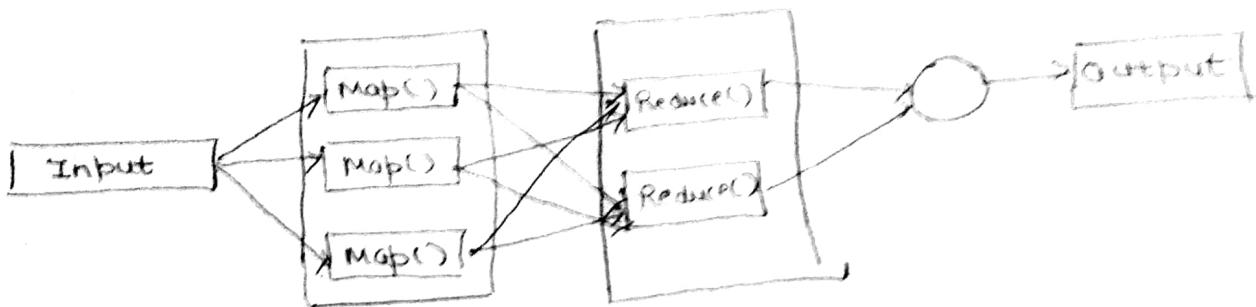


Map Reduce

It is a programming framework that allows us to perform distributed and parallel processing on large data sets in a distributed environment.



1. Two phases - Map & Reduce. Reducer phase takes place after mapper phase.
2. Map Job - Block of data is read and processed to produce key-value pairs.
3. O/P of Mapper (key-value pair) is IF to adder.
4. Reducer receives key-value pair from multiple mappers (shuffle stage) and aggregates those key-value pairs into smaller lot of tuples (key-value pairs) which is final O/P.

Advantages of Map-Reduce

1. Parallel Processing - We are dividing the job among multiple nodes and each node works simultaneously.
Data is processed by multiple machines instead of a single one and hence time taken to process data gets reduced.
2. Data Locality - Instead of moving data to processing unit, we are moving processing unit to data.
Data is distributed among multiple nodes & all the nodes work with their part of data in parallel. No chance of a node getting overburdened.

Terminology Related with Map-Reduce

- Payload - Applications implementing the Map-Reduce function.
- 2. Mapper - maps the I/P key-value pairs to a set of intermediate key-value pair.
- 3. NameNode - Node that manages the HDFS.
- 4. DataNode - Data is presented in advance before any processing takes place.
- 5. MasterNode - Node where JobTracker runs and accept job request from clients.
- 6. SlaveNode - Node where Map-Reduce program runs.
- 7. JobTracker - schedules jobs and tracks the assigned job to Task Tracker.
- 8. TaskTracker - ~~sched~~ Tracks the task & reports status to Job Tracker.
- 9. Job - A program is an execution of a mapper & reducer across a dataset.
- 10. Task - An execution of a Mapper / Reducer on slice of data.
- 11. Task Attempt - A particular instance of an attempt to execute a task on Slave Node.

HADOOP

- open source platform
- provides analytical technologies + computational boost to work with large volume of data
- Apache open source mechanism written in JAVA.
- allows distributed processing of large dataset
- client - server model.
 - Slaves - carry out computational task.
 - Master → responsible for data-distribution among clients.

HADOOP common - Java libraries & utilities required by other components
necessary files to start Hadoop.

DFS - cluster of storage solutions

Hadoop Yarn - platform for job scheduling & resource management

Hadoop Map-Reduce - perform mathematical computations

Hadoop Map-Reduce

- software platform - easily writing of applications which process big amounts of data
- limits amount of communication b/w processes as each individual record is processed by a task in isolation

Single Master Job Tracker

- resource management
- scheduling job component-task on slaves.

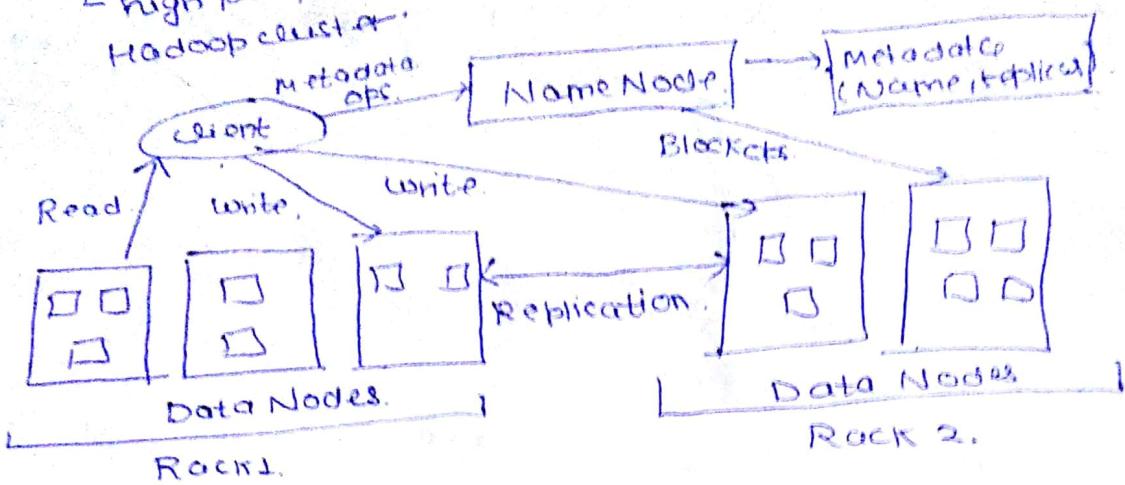
Slave Tracker

- executes task as directed by master
- provides task status info to master

Hadoop Distributed File System (HDFS)

• Hadoop Distributed File System.

↳ primary data storage system.
↳ employs a NameNode and DataNode architecture
↳ high performance access to data across highly scalable
Hadoop cluster.



• It follows master-slave architecture.

• NameNode - commodity H/w that contains GNU/Linux OS & NN software.
↳ system having NN acts as master server,

→ Manage file system namespace.

→ Regulate client's access to files.

→ Execute file system operations - renaming, closing & opening files etc.

• DataNode - commodity H/w having GNU/Linux OS & DN software.

↳ For every node in cluster, there will be data node.

→ Perform read-write operations on file systems.

→ Operations → block creation, deletion, replication

according to instructions in NameNode.

Goals of HDFS

Fault Detection & Recovery

HDFS have large no. of commodity HW, failure is frequent
↳ should have mechanism for quick & automatic fault detection & recovery.

2. Huge Datasets - HDFS should have 100s of nodes per cluster to manage large datasets.

3. Highly Fault-Tolerant - file system replicates each piece of data multiple times & distributes copies to each node. As a result, data on nodes that crash can be found elsewhere within a cluster.

Analytic Process - technique for analyzing & organizing complex decisions, based on mathematics & psychology

Business Understanding [Problems to be solved / Decisions to be made.
Once the objectives are identified & determined
└ identify the goals.

Data collection - Data is collected from different sources.
(structured, semi-structured / unstructured).

Data Preparation - Remove unnecessary / unwanted data.

Data Modelling - A model is created to analyze relationship b/w different objects of data.

Data Evaluation [Results are obtained from different test-cases.
Data is evaluated & rechecked.

Deployment - constantly checked for errors & data is maintained.

Characteristics of Big Data

1. Programmatic - Analysts use code to handle raw data.
└ to manipulate & explore it.
2. Data Driven → BD uses high volume of data to drive analysis.
└ critical characteristic for big data analytics.
3. Iterative → iterations on models provides more computational powers.
4. Use Many Attributes → used for describing relationship b/w entities.

• Categories of Big Data Analytics.

1. Prescriptive - reveals what actions should be taken. (most valuable)
[results in rules / recommendation for next steps.]
2. Predictive - analysis of scenario what might happen - predictive forecast
3. Diagnostic - looks at past performance & determine what happened/why.
[Result - an analytic dashboard.]
4. Descriptive - What is happening now based on incoming data.
[Use real-time dashboard or email reports.]

• Types of Big Data Analytics.

1. Basic Analysis - explore data where value is unknown (data may be useful or not)
Only used for small amount of data.
[Used when lot of disparate data need to be analyzed.
Visualizations of simple statistics.]

Slicing & Dicing - breaking down your data into smaller sets - easier to explore.

Basic Monitoring - monitor large volume of data in real time.

Anomaly Identification - identify anomalies - an event where observations differ from what's expected
[gives a clue - that something is going wrong]

2. Advanced Analytics - complex analysis of structured/unstructured data
[Map Reduce]
[sophisticated statistical models - ML, Neural Networks
Text Analysis are used]

Predictive modeling - statistical / data mining & ML (containing algos.)
↓
to determine future outcome.

Text Analytics - analyze unstructured text

extract relevant info

transform to structured info that can be leveraged in many ways.

Used computational linguistics / statistics etc.

3 | Operational Analytics } - Type of Business analytics.

- | Focuses on improving existing operations.
- | have different software & packages.
- | which offers various models for showing what happens within business.

4 | Monetizing Analytics } - can be used to optimize your business for better decision making & thus increasing revenues.

Eg Credit card providers use their assembled data to provide value-added services.

Able to assemble unique data set that is valuable to company.

Text Analysis - Analyzing unstructured data, extracting relevant info and transforming it into structured info that can be leveraged in various ways.

Deriving info from text sources.
unstructure & semi-structured.
↓
web logs, blogs, email, social media.

Statistical Models + Linguistic Theories (NLP+ML+Stats).
are used to capture patterns in human language such that machines can understand text & perform analytic task.

Entity extraction - identifies an item/person/geographical location/contact info/date/time/currencies (named entities).

E.g X is identified as a person in a text to analyze.

Facts extraction (Relationship) Relationship b/w two entities.
statement about something that exist / happens
E.g X is the CEO of company Y.

Events - usually contain a time dimension - causes facts to change
Eg- Status of a sales process.

Concepts - identify an event/process/trend/behaviour.
set of words/phrases → indicating a particular idea/topic to which a customer

Sentiments: identify viewpoints / emotions in the underlying text
Techniques used - ML / NLP

Social Media Analytics

Gathering data across the internet - from various blogs / social media, news articles / online text forums etc.

This huge stream of data - is analyzed (often using text analytical) to make better business decisions.
→ mine customer sentiments to support marketing & customer service activities.

- Used by health agencies - to identify public health threats.
- used by govt - to look terrorist activities.

Level 1 : Internal / within Enterprise - does not understand usefulness of social media.

Level 2 : Monitoring - minimal engagement ~ no intention to address any issue being discussed.

Level 3 : Broadcast - interaction with L2 - social media - broadcasts marketing message.

Level 4 : Viral - limited engagement

Level 5 : campaign - engages with community - run contests.

Level 6 : Collaborative - begins listening / analyzing / interacting / provides feedback.

Level 7 : People Powered ; Top level engagement

monitor conversation around their profile.
use SM - to actively locate & address problems

Text Analytics Tools for Big Data

Attensity

- original text analytics company - began developing & selling products 10 years ago
- world's largest NLP development group
- provides several engines for text analytics
 - Auto-classification
 - Entity-extraction
 - Exhaustive Extraction:
 - Attensity's flagship technology that automatically extracts facts from parsed text & organize this info

- Company is focused on social and multichannel analytics, and tailoring it to business users for engagement.
- Developed a grid computing system - provides high performance capabilities for processing massive amount of real-time data
- Uses Hadoop framework (MapReduce, HDFS, HBase) to store

Clatabridge

- pure-play text analytics vendor (deals with unstructured data)
- goal - help companies to take better business decisions based on customers' experiences & issues
 - Helping everyone in organisation to take actions & collaborate in real time.

- Includes
 - real time sentiment determination
 - classification of customer feedback data / text
 - staging verbatim for future processing into clatabridge

Features - Single-click root cause analysis.

- identify what is causing a change in volume of text feeds, sentiment etc.

Open Text

Canadian-based company.
best known for - its leadership in enterprise information management (EIM).

Vision - managing, securing, extracting value from unstructured data of enterprises.

Semantic Middleware

- Enables real-time analytics with high accuracy on large datasets across languages, formats and industry domains.
- Semantics can be exposed at different levels and work with different technologies to address business issues.

Streaming Data

Streaming data [analytic computing platform focused on speed. is useful when analytics [these application requires need to be done in real time - when data is in motion] continuous stream of unstructured data. → Data is continuously analyzed and transformed in memory before it is stored on disk.

→ Processing streams of data - by processing time windows of the data in a memory across several clusters of servers

Streaming data can be used in following scenarios.

- collecting info about movement around a secure site.
- To be able to react to an event - that needs immediate response. (patient's medical condition)
- Real time ~~cost~~ calculations of cost that depend on various variables - usage & available resources

Complex Event Processing (CEP)

(4)

→ Streams & CEP both are used to manage data in motion but the uses of both are quite different.

→ CEP is a technique for tracking, analyzing & processing data as event happens.

Info is then processed & communicated based on business rules & processes.

Idea — To be able to establish a correlation b/w stream of information & match the resulting pattern with defined behaviors.

→ It is advanced approach based on simple event processing.

that collects & combines data from different relevant resources to discover events & patterns. that can result in action

E.g. credit card companies uses CEP for better fraud management when a pattern of fraud emerges - Company shuts down credit whether reporting applications / sales management application

Difference b/w CEP and streams.

Streaming Data

Applied to analyze vast amount of data in real time.

Streaming data applications manage a lot of data & process it at a high rate of speed.

Managed in highly distributed clustered environment.

CEP

Focused on solving specific use cases based on events & actions.

Does not manage as much data. Runs on less complex hardware.

Operationalizing Big Data

Making Big data a part of your business process.

→ planning an integration strategy.

Data (traditional / big data) - needs to be integrated
as a seamless part:

Integrating Big Data

- |- understand the problem - - - - -
 - |- Identify the process involved - - - - -
 - |- Identify the info required to solve problem - - - - -
 - |- Gather data, process it & analyze the results - - - - -
- Patient eg.
- 1. Need to treat patient
 - 2. Diagnosis & Test
 - 3. Results of test
Past history
 - 4. Start treatment
monitor patient

Applying Big Data Within Your Organization

(Not complete).

1. Figuring the Economics of Big Data

To look at various methods for putting big data to work for u/f org
cost may vary due to size of org, purchasing power, vendor relations

B-D economics should be analyzed in following areas.

- |- Identification of data type & sources
- |- Technology changes / New technologies for bB
- |- New talent acquisition / upgrades to existing.

2. Identification of data type & sources

Security & Governance For Big Data Environments

(5)

- If collecting data from unstructured data sources, social media sites.
└ make sure that viruses/bogus links are not buried in content

Data Protection Options

1. Data anonymization
 - └ remove all data that can be uniquely tied to an individual (person's name / SNo / credit card No.)
 - └ can protect personal identification
2. Tokenization
 - └ protects sensitive data by replacing with random tokens / alias values that mean nothing to someone who gets unauthorized access to this data.
 - └ can protect credit card info, passwords, personal info
3. Cloud Database Control
 - └ access controls are built into the db.
 - └ protects within db so that each block of data doesn't need to be encrypted.

Data Governance

- ① visibility
 - └ You may not have control how your data is used and controlled.
 - └ No control over your visibility into your resources that are running outside your control.
- ② Unvetted Employees
 - └ company should go through extensive background check of all of its employees.
 - └ trust cloud-provider.

③ Auditing your big data process

- ↳ claim that you are meeting the rules necessary to support the operations of business.
- ↳ Need to show logs / evidence - that data you are using is secure & clean.
- ↳ Need to explain the sources of data.

④ Putting the Right organizational structure in place

⑤ Preparing for management of risk

⑥ Setting the right governance & quality policies

⑦ Developing a well governed & secure Big Data Environment