

- Q.1 characteristics of big data (5V's)
- 1) Volume:- The large amount of data generated, measured in Petabytes or exabytes.
- 2) Velocity:- The speed at which data is generated & processed in real time.
- 3) Variety:- The different types of data
- ① Structured - SQL DB
 ② Unstructured - Predefined, text
 ③ semi-structured - JSON, XML files
- 4) Veracity - The accuracy & trustworthiness of data.
- 5) Value - The meaningful insight derived from process data.

- Q.2 Q.3. Features of HDFS (Hadoop distributed file system)
- ① Scalability - Easily scalable by adding more nodes to the cluster.
- ② Fault tolerance - ensure reliability.
- ③ High throughput - optimize for large datasets enabling high data access.
- ④ Streaming data access:- Designed batch process allowing streaming access.
- ⑤ Large Data storage - Capable of storing very large files across multiple machines.

- Q.4. Challenges of Big data
- ① Storage ④ Transfer
② Capture ⑤ visualization
③ Analysis ⑥ privacy

Q. 5 Applications of Big data Analytics

→ ① Smarter health care

② Manufacturing

③ Search quality.

④ Trading Analytics

⑤ Traffic control

⑥ Telecom

⑦ Homeland security

Q. 6 Hadoop - i) Hadoop is an open source framework designed for distributed storage & distributed processing of large datasets using a cluster of computers.

ii) It enables scalable, fault tolerant & efficient data processing through its Core Components:

- ① HDFS
- ② map-reduce
- ③ YARN
- ④ common Hadoop.

• Why it matters:-

features

① Scalability

② Cost effective

③ Flexibility

④ Reliability

⑤ Wide Adoption

Q. 7 Advantages & limitation of hadoop.

→ ① Cost reduction

② Flexible

③ Enhanced scalability

④ minimal traffic

- Limitations of hadoop -

- ① Complex programming model - Java
- ② Not suitable for Real time processing - Batch processing design.
- ③ Small file problem:- inefficient at handling small files. due to excessive metadata storage in namenode.
- ④ Latency issues - High latency due to batch oriented nature of mapreduce, making it less suitable for time sensitive data processing.
- ⑤ Resource intensive - require significant hardware resources for setup & maintenance.

8# Hadoop Ecosystem & its Components.

Data management	oozie monitoring	chukwa monitoring	flume monitoring	zookeeper management	
Data Access	Hive SQL	pig dataflow	mahout ML	Avro RPC	Sqoop RDBMS connect
Data Processing	MAP Reduce cluster management		YARN Elastic resource manager		
Data storage	HDFS filesystem		HBASE column DB structure		

Components of Hadoop

- i) Hadoop Ecosystem holds variety of tools & frameworks to support big data.
- ① Data management - zookeeper
 - ② oozie - job scheduling & monitoring
 - ③ chickwa & flume - monitoring
 - ④ Data processing - pig, Hive - query SQL, dataflow Access, sqoop RDBMS connector.
 - ⑤ mahout, spark - ML
 - ⑥ Data processing - mapreduce - programming by YARN - yet another resource negotiator
 - ⑦ HDFS - Hadoop Distributed File System
 - ⑧ HBASE - noSQL Database
 - ⑨ spark - in memory data processing

Example - Imagine Company Analyzing Customer behaviour data.

- The data is stored in HDFS
- Data processed using mapreduce
- Data managed by YARN
- Supported by hadoop common libraries.

Hadoop Using Companies:

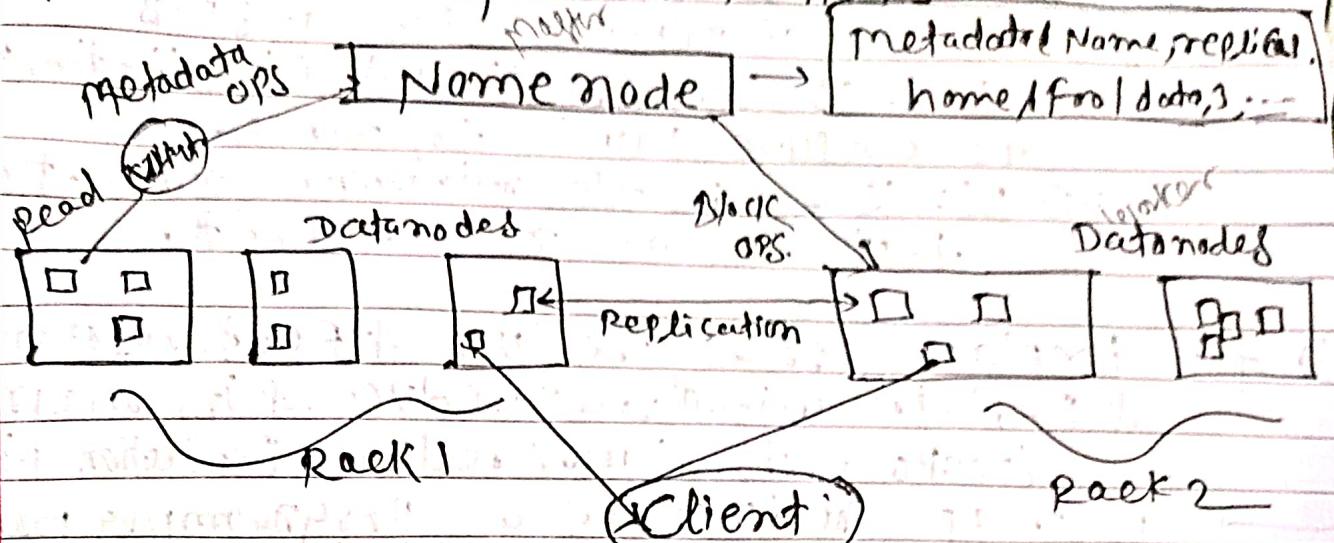
- ① Facebook
- ② yahoo
- ③ Twitter
- ④ ebay
- ⑤ AWS
- ⑥ netflix
- ⑦ Spotify
- ⑧ BANK OF America

- ① Name node
- ② Data nodes
- ③ Secondary data nodes

④ Block storage

Data
Page

g) Architectures of Hadoop HDFS.



- ① Name node - i) Manage the metadata of the file system [Master node]
 ii) Manager & Controller of HDFS
 iii) keep track of the location of data blocks across data nodes.
 iv) manage file system namespace & control access to files.

- ② Data nodes - i) [Worker node] that store actual data blocks.
 ii) Perform read & write operations as instructed by the name node.
 iii) periodically send heartbeat signal to the name node to confirm their status.

- ③ Secondary name node - i) It is separate physical machine which act as [helper of namenode]
 ii) It perform periodic check points. It communicate with namenode and take snapshot of meta data which help to minimize the down time.

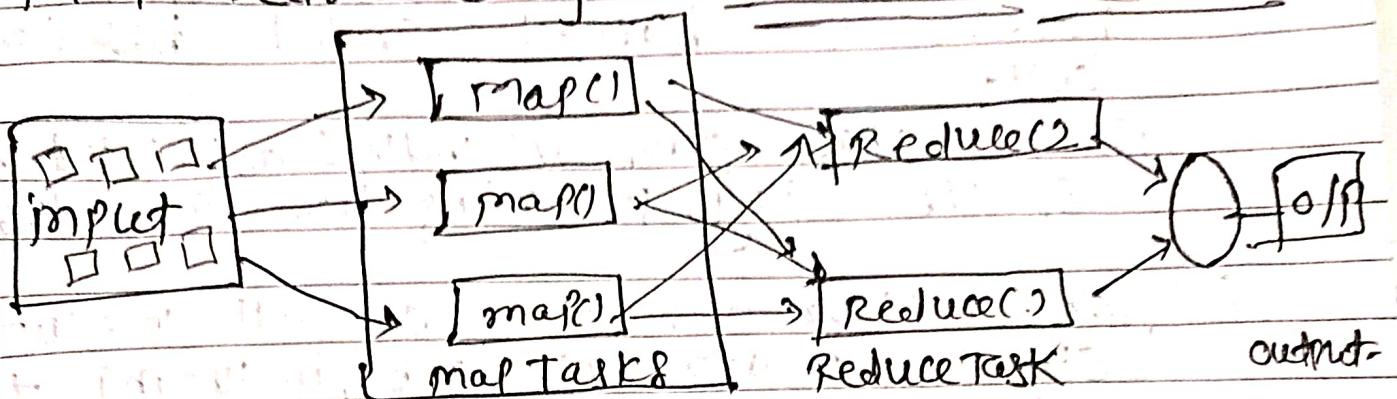
- Q) Block storage - i) Data in HDFS is stored in blocks
 default size 128MB
 ii) Blocks are replicated (replication factor is 3)
 to ensure fault tolerance.
 iii) Block - is the minimum amount of data
 that it can read or write.

10. # Explain map reduce & features of map reduce
 & Explain with diagram reduce & its main components & map reduce execution pipeline.
- i) map reduce is a programming framework that allows us to perform distributed & parallel processing on large data sets in a distributed environment.
- ii) Data processing - Breaking down tasks into two main functions: ① map ② reduce.
- iii) map - [Sorting & filtering the data]
 - organizing them into form of group - key-value pair based result which is later on processed by the reduce() method.
- iv) Reduce - (Summarization by aggregating the mapped data.)
 reduce - takes o/p of the map as input, reduce & combine them into smaller tuples.

Features of map reduce

- ① Scalability
- ② Fault tolerance - Reliability.
- ③ Parallel Processing
- ④ Data Locality
- ⑤ Simplicity,
- ⑥ Extensibility

Map reduce Diagram & Execution Pipeline



Components

- ① Input Splits
- ② map function
- ③ Reduce function
- ④ output format
- ⑤ Shuffle & Sort
- ⑥ Job & Task Tracker

① Input Splits - i/p data is split into small chunks called i/p splits. Each split processed independently.
i) splitting helps in parallel processing by distributing the splits across different nodes.

② map function - map function process each i/p split & generate a set key-value pairs.
ii) It perform data filtering, sorting & X'er based on the logic defined in the map function.

③ Reduce function - Takes the grouped key-value pairs & perform aggregation or summarization.
ii) o/p is final set key value Pairs, which represent the result of the map reduce job.

④ output format - HDFS, o/p format defined how the key-value pairs are stored & represented.

⑤ Shuffle & Sort - After map phase key value pairs are shuffled & sorted based on the key.

iii This phase groups all values associated with specific key, preparing them to reduce phase.

⑥ Job Tracker - manage overall execution of the mapreduce job; including task assignment, monitoring & job completion.

⑦ Task Tracker - Runs on each node in the cluster & responsible for executing the task assigned by the job tracker.

②
NoSQL

① TYPES OF NoSQL DB.

- ① Document Store - json, XML, BSON, eg. MongoDB
- ② Key-value store - key-value pair, eg. Redis, DynamoDB
- ③ Column family store - eg. Apache Cassandra, HBase
- ④ Graph Database - nodes, edge, eg. Neo4j, Amazon Neptune

② Features of NoSQL DB.

- ① Schema flexibility - dynamic schema
- ② Scalability
- ③ High performance
- ④ Distributed Architecture
- ⑤ High Availability
- ⑥ Support for unstructured data
- ⑦ Eventual Consistency
- ⑧ Fault tolerance

- Q. 3 Advantages of NoSQL
- ① Scalability
 - ② flexibility
 - ③ Availability
 - ④ performance
 - ⑤ cost effective
 - ⑥ Suitability for Big data

- Disadvantages of NoSQL
- ① Complexity
 - ② Lack of standardization
 - ③ Limited Acid support
 - ④ Less mature tools
 - ⑤ Eventual consistency

Q. 4

SQL

- ① Relational DB (Rows & Column Tables)
- ii) fixed schema
- iii) vertical scalability
- iv) Structured query language
- v) full support ACID properties
- vi) Example - MySQL, PostgreSQL

NoSQL

- i) non-Relational (document key-value, graph, column)
- ii) dynamic schema
- iii) horizontal scalability
- iv) JSON query
- v) limited support ACID
- vi) MongoDB, Redis, Neo4j

Q. 5

DBMS

- i) Database management system
- ii) handle structured data
- iii) store static data in tables
- iv) support SQL
- v) Batch processing higher latency
- vi) Ex - MySQL, PostgreSQL
- vii) support CRUD operation

DISMS

- i) Data Stream management system
- ii) Handle streaming data
- iii) process continuous data
- iv) support continuous queries
- v) real time processing, low latency
- vi) Apache Flink, Apache Storm
- vii) support read operation

Q.1

DGIM Algorithm [- Datar - Gionis - Indyk - Motwani Algorithm]

→ Used for Counting the number of 1's in a binary stream over a Sliding Window.

We modify the buckets after a new 1's arrives.

- It uses $O(\log_2 N)$ bits to represent a window of N bits.
- Error rate is no more than 50%.

Elements

- ① Timestamp
- ② Bucket (0 & 1)

Rules for Bucket

- ① Every bucket should contain at least a single 1 in it.
- ② Right side of the bucket should strictly start from 1.
- ③ Length of the bucket is equal to the number of 1's in it.
- ④ Every bucket length should be in power of 2.
 $2^0 = 1, 2^1 = 2, 2^2 = 4, 2^3 = 8, 2^4 = 16 \dots$
- ⑤ As we move to left, the bucket size should not decrease.
- ⑥ No more than two bucket can have same size.

For Data Stream mining.

②

Bloom Filter - is a probabilistic Data Structure used to test whether an element is part of a set.

ii) It uses multiple hash function and bit arrays, allowing space-efficient storage but with a possibility of false positives.

Application:-

Used in Stream Data mining to quickly check membership of elements without storing the entire dataset.

③

Explain issues in Data stream query processing.

→ ①

High throughput - Stream generate large amount of data rapidly, require efficient processing techniques.

②

Approximate Answers - full accuracy is often sacrificed for performance, require approximate query answers.

③

Memory Constraints :- Storing & Processing an entire data stream in memory is impractical.

④

Data skew :- Un-even data distribution in streams can impact performance and accuracy.

-2

Given - Bloom Filter

$$n = 15$$

$$k = 3$$

$$S = \{25, 36, 49\}$$

check membership of 16 & 81

- $h_1(x) = x \% 15$

- $h_2(x) = (x+2) \% 15$

- $h_3(x) = (2x+1) \% 15$

→ step ① Initialize bit array/vector

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

① S = 25 Insert it.

- $h_1(x) = 25 \% 15 = 10$

- $h_2(x) = (25+2) \% 15 = 12$

- $h_3(x) = (25 \times 2 + 1) \% 15 = 6$

$$\begin{array}{r} 25 \\ - 15 \\ \hline 10 \end{array}$$

Remainder

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	0	1	0	0	0	1	0	1	0	0	0

bit vector.

② Let insert $S = 36$

- $h_1(x) = 36 \% 15 = 6$

- $h_2(x) = (36+2) \% 15 = 8$

- $h_3(x) = (2 \times 36 + 1) \% 15 = 13$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	0	0	0	0	1	0	1	0	1	0	1	1	0	0

bit vector

③ Let insert $S = 49$

- $h_1(x) = 49 \% 15 = 4$

- $h_2(x) = (49+2) \% 15 = 6$

- $h_3(x) = (2 \times 49 + 1) \% 15 = 9$

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

Bit vector.

① Let's Lookup 16

$$\cdot h_1(x) = (16 \mod 15) = 1$$

$$\cdot h_2(x) = (16+2) \mod 15 = 3$$

$$\cdot h_3(x) = (16 \times 2 + 1) \mod 15 = 3$$

Since 1 & 3 are 0's 16 is not present in the set

② Let's Lookup 8

$$\cdot h_1(x) = 8 \mod 15 = 8$$

$$\cdot h_2(x) = 8+2 \mod 15 = 10$$

$$\cdot h_3(x) = 8 \times 2 + 1 \mod 15 = 13$$

1	2	3	4	5	6	7	8	9	10	11	12	13	14
1	1	1	1	1	1	1	1	1	1	1	1	1	1

Since 6, 8 & 13 are 1's, 8 is present in the set (False +ve)

* Applications of Bloom Filters

① Aurora

② google Chrome

③ medium

④ Big Table / Hbase / Cassandra / postgresql

P-3

$$n=12$$

$$k=3$$

Practice

$$S = \{17, 23\}$$

check membership of 13 & 29

Flajolet - Martin Algorithm - estimates the number of distinct element in a data stream using hash function & bit patterns.

Given : Stream - 4, 2, 5, 9, 1, 6, 3, 7
Hash Function

$$h(x) = ax + b \bmod 32$$

$$h(2x) = 3x + 1 \bmod 32$$

$$h(3x) = 2x + 6 \bmod 32$$

① $h(x) = 3x + 1 \bmod 32$

$$= 3(4) + 1 \bmod 32 = 12 + 1 \bmod 32 = 13 = 01101$$

$$= 3(2) + 1 \bmod 32 = 6 + 1 \bmod 32 = 7 = 00111$$

$$= 3(5) + 1 \bmod 32 = 15 + 1 \bmod 32 = 16 = 10000$$

$$= 3(9) + 1 \bmod 32 = 27 + 1 \bmod 32 = 28 = 11100$$

$$= 3(1) + 1 \bmod 32 = 3 + 1 \bmod 32 = 4 = 00100$$

$$= 3(6) + 1 \bmod 32 = 18 + 1 \bmod 32 = 19 = 10011$$

$$= 3(3) + 1 \bmod 32 = 9 + 1 \bmod 32 = 10 = 01010$$

$$= 3(7) + 1 \bmod 32 = 21 + 1 \bmod 32 = 22 = 10110$$

$$\gamma(a) = \text{Trailing zeros} = \{0, 0, 4, 2, 20, 1, 1\}$$

$$R = \max[\text{Trailing zeros}] = 4$$

$$\text{Estimated output} = N = 2^R = 2^4 = 16$$

② $h(x) = 2x + 6 \bmod 32$

$$= 4 + 6 \bmod 32 = 10 \bmod 32 = 10 = 01010$$

$$= 2 + 6 \bmod 32 = 8 \bmod 32 = 8 = 01000$$

$$= 5 + 6 \bmod 32 = 11 \bmod 32 = 11 = 01011$$

$$= 9 + 6 \bmod 32 = 15 \bmod 32 = 15 = 01111$$

$$= 1 + 6 \bmod 32 = 7 \bmod 32 = 7 = 00111$$

$$= 6 + 6 \bmod 32 = 12 \bmod 32 = 12 = 01100$$

Trailing zero's = The zero's at the right must end
Counting right side zero



$$376 \bmod 32 = 9 \bmod 32 = 9$$
$$746 \bmod 32 = 13 \bmod 32 = 13$$

1	1	6	8	4	2
0	1	0	0	1	0
0	1	1	0		

$$\gamma(a) = \text{Trailing zero's} = \{1, 6, 8, 4, 2, 0, 0\}$$
$$l = \max[\text{Trailing zero}] = 2$$
$$\text{Estimated output } N = 2^R = 2^{l+1} = 2^{2+1} = 8$$

(P2) Practice. - Stream = 1, 3, 2, 1, 2, 3, 4, 3, 1, 2, 3, 1
Hash function $h(x) = 6x + 1 \bmod 5$
Flajole Martini Algorithm.

⑤ Finding Similar Items & Clustering

① CURE Algorithm - Clustering Using Representatives
It is a hierarchical clustering algorithm that handles large datasets & outliers better than traditional methods.

Steps - ① Sample Selection - select random sample of data points.

② cluster the sample - perform an initial clustering on the sample using a traditional algorithm [e.g K-means.]

③ Select Representative Points :- for each cluster, select representative points.

④ shrink Representative points \rightarrow toward the centroid using shrinking factor.

⑤ merge clusters - based on the proximity of their representative points.

⑥ label the remaining Data - Assign the remaining Data Points to the nearest clusters.

Example — In a dataset with scattered clusters, CURE selects a few representatives points from each clusters & moves from closer to the center, effectively handling non-spherical clusters & outliers.

④ Distance measures - ① Hamming distance

② Cosine angle

③ Jaccard distance

④ Jaccard similarity

④ Euclidean Distance.

P-1). Explain with example two major classes of Distance measures!

① Euclidean distance — measures the straight line distance between two points

$$d(x, y) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2}$$

$$\text{Ex: } (1, 2) \text{ & } (4, 6)$$

$$(x_1, y_1) (x_2, y_2) = \sqrt{(4-1)^2 + (6-2)^2}$$

$$= \sqrt{9 + 16}$$

$$= \sqrt{25} = \underline{\underline{5}}$$

② Cosine angle similarity — measure cosine angle between two vectors, indicating their similarity in terms of direction

$$\text{Cosine similarity} = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \cdot \|\mathbf{B}\|}$$

$$\text{Eq. } A = (1, 2) \quad B = (2, 4)$$

$$A = (1 \times 2) \cdot (2 \times 4) = 16$$

$$\|\mathbf{A}\|^2 = 1^2 + 2^2 = \sqrt{5} \quad \|\mathbf{B}\|^2 = 2^2 + 4^2 = \sqrt{20}$$

$$= \frac{16}{\sqrt{5} \times \sqrt{20}} = \underline{\underline{1.6}}$$

$$\cos^{-1}(-0.6) = 126^\circ$$



Cosine distance = $1 - \text{Cosine similarity}$

$$= 1 - 1.6 = -0.6$$

③ 1 indicate \rightarrow vector A & B pointing in exactly same direction.

0 indicate \rightarrow they considered "close" to each other in terms of direction.

③ Jaccard distance — measure the dissimilarity between two datasets & is calculated

Jaccard distance = $1 - \text{Jaccard similarity}$.

example = JS = 80%

$$JD = 1 - 0.8 = 0.2 \text{ or } 20\%$$

④ Jaccard similarity

$$JS(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

P-1

Example = $A = \{0, 1, 2, 5, 6, 8, 9\}$

$B = \{0, 2, 3, 4, 5, 7, 9\}$

\rightarrow

$$JS(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

$$= \frac{\{0, 2, 5, 9\}}{\{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}} = \frac{4}{10} = 0.4$$

$$T(A, B) = 0.4 \approx 40\%$$

④ Hamming distance - Difference between
Corresponding bits.

e.g- ① 110011 & 010101
② 11010 & 01011

→ ① → Compare bit by bit 110011 & 010101

110011
 010101
 $\downarrow \downarrow \downarrow \downarrow \downarrow$
Position 100110

→ 1 4 5 Difference position $1, 4, 5 = 3$

Hamming distance = 3

② 11010
 01011
 $\downarrow \downarrow \downarrow \downarrow \downarrow$
Position 10001
1 5 Difference Position = $1, 5 = 2$

Hamming distance = 2

-2 cosine angles between

1. $(3, -1, 2)$ & $(-2, 3, 1)$

A · B

$$\begin{aligned} & \|A\| \cdot \|B\| \\ &= (3x1 + (-1)x2) \times (-2x3 + 3x1) = (3x-2) + (-1x3) + (2x1) \\ &= (\sqrt{(3)^2 + (-1)^2 + (2)^2}) \times (\sqrt{(-2)^2 + (3)^2 + (1)^2}) \end{aligned}$$

$$\begin{aligned} &= \cancel{-6} - 7 = -0.5 \\ &\sqrt{14} \times \sqrt{14} \end{aligned}$$

$$\cos^{-1}(-0.5) = 120^\circ$$