytics describes what is likely to happen.

Prescriptive - analytics prescribes what steps to take to avoid a future problem.



→ visualization tools (weka Tool)

	<u>Structured Data</u>	<u>Semi Structured Data</u>	<u>Unstructured Data</u>
Analogy	Relational Database Table Structured query allows complex joins	XML / RDF Queries over anonymous nodes are possible	Character / binary data Only textual queries possible
Scalability	Scaling DB scheme is difficult	Simple	Very scalable

3 Dimensions / characteristic of Big data

3V are three defining properties or dimensions of big data. Volume refers to amt. of data, variety refers to the number of types of data and velocity refers to the speed of data processing.

Volume - Size of available data has been growing at an increasing rate.

Velocity - Data is increasingly accelerating the velocity at which it is created and at which it is integrated.

Date / /

Variety - presents an equally difficult challenge. The growth in data sources has fuelled the growth in data types. In fact, 80% of world's data is unstructured.

Variety - refers to the biases, noise and abnormality in data. Is the data that is stored meaningful to the problem being analyzed. (inconsistencies and uncertainty in data)

Value - Bulk of Data having no value is of no good to the company, until you turn it into something useful.

Problem with Big Data

- i) Lack of Understanding ii) Security Gaps
- iii) Low Quality and Inaccurate Data iv) Lack of skilled workers

Dataset - is a collection of data. In case of tabular data a dataset corresponds to one or more database tables, where every column of a table represents a particular variable and each row corresponds to a given record of the data in question.

Case Study

Manufacturing Big Data Use Cases

i) Predictive Maintenance

- helps in predict equipment failure.
- deploy maintenance more cost effectively.
- data can be get from sensor data, engine temp etc
- to predict the remaining optimal life of systems and components to ensure that they perform within specifications.

Challenges

- integrate data coming from different formats and identify the signals that will lead to optimizing maintenance.

ii) Operational Efficiency

- analyze and estimate production processes, proactively respond to customer feedback, and anticipate future demands.

Challenges

Data teams must balance the data volume with the growing number of sources, users, and application.

iii) Production Optimization

- understanding the flow of goods in production!
- Using Data analysis to increase productive time.

Challenges

Combining different kinds of data.

b) Retail Big Data Use Cases

a) Product Development -

By classifying key attributes of past and current products and then modeling the relationship b/w the attributes, building predictive models for new products and services. Deeper analysis to plan, produce and launch new products.

Challenges:

- Creating segments acc to customer behaviour identify sophisticated use patterns and behaviour and map them to potential new offerings.

b) Customer Experience

- gathering data from social media, web visits, calls and other ^{data} source. Big data analysis can be used to deliver personalized offers, reduce customer churn, and proactively handle issues.

Challenges high val

- gathering data from various resource can be difficult. Once the data is integrated, path only can be used to identify experience paths and correlate them with various set of behaviour

c) Customer Life Time Value

- identify your best customers, marketing them with special offers. Sales team devotes more time to them.

Challenge

- Identify your high-value customers, analysing a high vol. of transaction data and create sophisticated model that examine past behaviour and predict future actions.

Healthcare Big Data Use Cases

a) Genomic Research

Using big data, researchers can identify disease genes and biomarkers to help patients pinpoint health issues they may face in the future.

Challenge

- running complex algo. on the data is complicated and can require long processing times

b) Patient Experience and Outcomes

With big data, health care organization can create 360° view of patient care as patient moves through various treatments and departments.

Challenge

- Improving the patient experience requires a large volume of patient data (it can be multi-structure)

c) Claims Fraud

Big data helps healthcare organizations detect potential fraud by flagging certain behaviours for future examination.

Challenge

- Claims fraud analytics is a complex process that involves integrating different datasets, analyzing the claims data, and identifying complex fraud patterns.

Oil and Gas Big Data Use Case

a) Predictive Equipment maintenance

- predict the remaining optimal life of their system and components, ensuring that their assets operate at optimum productive efficiency.

Challenge

Collecting data of varying format.

b) Oil Exploration and Discovery

Data generated from seismic monitors can be used to find new oil and gas sources by identifying traces that were previously overlooked.

Challenge

- integrate and analyze an enormous volume of structure data

c) Oil Production Optimization

- measures well production to understand ~~usage~~ rates. Analyse why well outputs aren't tally with their predictions.

Challenges

- analyzing large volume of data

e) Telecommunications Big Data Use Cases

a) Optimize Network capacity

Network usage analysis can help companies identify areas with excess capacity and resolute bandwidth as need, BD helps in planning for infrastructure investment and design new services that meet customer demand.

Challenges

- for network analytics requires analyzing a high volume of call detail records.

b) Telecom customer churn

- predict overall customer satisfaction by considering service quality, convenience.
- Set up alert when customers are at the risk of churning - and take action with retention campaign and proactive offers.

Challenges

- require to analyzing past and current data to create new model to predict churn.

C) New product Offerings

- design new products and features based on customer behavior. (for future offerings)

challenge

analyzing high-volume product-log data in different format.

Crowdsourcing is the practice of engaging a 'crowd' or group for a common goal - often innovation, problem solving or efficiency. It is powered by new technologies, social media. Crowd sourcing can take place on many different levels and across various industries.

- collectively contribute - whether with ideas, time, expertise or funds - to a project or cause,
- powerful business marketing tool - promoting and growing the company by creativity.

Advantages

- i) Numerous idea from numerous people
- ii) cheap - similar to outsourcing, It is used to cut costs. when you don't have to employ people and pay them a wage, whether they are working or not, it definitely lowers costs.
- iii) Fast - take less time to find right person to do the job

Disadvantage

- i) Quality could be questionable - when you hire numerous people to do a job if you aren't lucky enough to hire the professionals.

ii) Unreliable way to get a job done

iii) Confidentiality

Inter Firewall & Sans Firewall Analysis

Cloud Computing

Definition • Provides resources (storage, computing dB, monitoring tools etc) on demand.

Reference • It refers to internet services from SaaS, PaaS to IaaS.

Use for • Use to store data and info. on remote servers.

(range of network of cloud servers over the internet)

Big Data

• Provide a way to handle huge volumes of data and generate insights.

• It refers to data, which can be structured, semi-structured and unstructured.

• used to describe huge vol. of data and info (discover undiscovered patterns)

Cloud Computing Services

i) Infrastructure as a Service (IaaS) - Here the service provider offers entire infrastructure along with the maintenance related tasks. Organization make use of unlimited storage potential of the cloud infrastructure. They can expand and shrink their storage space as needed w/o worry.

ii) Platform as a Service (PaaS) - in this service, the cloud provider offers resources like queues, databases etc. However, the responsibility of configuration and implementation related tasks depends on the customer.

iii) Software as a Service (SaaS) - This service is the most facilitated one which provides all the necessary setting and infrastructure. provider IaaS for the platform and

infrastructure are in place. (provides self service capabilities to users with scalable features to increase usage on requirement)

Cloud computing types

- i) private - used for a single organization, can be internally or externally hosted.
- ii) public - provisioned for open use for the public by a particular organization, who also hosts the service.
- iii) Hybrid - composition of 2 or more clouds (private, public or community), that remain unique entities but they bound together.
- iv) community - shared by several organization, typically externally hosted.

GZG, GZB, GZC, GZE

R-language
SPSS S/W
IBM Watson

] virtual Tool.

Mid Sem

RDB
ACID

NoSQL
BASE
Basically Available, Soft-state, Eventually Consistent

Type of db

Schema

Db categories

Complex queries

Hierarchical Data

Storage

Scalability

Language

Basic Properties

SQL
Relational database

Pre-defined

Table based dB

Good for CQ

Not the best fit

NoSQL

Non-relational Database

Dynamic

Document-based dB,

key-value stores,

Not good for CQ
best fit

Vertically scalable

Structured query language

Based on ACID property

Horizontally Scalable

Unstructured Query lang.

Based on CAP Theorem

NoSQL database provides a mechanism for storage and retrieval of data that is modelled in means other than the tabular relations. (Not only SQL)

Theory of NoSQL : CAP

• Many Nodes

• Nodes contain replicas of partitions of the data.

Consistency

→ All replicas contain the same version of data.

→ Client always has the same view of the data (no matter what node)

Availability

→ System remains operational on failing nodes.

→ All clients can always read and write.

Partition tolerance

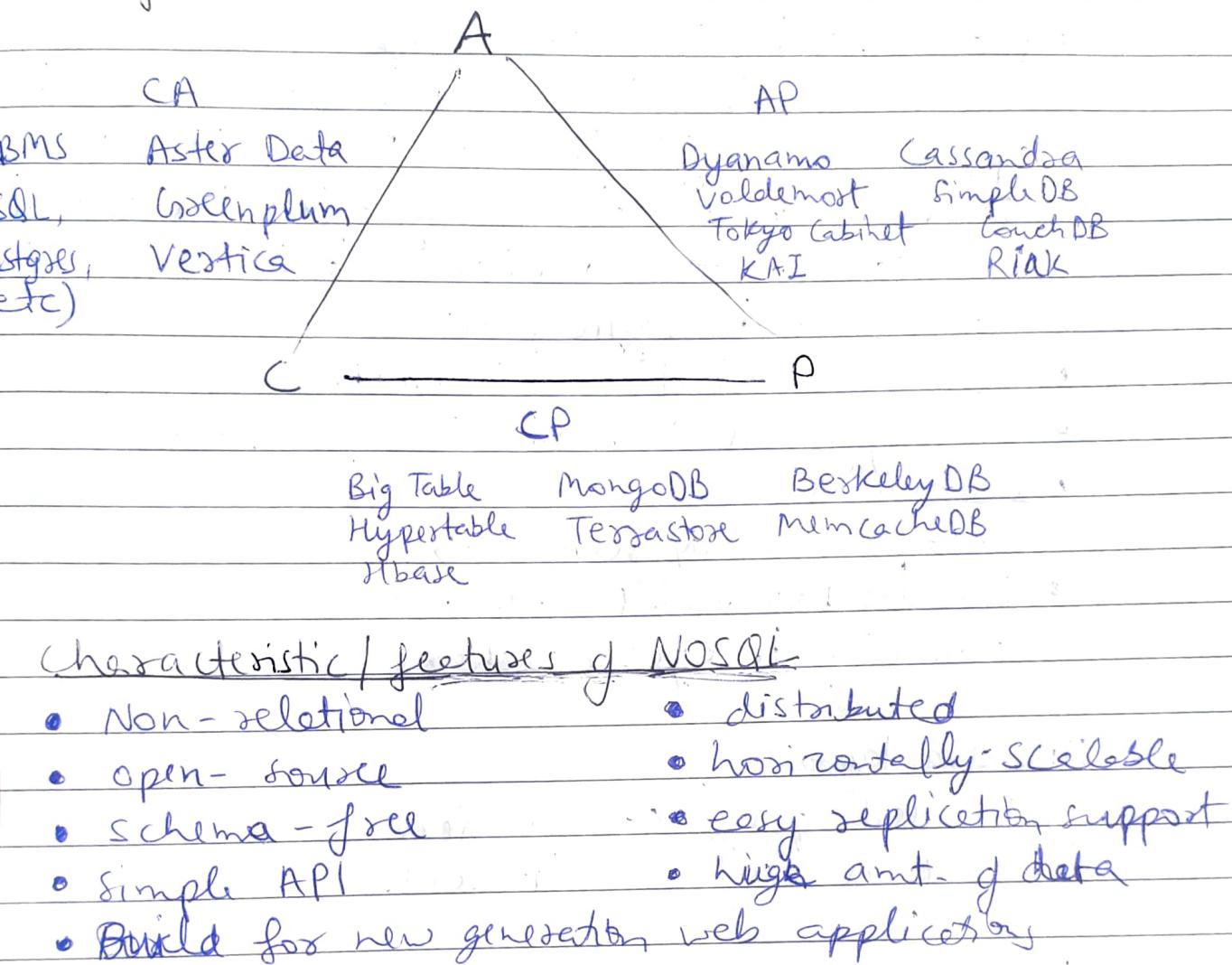
→ multiple entry points

→ System remains operational on system split (communication malfunction)

→ System works well across physical network partitions.

Types of NoSQL

Use of NOSQL



Characteristic / features of NOSQL

- Non-relational
- open-source
- schema-free
- simple API
- built for new generation web applications
- distributed
- horizontally scalable
- easy replication support
- huge amt. of data

NOSQL is an approach to database design that can accommodate a wide variety of data model.

NOSQL which stands for 'Not only SQL' is an alternative to traditional relational database.

NOSQL is specially designed for large amt. of data stored in distributed environment.

Need - bcz data are not structured e.g. data at Facebook, googl.

- Advantage -
- i) supports all type of data
 - ii) No static scheme
 - iii) Faster data processing
 - iv) low operational cost
 - v) support Distributed system

- Elastic scaling - is the ability to automatically add or remove network infrastructure based on changing application traffic patterns. Elastic load balancers will scaling is used to automatically adjust the amt of resources that are allocated to deliver an application in response to changes in traffic patterns.
- Big Data - Vol of data that are being stored have increased.
- Goodbye DBA's - Automatic report, distribution, tuning...
- Economics - cheap commodity servers \rightarrow less costs per transaction
- Flexible Data Models - allows the direct representation of data with irregular schema.

Disadvantages

- i) Limited query capabilities.
- ii) Not a defined standard.

Challengers

- i) Maturity \rightarrow still in pre-production phase
 \rightarrow key features yet to be implemented
- ii) Support \rightarrow limited resources or credibility
- iii) Administration
 - Require lot of skill to install and effort to maintain.
- iv) Limited ad-hoc querying
- v) Expertise - Few no. of NoSQL experts available in the market.

RBAMs

- regular backup
- access through master server.
- non-programmers writing queries.
- data updates are frequent

NoSQL

- replication
- sharding across multiple nodes
- only programmers writing queries
- write-once & read multiple

predictable, linear growth

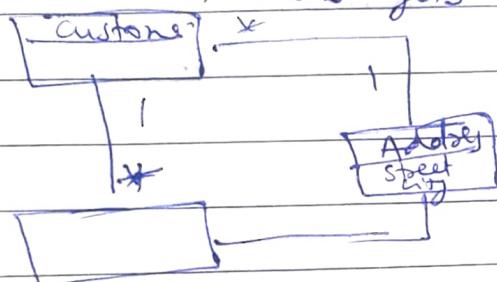
• only programs writing code.

DOMS Page No.

Date / /

Data Model

A aggregate - is collection of related objects that we treat as a unit. When a unit of data is asserted from the db gets all the selected data along with it.



These units of data form the boundary for ACID operation with the db.

Type of NoSQL database

key values	Document oriented.	wide column family	Graph databases.
key value ↑ lookup	eg - id age 1 10 2 20 ↓ document & collections (tables) eg - MongoDB	columns collect ↓ column family key [personal] [professional] emp-id → home → exp → loc → visa → age ↓ selections	eg - Neo4j, TitanDB P1 ① (V1) → Post tagged ② V2 ③ V3 ↓ selections
Every single item in the db is stored as an attribute name (or "key") together with its value stores	Documents can contain many different key-value pairs	optimized for queries over large data sets and stores cols of data together, instead of rows	are used to store info about network

BigTable - is a distributed storage system for many structured data at Google.

→ scalability, high performance

→ sparse, distributed, multidimensional sorted

API - • functions for creating and deleting tables
column families etc.

Master Server - responsible for assigning tablets to tablet servers and balancing the tablet server load, garbage collect

Tablet - set of consecutive rows of a table and is the unit of distribution.

Tablet server - • Each tablet server manages a set of tablets
• clients communicate directly with the tablet server
• Tablet server handles read and write requests
• Also split tablets that has grown too large

Materialized view - is a database object that contains the results of a query.

MongoDB → CP

It is a free and open-source cross platform document-oriented database program.

It uses JSON-like documents with schemas.

Uses BSON format

It is a binary form for representing simple data structures

BSON → "Binary JSON"

Which PL can be used with MongoDB?

JS, Python, Ruby, Perl, Java, C#, C++.

Functionality of MongoDB

- Dynamic Schema
- Document-based database
- Built-in horizontal scaling
- Secondary indexes
- Query language via an API
- Sharding
- Fast In-place Updates

Why use MongoDB?

- Simple queries
- Rich queries
- Easy & fast integration of data
- Document Oriented Storage
- No ER diagram

RDBMS

Table, view

Row

Col

Join

MongoDB

collection

Document [BSON]

Field

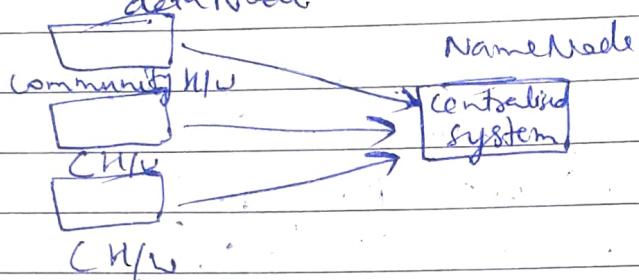
Embedded Document

Map Reduce

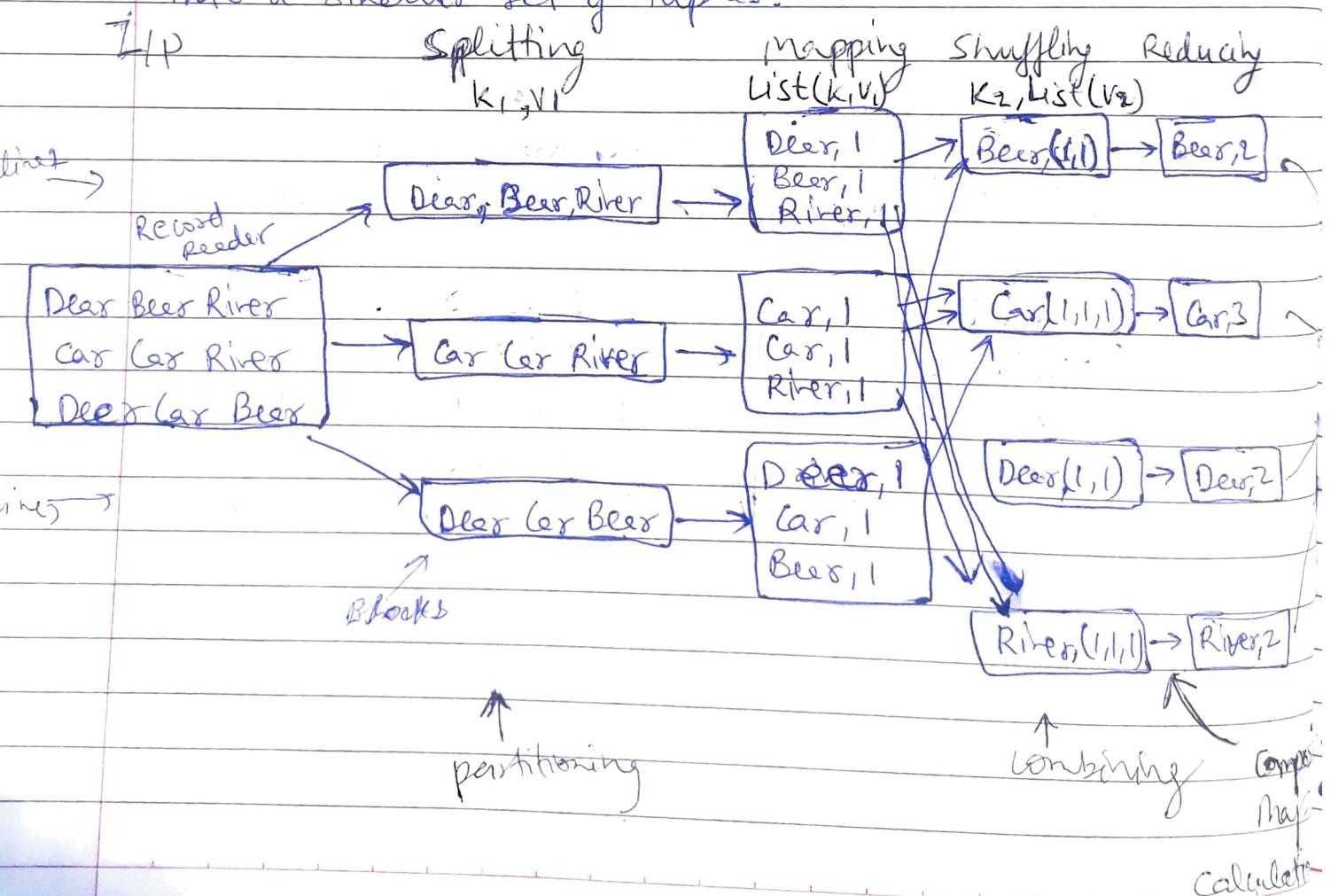
commodity clusters

Map Reduce

(Simple programming model for data-intensive computing on large Map Reduce divides a task into small parts and assigns them to many computers. Later, the results are collected at one place and integrated to form the result data.



- The Map task takes a set of data and converts it into working another set of data , where individual elements are broken down into tuples (key - value pairs)
 - The Reduce task takes the o/p from the Map as an input and combines those data tuples (key - value pairs) into a smaller set of tuples.



Value	Map function	Stage	W	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{A, 0, 1\}$	(I, K)	(4)	$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{A, 1, 2\}$	$(0, 0)$	$(A, 0, 1), (A, 1, 2)$	$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{A, 0, 3\}$		$(B, 0, 5), (B, 0, 7)$	$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{A, 1, 4\}$	$(0, 1)$	$(A, 0, 1), (A, 1, 2)$	$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{B, 0, 5\}$		$(B, 1, 6), (B, 1, 8)$	$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{B, 0, 7\}$			$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{B, 1, 6\}$			$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$
$\{B, 1, 8\}$			$(1, 1)$	$O_1 O_2 1 \times 5 + 2 \times 7 = 19$

MapReduce used for?

- Index building for Google search
 - Analyzing search logs (Google Trends)
 - Statistical language translation
 - Data mining
 - Ad optimization
 - Machine learning

Map Reduce Goals

- Scalability to large data volumes
 - Cost-efficiency.

Commodity cluster computing - involves the use of large numbers of already-available computing components for parallel computing, to get the greater amount of useful computation at low cost.

↳ challenge

- Programming in a cluster is hard.
 - cheap nodes may fail
 - commodity network. → low BW

Applications

- Google Translator
 - Geographical Data
 - ~~Artificial~~ AI
 - Pdf Generation

Matrix Multiplication

1 Step 2 Step

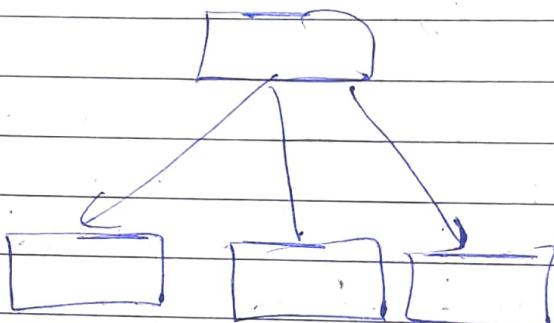
- i) Map
 - ii) Reduce

• Pdf Generation
Perform matrix multiplication (one step on following matrices)

Hadoop

- Open source framework that allows distributed processing of large datasets on the cluster of commodity hardware.
- Hadoop is a data management tool and uses Scalable storage.

→ is network-attached storage (NAS) architecture in which the total amount of disk space can be expanded through the addition of devices in connected arrays with their own resources. (H/W added as per need arises)



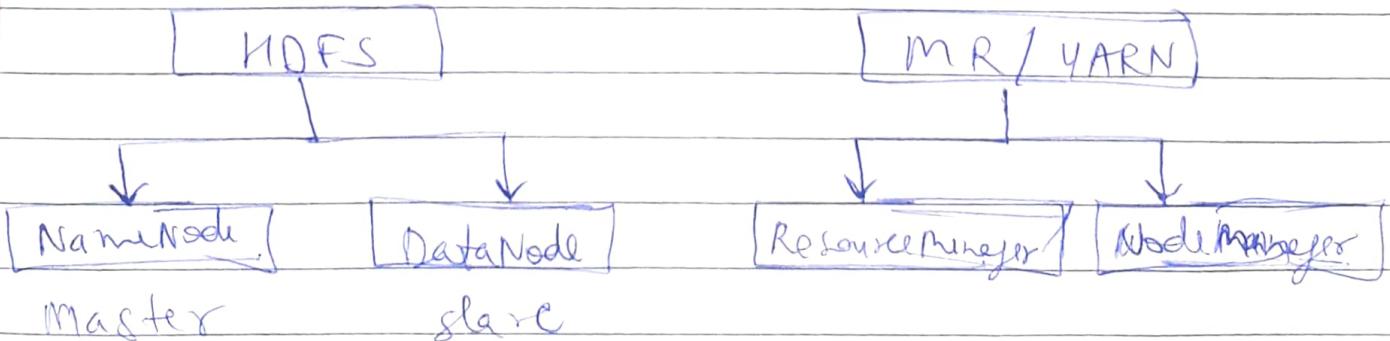
Hadoop Components

- Hadoop 1 components
 - HDFS (Hadoop Distributed File system)
 - ↳ It is storage unit of Hadoop.
 - MapReduce:
 - ↳ It is the processing unit of Hadoop.
- Hadoop 2 components
 - HDFS
 - YARN / MRV2
 - ↳ resource management unit of Hadoop

5 Hadoop Daemons

- i) NameNode - used to hold the metadata (info. about location, size of file/block). There will always be only 1 NameNode in a cluster.
- ii) Secondary NameNode - used as a backup for NameNode. It holds pretty much same information as of that of NameNode. If NameNode fails, this one comes into picture.
- iii) DataNode - stores actual HDFS data blocks. The no. of DataNode depends on your data size. The DataNode communicates to NameNode on segment basis.
- iv) Job Tracker - manages Map Reduce jobs, distributes individual tasks to machines running the TaskTracker. coordinates MapReduce stages.
- v) Task Tracker - The jobs given by Job trackers are actually performed by Task Trackers. Each DataNode will have one task tracker. Task trackers communicate with Job trackers to send statuses of the jobs.

Note - Hadoop HDFS Architecture follows a Master/ Slave Architecture, a cluster comprises of a single NameNode (Master Node) and all the other nodes are DataNodes (Slave nodes).



Running

Modes of Hadoop

- i) Standalone Operation - By default, run in a non-distributed mode (single JVM). Suitable for running MapReduce programs during development.
- ii) Pseudo-Distributed Mode - Hadoop daemons run on local machine. Each Hadoop daemon will run as a separate java process.
- iii) Fully-Distributed Mode - Hadoop daemons run on a cluster of machines. Fully distributed with minimum 2 or more machines as a cluster.

File Formats in Hadoop

- | | |
|---------------------|--------------------|
| i) Text / CSV Files | v) RC Files |
| ii) JSON File | vi) ORC Files |
| iii) AVRO Files | vii) Parquet Files |
| iv) Sequence File | |

Hadoop Data Analysis Techniques

to analyze the huge stock data being generated frequently.

- MapReduce

→ Powerful model for parallelism

→ Based on a rigid procedural structure

- Pig

→ Procedural data-flow language

→ Used by programmers and researchers

- Hive

→ Declarative SQL-like language

→ Used by analysts for generating reports.

Scale out - is a type of capacity expansion ~~consisting on the addition~~ of new H/W resources instead of increasing the capacity of already available H/W resources, such as storage.

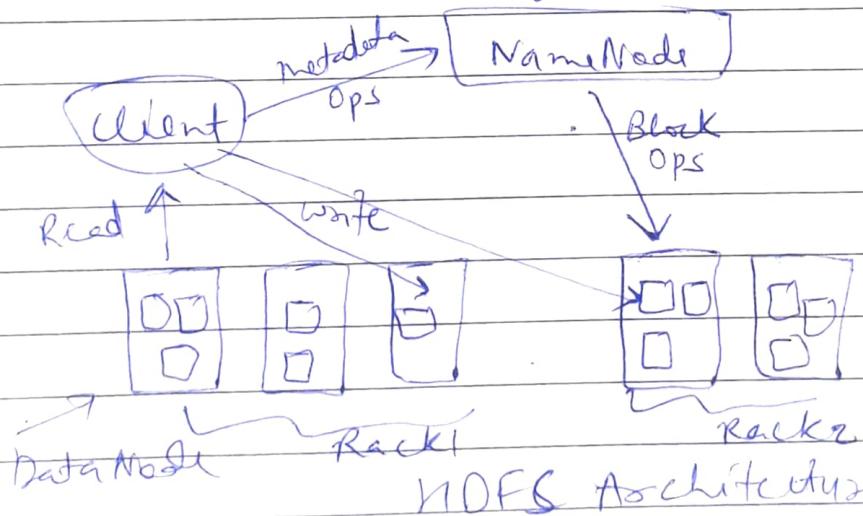
JOMS Page No.
Date

Hadoop streaming - is a utility that comes with Hadoop distribution. This utility allows you to create and run Map/Reduce jobs with any execute or script as the mapper and/or the reducer.

Hadoop Pipes - is the C++ interface to Hadoop Reduce. Hadoop Pipes uses sockets to enable task backends to communicate processes running the C++ map or reduce functions.

HDFS

- Scalability for large data set.
- Reliability to cope with H/W failure.
- Suitable for the distributed storage & processing
- good for streaming data.



Goal

- Auto fault detection and recovery
- Huge datasets.
- easy to scale

Replication of blocks for fault tolerance.

Execution

[calling `create()` by client]

↓ RPC call

Name node

→ Check for file already open →
↓ permissions

↓ Yes

written by
client

↓ [Create new file]

Client - set of programs

DOMS

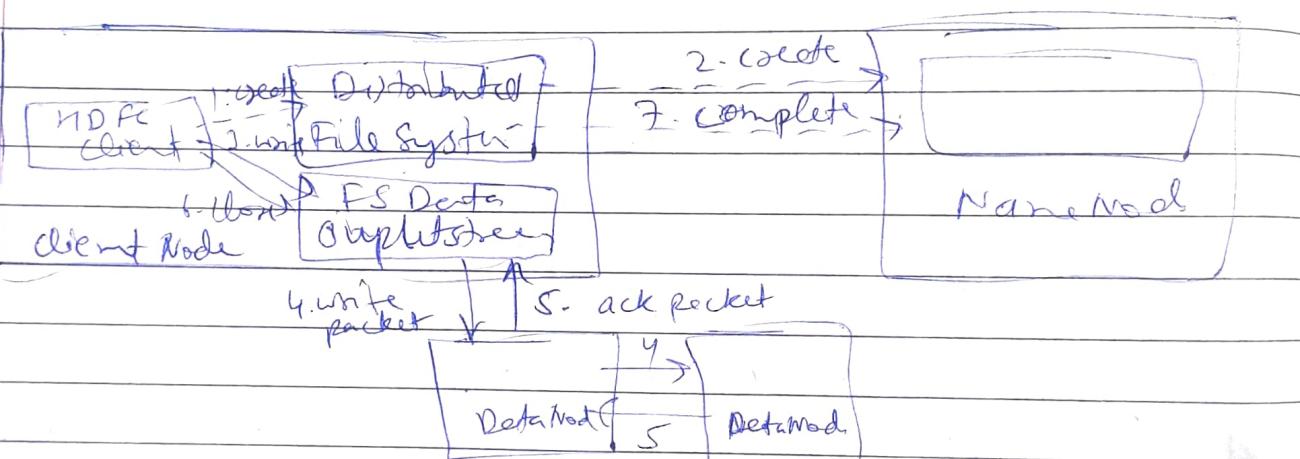
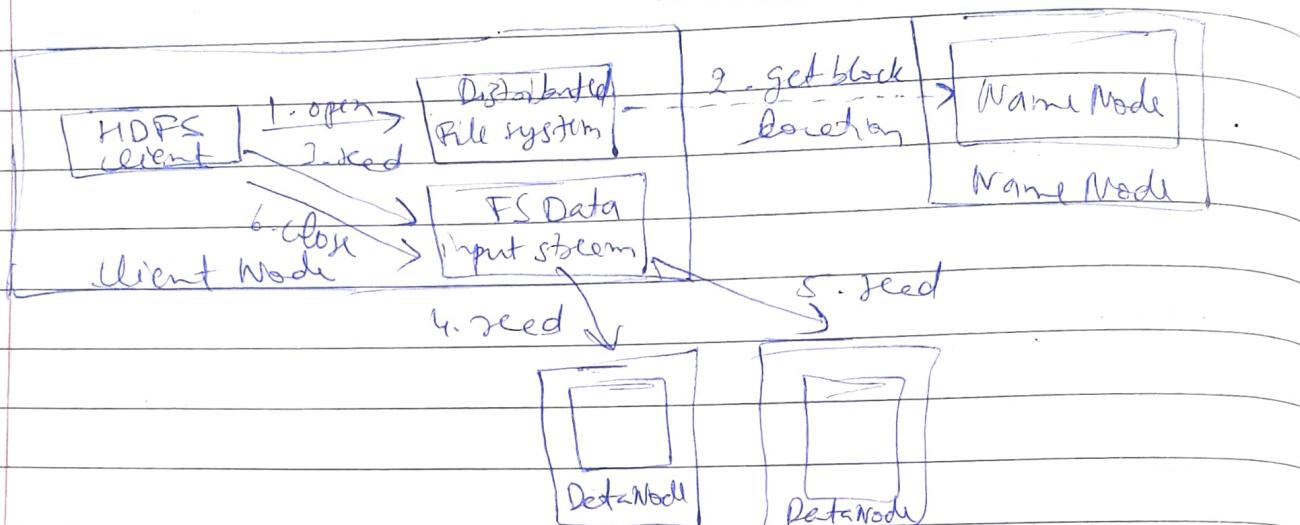
Page No.

Date

/ /

I/O in Hadoop

File Read



Data Integrity

- A HDFS client creates the checksum of every block of its file and stores it in hidden files in the HDFS namespace.
- When a client retrieves the contents of file, it verifies that the corresponding checksums match.
- If does not match, the client can retrieve the block from a replica.

Compression

- Store data in a format that requires less space than original.
- Useful in storing & transmitting the data.
- lossless, ~~lossy~~ compression.

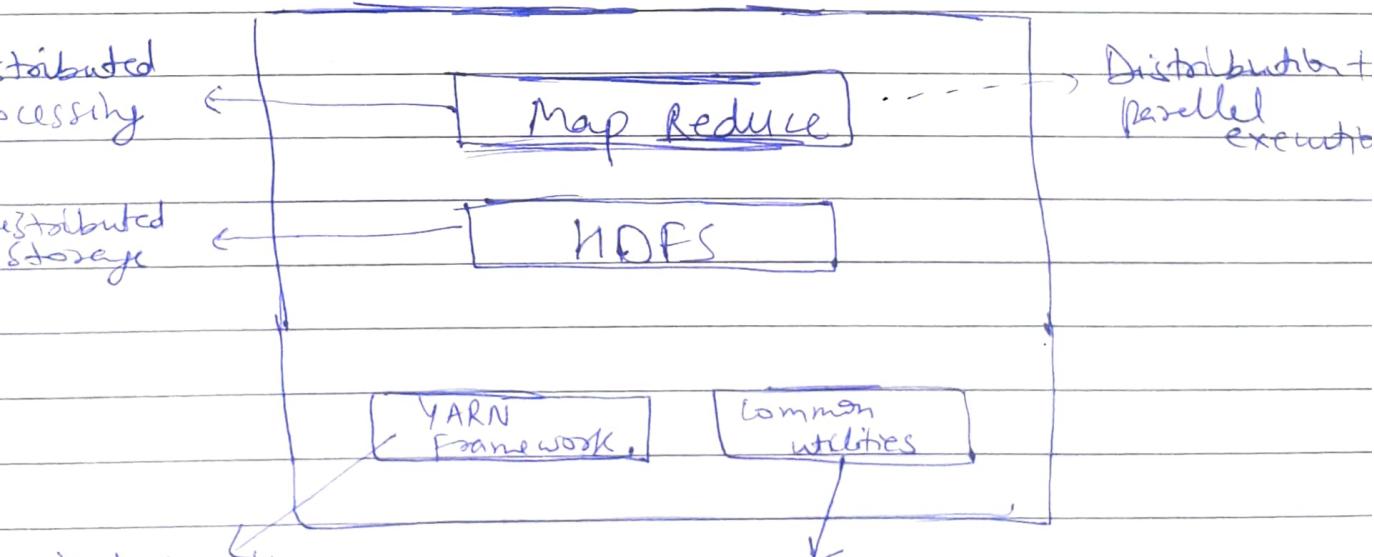
Data serialization

It is a process to format structured data in such away that it can be reconverted back to the original form. It is done to translate data structures into a stream of data.

Avro File-based data structures

Avro is a row-based storage format for Hadoop which is widely used as a serialization platform. Avro stores the data definition (schema) in JSON format making it easy to read and interpret by any program. The data itself is stored in binary format making it compact & efficient.

Hadoop Architecture



Yet Another
Resource Negotiator

(Job scheduling
& Resource Management)

divide task into jobs

in distributed system.

(Scalability, compatibility,
cluster utilization)

Java lib

& utilities (Java files &
scripts)

(are the files needed by all
other components)

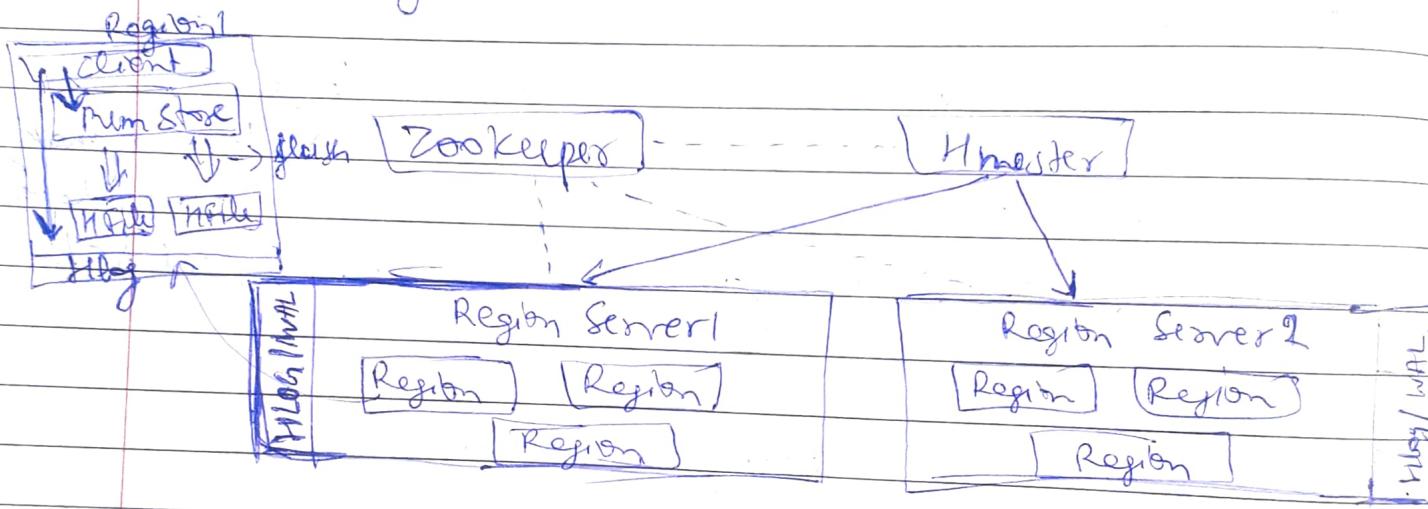
SQL → Rigid schema, consistency, transaction
NoSQL → speed, flexibility, scalability

DOMS	Page No.
Date / /	

Hadoop ke upper layer
Kafka ke 2
Banega +

Hbase

- NoSQL data store build on top of HDFS.
- column-Oriented data store.
- Store large amount of data (TB, PB)
- Big Data with random read and writes.
- handle various types of data.
- Horizontally scalable



Hmaster - to monitor Region Server instances present in the cluster.

Region Server - receives writes and read requests to a specific region, where the actual column family resides.

Region - consists of the distribution of tables and are comprises of column families.

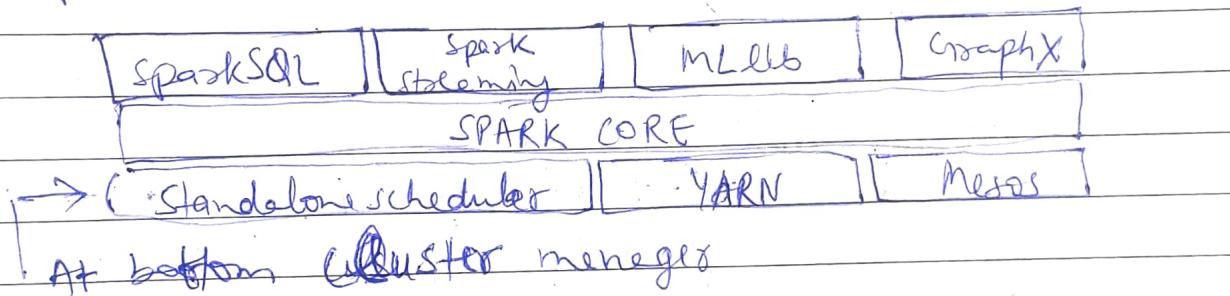
Zookeeper - centralized monitoring server which maintains configuration information & provide distributed synchronization.

Apache Spark (Cluster computing platform)

- supports data analysis, ML, graph, streaming data etc.
- can read/write from a range of data types and allows development in multiple languages.
- integrates closely with other Big Data tools
- memory computation.

Basic knowledge of Python, scala, SQL, Linux, Hadoop, statistics

Components



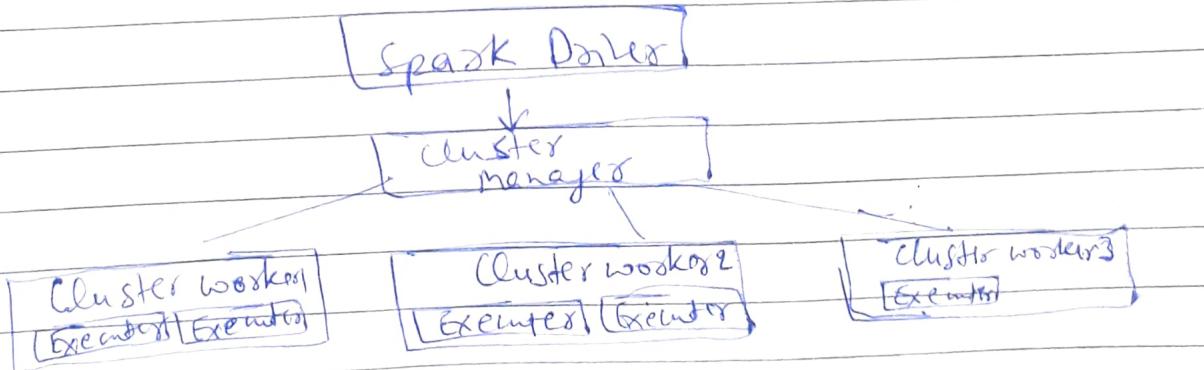
SPARK CORE - basic functionality of spark (task scheduling, MM, fault recovery, interacting with storage systems)

Spark SQL - spark's package for working with structured data.

Spark Streaming - Enable processing of live streams of data

ML lib - provides multiple types of ML algorithms.

GraphX - Library for manipulating graphs.



SPARK ARCHITECTURE

Basics

- RDD - Resilient Distributed Dataset
- DAG - Directed Acyclic Graph

RDD

- Fault tolerant
- Immutability
- Lazy Evaluation
- Caching
- Partitioning

Transformations:

```

    var a = sc.textFile("hdfs://.../abc.txt")  

    || RDD1
    var b = a.filter(...)  

    var c = b.distinct(...)  

    || RDD2
  
```

Actions

collect()



Not take action till this
action on cluster

Cassandra

It is a highly scalable, high-performance distributed database designed to handle large amount of data across many commodity servers, provide high availability with no single point of failure. It is a type of NoSQL database.

- column-oriented database
- created at Facebook

Features

- i) Elastic scalability - add more H/w to accommodate more data as per requirement.
- ii) No single point failure
- iii) Linear-scale performance - though put increases as you add more node, thus increase quick response time.
- iv) Flexible data storage - accommodates all possible data formats include structured, semi-structured and unstructured
- v) Supports ACID property
- vi)

HIVE

- Data warehousing package built on top of hadoop.
- Used for Data Analysis.
- created for users comfortable with SQL.
- Query Language - HQL or HiveQL
- Used for managing and querying Structured Data
- No Need to Learn Java.
- Uses Map-Reduce for execution
- Developed at facebook
- Not designed for OLTP and Does not offer real-time queries

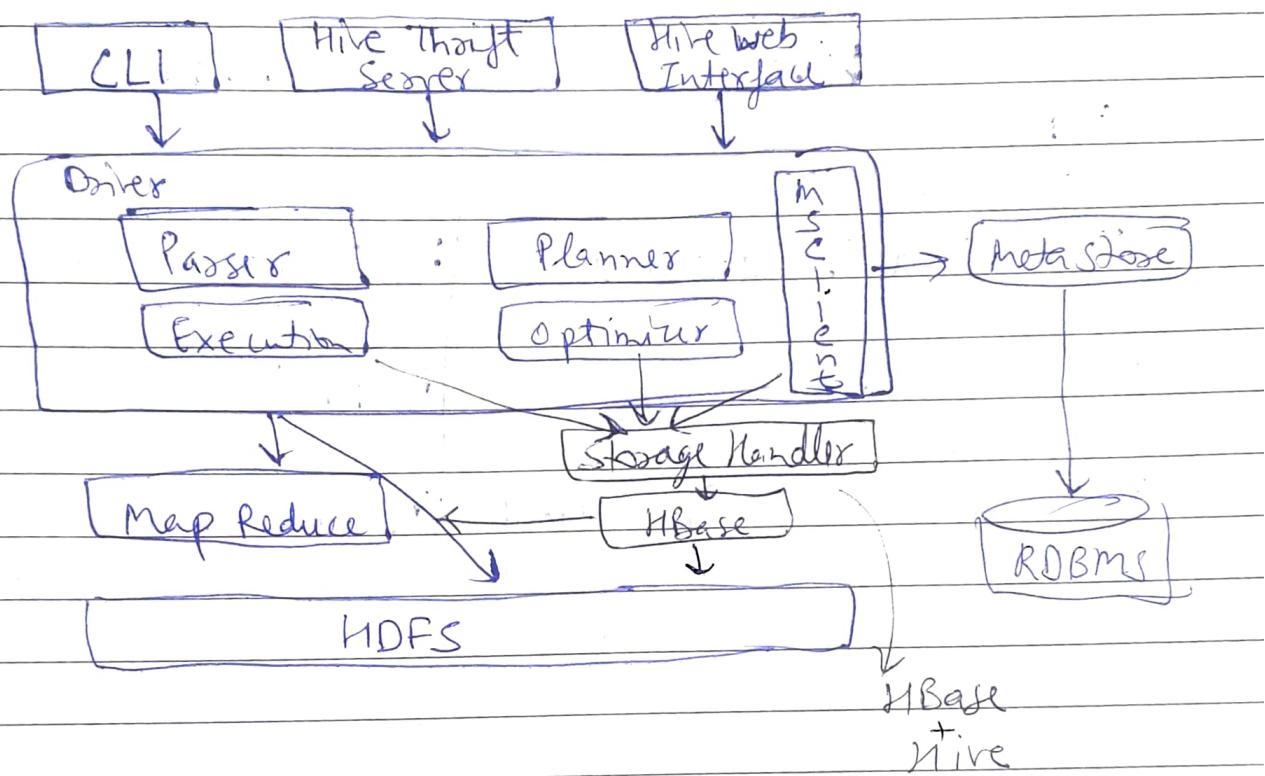
HIVE Features

- Creating Tables (Partitioning and Bucketing) \Rightarrow speed
- JDBC / ODBC Drivers are available
- Data is stored on HDFS
- Uses hadoop for fault tolerance . -

Where to use Hive

- Data Mining - to observes the pattern
- Document Indexing .
- Predictive modelling - to predict the user behaviour
- Custom facing UI - to see in different way
- Spam Detection
- Ad Optimization

Hive Architecture



Integrating with HBase

- Now single Hive query can now perform complex operation such as join, union.
- A lot of data sitting in HBase due to its usage in a real-time environment, but never used for analysis.

Compare Hive & RDBMS

Hive

- i) Focus on analytics and big data
- ii) Limited transaction support
- iii) Distributed processing via Map Reduce
- iv) Limited indexing support
- v) No triggers

RDBMS

- i) focus on online or analysis.
- ii) Transaction are usually supported
- iii) Distributed processing varies by vendor.
- iv) Full indexing support
- v) trigger support.

Difference b/w HIVE, PIG, SPL

PIG

- Developed at Yahoo
- Apache pig raises the level of abstraction for processing large datasets.
- Pig is used to analyze larger sets of data by representing them as dataflows.
- Works on top of the Hadoop

Why PIG?

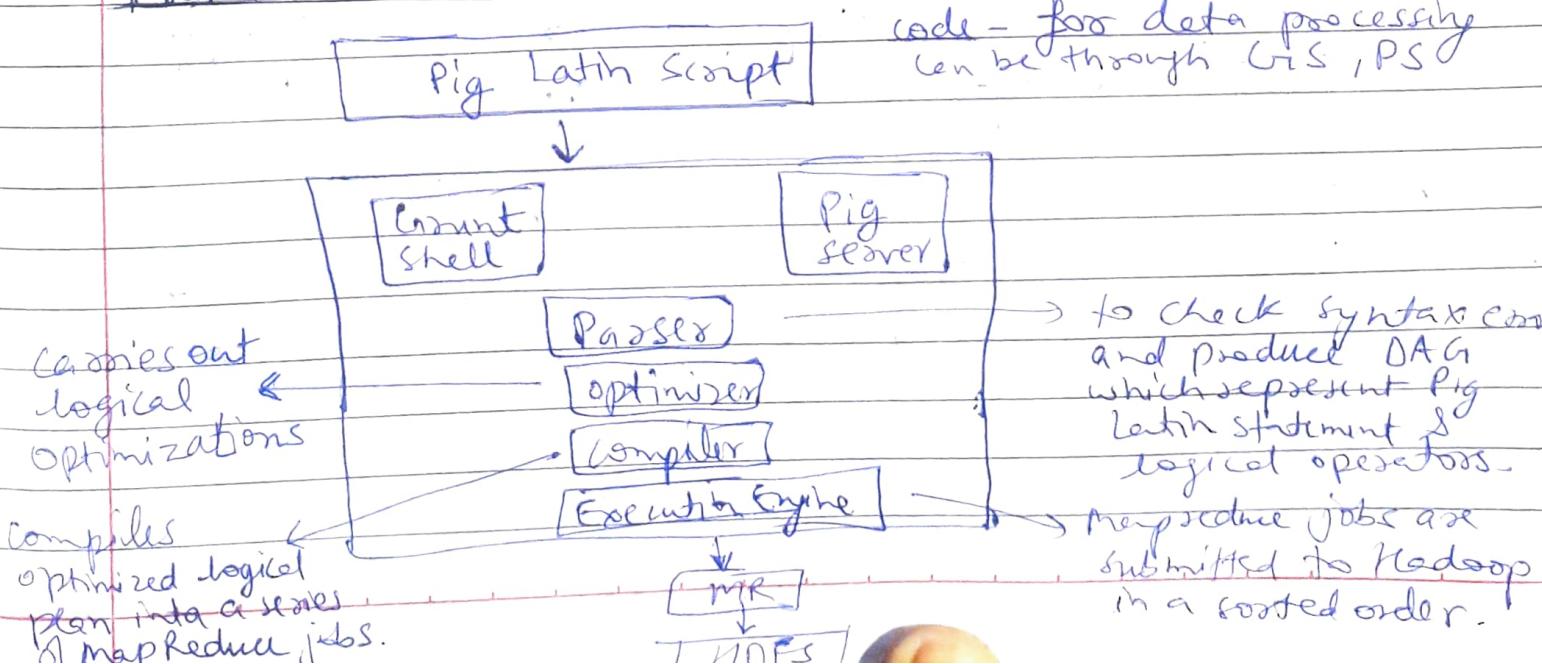
- Programmer who are not so good with Java.
- Pig uses multi-query approach.
- Built in operators :- join, filter, order
- Provides nested datatypes.

→ Pig is made up of two pieces

1. Pig Latin - Used to express dataflow

2. Execution Environment - To run Pig Latin programs

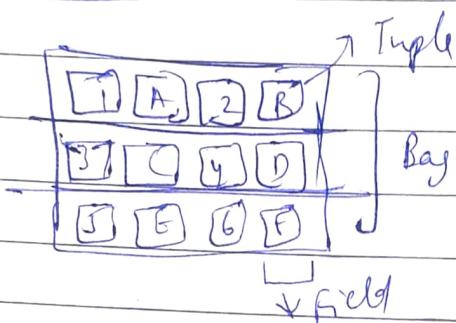
PIG Architecture



Features of Pig

- Rich set of operators
- Easy to program
- UDFs
- Optimization
- Ability to handle all kind of data (eg - S,

Data Model



Apache Pig Vs MapReduce

Apache Pig

- It is a Data Flow language
- HLL
- Join operation is simple
- basic knowledge
- use multi-query approach
⇒ reducing the code size
- no need of compilation

Map Reduce

- It is a data proc paradigm
- LLL and rigid
- quite difficult
- Exposure with Java is must to work with
- 20-times more n of lines of code.
- requires long comp

Pig

- procedural language
- schema is optional
- data model is nested relational
- submit -

SQL

- declarative language
- schema is mandatory
- data model is flat relational
- more opportunity for query optimisation.

Pig

language • Pig Latin
wid

originally • Yahoo
created

- data flow language
- procedural language
- handle S, US, semi-s

Hive

- HiveQL
- Facebook
- query processing language
- declarative language
- handle structured data

Application

- To perform data processing for search platforms
- To process time sensitive data loads.

Time Series

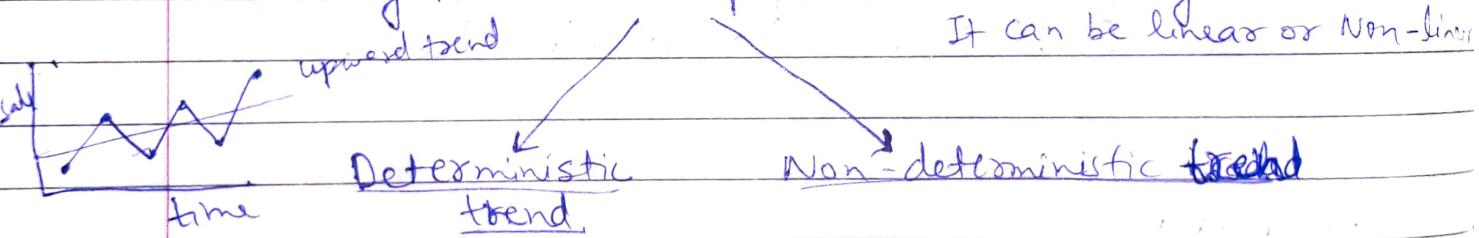
The set of data collected on the basis of time is called time series.

Importance

- Study Past Behaviour of data
- Forecast Future
- Estimates Trade Cycle
- Comparison
- Universal Utility

Components of Time Series

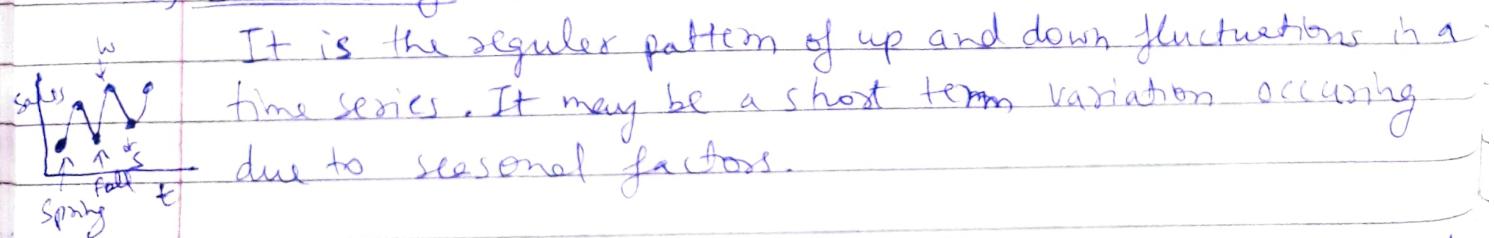
a) Trend - is a long-term increase or decrease in the series of time which persists over a long time.



Trend in which the value of the time-dependent variables increases or decreases inconsistently.

Trends in which the value of the time-dependent variable increases or decreases inconsistently.

b) Seasonality -



It is the regular pattern of up and down fluctuations in a time series. It may be a short term variation occurring due to seasonal factors.

c) Cyclicity - It can be defined as a variation caused by circumstances that repeat in irregular intervals.

d) Irregularity - variation occurs due to unpredictable factors and also do not repeat in particular patterns.

Time series Model

Non - seasonal

Model with Trend

Classical Decomposition Model

with Trend

$$X_t = m_t + Y_t$$

↑ ↑ ↑
 stochastic process trend random noise

$$X_t = m_t + s_t + Y_t$$

↑
 seasonal component

Basic methods to estimate by eliminating trend -

1. Trend estimation

- by Regression Analysis
- by Linear Trend Forecasting

Methods

1. Filtering
2. Differences
3. Joint-fit method

2. Trend elimination by differencing

Method of Measuring Trends

1. Freehand Curve Method
2. Semi - Average Method
3. Moving Average Method
4. Least Squares method

Data Science

Theories and techniques from many fields and disciplines are used to investigate and analyze a large amount of data to help decision makers in many industries.

Data Science

- It is a field that refers to the collective processes, theories, concepts, tools and technologies that enable the review, analysis and extraction of valuable knowledge and information from raw data.

Python, SAS, SQL

Tools &
Languages

Applications
Digital advertisements,
Recommender system

Data Analytics

- Data Analytics is the process of examining data sets in order to draw conclusions about the information they contain, increasingly with the aid of specialized systems & S/W.
- R, Tableau Public, Apache Spark

• Health care, Travel

Scientific Modelling

- Physics-based model
- Problem-structured
- Mostly deterministic, precise

Data-Driven Approach

- General interest engine replaces model.
- Structure not related to problem
- Statistical models handle true randomness, and un-modeled complexity.

Machine Learning

- Develop new (individual) models
- Prove mathematical properties of models
- Publish a paper

Data Science

- Explore many models, build and tune hybrids.
- Understand Empirical properties of model
- Take action.

Data Science

Approach	• Scientific (Exploration)
Problems	• Unbounded
Paths	• Iterative, exploratory, non-solution linear
Presentation skills	• Important
Research Experience	• Important
Programmatic skills	• Not as important
Data Skill	• Important

Data Engineering

Engineering (Development)
• Bounded
• Mostly linear
• Non Important
• Non as important
• Important
• Important