

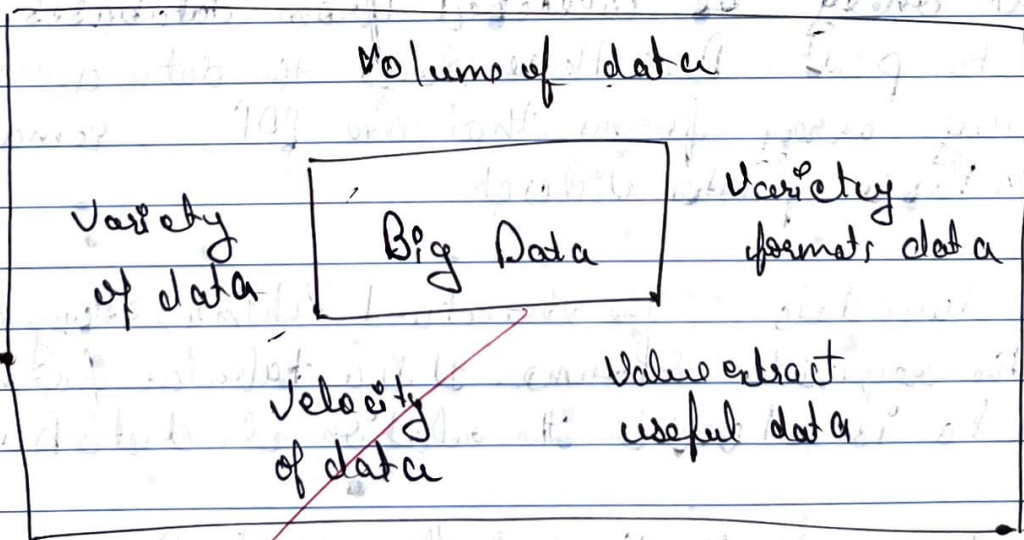
## Assignment No. 1

Q. Write on Briefly on Big data characteristics Hadoop architecture  
Hadoop ecosystem

⇒ Big data Characteristics

Big data contains a large amount of data that is not being processed by traditional data storage or the processing unit. It is used by many mathematical companies to process the data and business of many organizations. The data flow could be exceed 150 PB (petabytes) per bytes before replication.

There are five Vs of Big Data that explain the characteristics.



The 5 Vs of Big data are as follows:



1] Volume

The name of Big data itself is related to an enormous size. Big data is vast volume of data generated from many sources daily such as business process, machines, social media platform, networks human enterprises and many more.

eg :- Facebook can generate approximately a billion message 4.5 B. This that the 'like' button is touched & more than 350 million new post uploaded each day.

2] Variety

Big data can be structured, unstructured & semi-structured. Data are being collected from different sources. Data will only be collected from databases & sheets in the past. But these days the data comes in the array array form that are PDF's, emails, audios, SM Passes, photos, videos etc.

1] structured data :- The structured schema, along with all the required columns. It is a tabular form. Structured data is stored in the relational database.

2] semi structured :- The semi structured the schema is not appropriately defined eg: JSON, XML, S-  
OTP (online transaction processing) systems are built to work with semi-structured data. It is stored in real time i.e. table.



Unstructured data: All the unstructured files, log files, audio files, & image files are included in the unstructured data. Some organizations have such data available but they did not know how to derive the value of the data since the data is raw.

Quasi structured data: The data formats contain data with inconsistent data formats that are formatted with that and time with some tools.

eg: Web server logs is the log file is created & maintained by some server that contains a list of activities.

## Veracity

Veracity means how much the data is reliable. It has many ways filter or translate the data. Veracity is the process of being able to handle & manage data efficiently. Big data is also essential in business development.

eg:- Facebook posts with hashtags

## Value

Value is an essential characteristic of Big data. It is not the data that we process or store. It is valuable and reliable data that we store, process & also analyze.



5) Velocity:

Velocity plays an important role compared to others. Velocity creates the speed by which the data is selected in real time. It contains the linking of incoming datasets, speeds, rate of change, activity bursts. The primary aspects of Big data in the provide demanding data rapidly.

Big data velocity deals with the speed at the data flows from sources like application, logs Business process network & social media site, sensors mobile devices, etc.

Hadoop architecture:

Hadoop is an open source framework from Apache and is used to store process & analyze the data which are very huge in volumes. Hadoop is written in Java & is not OLAP. It is used for batch processing. It is used by Facebook, Yahoo, Google.

Modules of Hadoop

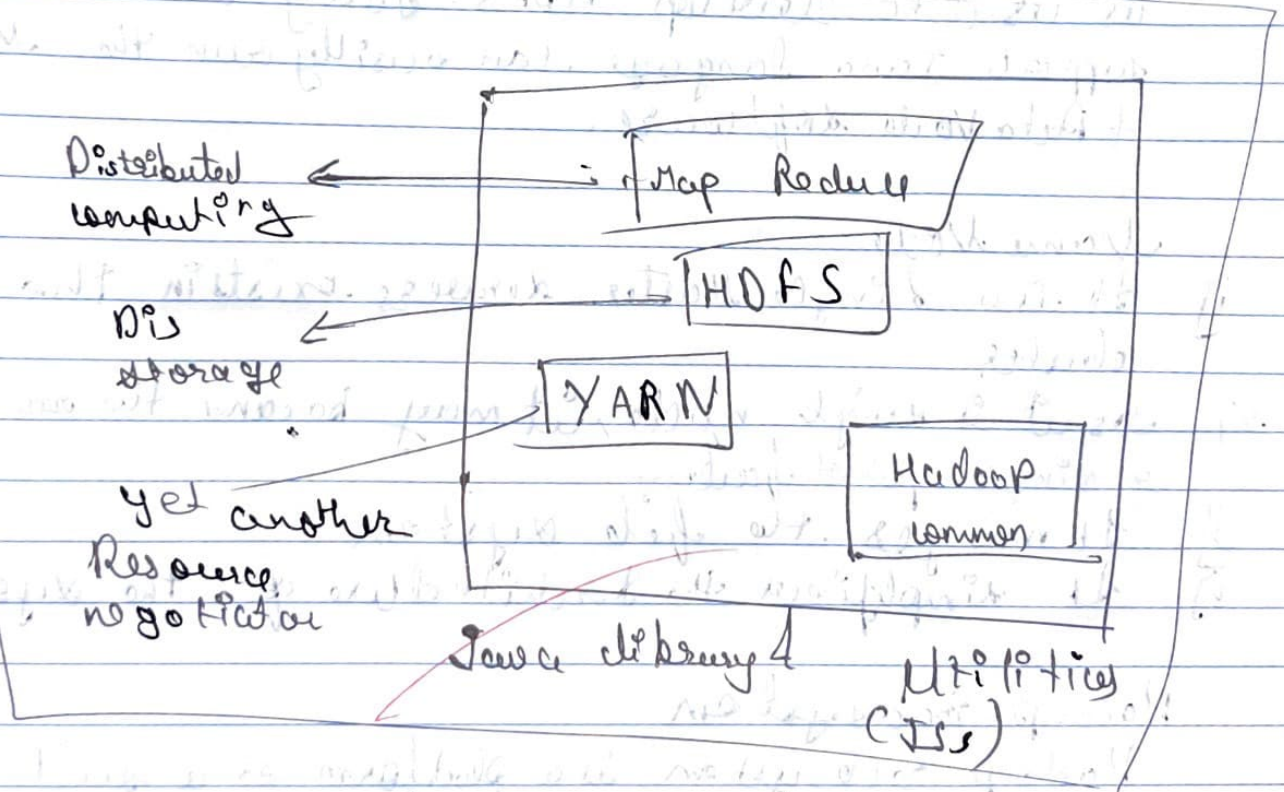
1) HDFS:- Hadoop distributed file system. It is published its paper by FS & on the basis of that HDFS was developed. The idea is that the file will be broken into blocks & it is stored in nodes. The distributed architecture.



2) Yarn :- Yet another Resource Negotiator is used for job scheduling & manage the cluster

3) Map Reduce :- this is a framework which helps with java program to do in the parallel computation on data using key value pair, the Map task takes input data and converts it into data put. Map task is consumed by reduce task & then the out of reducers gives the desired results.

4) Hadoop Common - These three developers are used to start Hadoop and are used by other Hadoop modules.



Hadoop architecture is a package of the file system, MapReduce engine & the HDFS (Hadoop Distributed File System)

A Hadoop cluster consists of a single Master Node & multiple slave Nodes. The Master Node includes Job tracker, task tracker, Nodes.

Hadoop distributed file system the Hadoop Distributed File System (HDFS) is a distributed file system for Hadoop. It contains a Master / slave architecture. This Java language is used to develop HDFS. So any machine that supports Java language can easily run the HDFS & DataNode software.

Master Node

1) It is a single Master server exists in the HDFS cluster

2) As it is single node, it may become the root of single point failure

3) It manages the file system

4) It simplifies the architecture of the system

Hadoop Ecosystem

Hadoop Ecosystem is a platform or a suit which provides various services to solve big data problems. It includes apache projects and various components tools & structure



## Components of Hadoop

- 1] HDFS : File system
- 2] Map Reduce :- Data Processing
- 3] Yarn :- Yet another Resource negotiator
- 4] Spark - In Memory data processing
- 5] Pig/Hive :- Query Based
- 6] HBase - No-SQL Databases

Oozie	Chukwa	Flume	ZooKeeper
-------	--------	-------	-----------

Hive	Pig	Mahout	Storm	Sqoop
------	-----	--------	-------	-------

Map Reduce	YARN
------------	------

HDFS	HBASE
------	-------

1] HDFS: Components Name Node (meta data) + Node (Actual Node) temporary function - stores large datasets across multiple nodes, ensuring fault tolerance

2] YARN  
Component Resource manager in Node Manager  
function - Process large datasets using distributed parallel algorithms



- 3] Map Reduce  
Components - Map Mode Hangover, Reduce  
(Aggregate & summarize data)
- 4] Pig :- Uses Pig Latin (SQL like language) to process & analyse data. Executes commands. MapReduce & stores results in HDFS
- 5] Apache HIVE Mahesant  
function - Provides Machine learning algorithms for clustering, classification, collaborative filtering
- 6] Hive :- Uses SQL (Hive Query Language) for SQL like querying of large datasets supports real time & batch processing
- 7] MapReduce :-  
In memory processing for batch, real time & interactive
- 8] HBase :- NoSQL database for real time & interactive data. Efficient for quick lookups in large data types
- 9] Zookeeper :- Coordinates & synchronizes distributed system, ensuring consistency
- 10] Oozie :- Workflow scheduler. Manages job execution in sequential & event Based coordinated job manner

15/11/2022



## Assignment 2

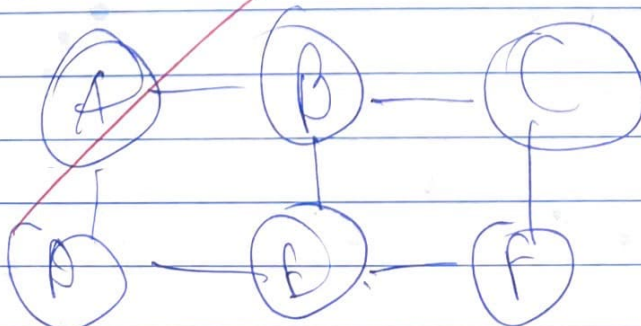


In Kruskal - Newman algorithm we are giving to divide the nodes of the graph into 2 or more clusters. This algorithm removes the edges with the highest betweenness until there is no edges remain or any specific no. of edges is present. Betweenness is the number of the shortest paths between pairs of nodes that are run through it.

### Algorithm.

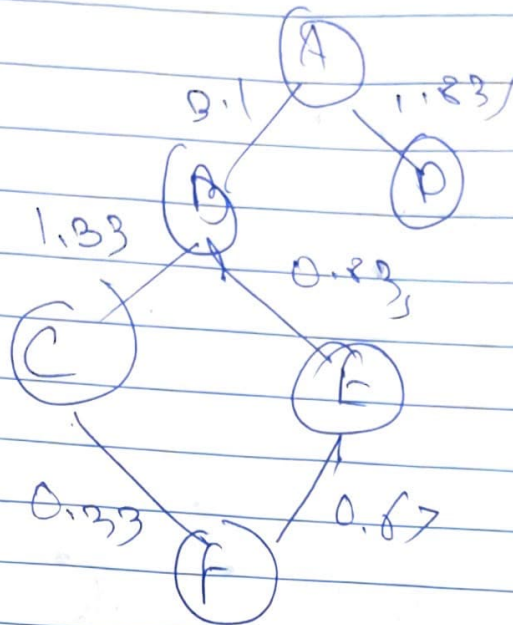
- Create a graph of  $N$  nodes & its edges or take an input graph like Facebook graph.
- Calculate the betweenness of all created graphs in the graph.
- Now remove all the edges with the highest betweenness.
- Recalculate the betweenness of all the edges that got affected by the removal of edges.
- Now repeat the step 3 & 4 until no. edges remain.

### Example

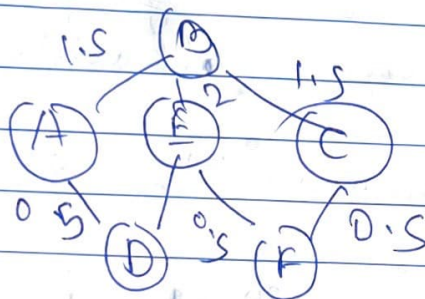




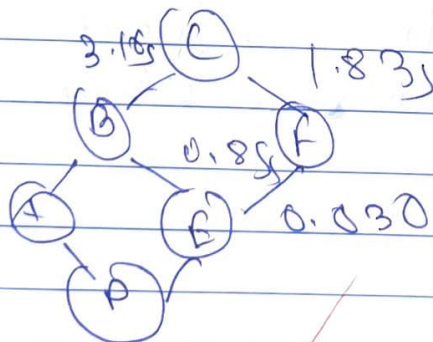
Starting with A



for B

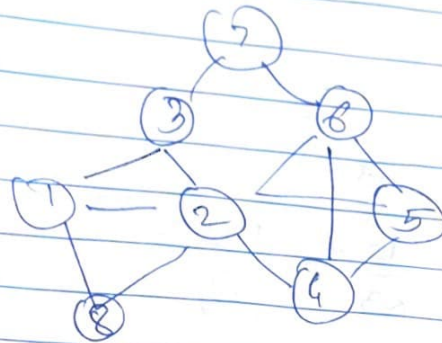


for C



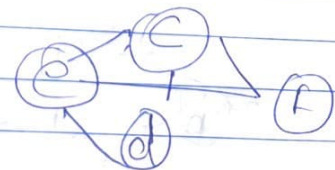


2) clique percolation method builds up the communities from  $k$ -cliques which correspond to subgraphs of  $R$ -nodes  
Ex



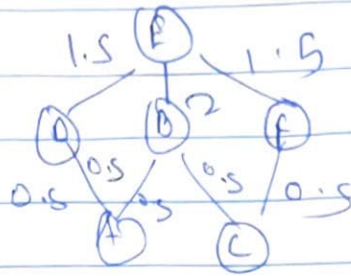
there are six  $k=3$  cliques

a	: 1, 2, 3	> community	9
b	: 1, 2, 8		
c	: 2, 4, 5	> community	1
d	: 2, 4, 6		6
e	: 2, 5, 6	> community	
f	: 4, 5, 6		

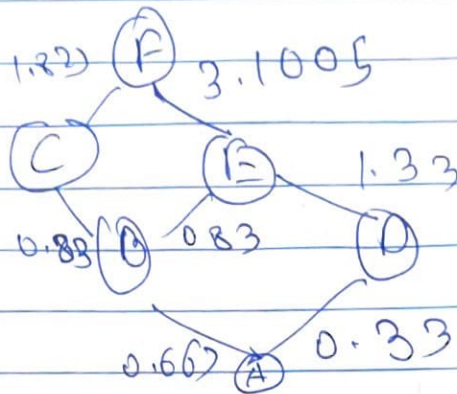




for E



for F

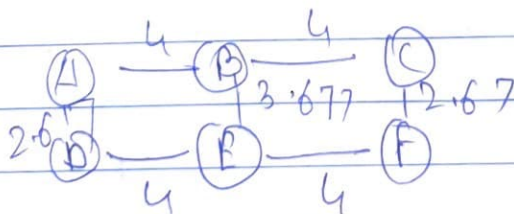


Edges

AB  
AD  
BC  
BE  
CF  
DE  
EF

Edge Between

E  
8.33  
8  
7.34  
5.33  
8  
8



~~AD~~  
15110724