

Module 1

PAGE NO.:

DATE:

Topics covered

- Introduction to big data analytics
- Big data
- Scalability & Parallel Processing
- Designing data architecture
- Data sources
- Quality
- Pre processing & storing
- Data storage & Analysis
- Applications
- Case studies.

What is Big-Data (Data in Petabytes).

: very complex and too large data sets are known as Big Data - which is too complex to handle by traditional data processing softwares.

Defn: A ^{high} volume, high velocity and/or high variety information asset that requires new form of processing for enhanced decision making, insight discovery and process optimization. (Gartner 2012)

3V's are necessary conditions to consider a data set as big data., these 3V's are also known as characteristics of Big Data

Volume: This refers to the amount of data generated by applications in terms of size. This is a measurable unit. The size also determines how to process ^{the data}, as well.

Velocity:

This refers to how fast the data

is getting generated by the applications. This also plays a huge role in processing and analysing the big data.

Variety : In traditional databases the data is well structured and can be stored in rows & columns but now data is available in different formats and forms. The format of the data is to be known to decide the type of Big Data.

One more V is also added along with 3V's, i.e. Veracity - this refers to the quality of data, it is important to know the accuracy of the data so veracity plays a crucial role.

Types of Data : In big data it is important to know the types of data.

There are different types of data generated by the software

1. Structured, 2. unstructured 3. semi-structured
4. multi-structured.

1. Structured Data : The data which is to the point, factual and conform and associates with data schemas and data models. eg: RDBMS tables.

SQL is used to query structured data. It is highly organized and understandable for machine language.

2. Unstructured data : This is raw data. It lacks any predefined model or format. This does not conform to any model or unstructured data needs more storage & providing security is hard. Managing analysing & searching is also very hard.

3. Semi structured data :

This is a structured data but does not follow tabular form. It uses other forms of markers or hierarchy within the data.

Eg: XML, email

4. Multi structured data :

This refers to data consisting of multiple formats of data. structured, unstructured, semi structured. They can have many formats. They are found in non-transactional systems.

Another way of classifying big data is as follows.

1. Social networks & web-data
2. Transactions data & business process
3. customer master data
4. Machine generated data.
5. Human generated data.

Big data classification :

Data sources - eg - RDBMS

Data formats - structured vs semi structured

Big Data sources - eg. distributed file systems.

Big data formats - unstructured, semi structured & multi structured.

Data stores structure - web, cloud servers.

Processing data sets - eg. Batch, real time.

Processing big data

ratio - High volume, velocity

Analysis types - Batch, scheduled

Big data processing

methods - Batch processing,
real time processing.

Data analysis

methods.

- statistical analysis
- predictive analysis

Data usage

- Human, Business
process.

Big Data Handling techniques:

1.3 : Scalability & Parallel processing

- Hundreds of terabytes & to peta bytes of data
- Needs hundreds of computing nodes.
- Requirement here is processing of data in short time at minimum cost.

Scalability enables increase or decrease in the capacity of data storage, processing and analytics. It is the capacity of the system to handle the magnitude of the work.

When the workload increases the capability needs to increase. When the workload & complexity exceed the system capacity scale it up & scale it out.

Vertical Scalability : it means scaling up the given systems resources & increasing the systems analytics, reporting and visualization capabilities.

Scaling up : Designing the algorithm according to the architecture that uses resources efficiently..

e.g. :

$$\begin{aligned} \text{Data} - & x \text{ terabytes} \\ \text{processing time} - & t \\ \text{code size} - & n \quad (\text{increased with complexity}) \end{aligned}$$

Scaling up means processing takes equal less or much less than ($n \times t$).

Horizontal Scalability

Increasing the number of systems working in coherence & scaling out - the workload.

Scaling out - using more resources & distributing the processing & storage tasks in parallel.

γ resources.

x terabytes of data in

t time.

then $(\gamma \times x)$ terabytes of data processed on γ parallel distributed nodes.

Easy way to scale up & scale out is buying faster CPU's bigger & faster RAM etc. but this is expensive.

Alternatively we can go for

MPP's (massively parallel processing platforms) cloud, grid, cluster & distributed computing s/w.

i) massively parallel processing platforms (MPP):

Many programs are so big and complex that they can't be run on a single CPU. The solution is either ~~not~~ scale up the computer system or. up MPP platforms parallelization can be done at 3 levels

i) ~~not~~ thread level

ii) multiple CPU's

iii). Separation computing

Distributed computing model :

it uses cloud, grid or clusters which processes & analyzes big & large data sets on distributed computing nodes connected by high speed networks.

2. Cloud computing : Internet based computing that provides shared processing resources & data to the computers and other devices on demand.

benefits : perform parallel & distributed computing in a cloud computing environment

- multiple nodes perform automatically & interchangeably.
- offers high data security
- Failure of one node does not cause failure of entire system.

Eg: AWS - Amazon web services. Elastic compute cloud. (EC2)

Microsoft Azure or Apache cloudstack.

features of cloud computing -

- a) On demand service
- b) Resource pooling
- c) Scalability
- d) Accountability
- e) Broad n/w access -

Types of cloud computing - H/IAD , n/us connection
database storage, data center etc.

① Infrastructure as Service (IaaS)

- eg: Tata communications

Amazon data centers & virtual servers

Apache openstack — open source.

② Platform as service (PaaS) — Provides run time environments to allow developers to build applications & services.

eg: Hadoop cloud services.

③ Software as service (SaaS) — Providing software applications as a service to end users.

eg: SQL geogSQL, IBM BigSQL,
microsoft polybase oracle big data SQL.

④ Grid and Cluster Computing :

group of computers from several locations are connected with each other to achieve a common task.

features of grid :

- scalable .

- sharing of resources

- distributed nature of resource utilization .

drawbacks — single point failure leads to failure of entire grid

cluster computing :

- Group of computers connected by a n/w to accomplish a single task.
- It is mainly used for load balancing.
- Among the group of computers the processes are shifted between the nodes for load balancing.
eg: Hadoop architecture.

Differences

<u>Distributed computing</u>	<u>cluster computing</u>	<u>Grid computing</u>
① Loosely coupled	Tightly coupled	large scale.
② Heterogeneous	Homogeneous	Cross organizational
③ Single administration	cooperative working	Geographical distribution
		Distributed management

④ Volunteered Computing

- Distributed computing paradigm which uses resources of volunteers - they are organizations or members who own personal computers.

issues

- heterogeneity
- Drop outs from the n/w.
- Sporadic availability
- Incorrect results are unaccountable due to anonymity.

Designing Data architecture

Layer 5 Data consumption	Export of data sets to cloud web etc	Data Sets usage BPs, BI's Knowledge discovery	Analytics. (real time, near-real time).
Layer 4 Data processing	Processing technology: Map reduce, Hive, Pig, Spark	Processing in real time: scheduled batches or hybrid	Synchronous or asynchronous processing
Layer 3 Data storage	Considerations of type, formats compression frequency	Hadoop distributed file system	NO SQL databases Ingestion of structured data MongoDB, MySQL
Layer 2 Data ingestion & acquisition	Ingestion using Extract load & Transform (ELT)	Data Semantics such as replace, append ...)	Ingestion of data sources in batches or real time
Layer 1 Identification of internal & external sources of data	Sources for ingestion of data	Push or Pull of data from sources for ingestion	Data types for files, web - Data formats structured, unstructured

Design of logical layers in data processing architectures & functions in the layers.

Data sources, quality, pre processing & storing :-

- Applications
- programs & tools use data.

Sources can be external

Eg: sensors, trackers, weblogs, computer systems logs & feeds.

Sources can be machines, which source data from data creating programs.

- Data sources can be structured, semistructured or unstructured.

- can be social media.

- can be internal

- can be data repositories such as database, tables, flat file, spread sheet, mail server, web server, CSV

- Sources may be data store for applications.

- structured.

- unstructured.

- Sensors signals & GPS.

Data Quality :

High quality data - enables all the required operations, analysis, decisions planning & knowledge discovery correctly.

A high quality data has 5 R's.

- Relevancy - this is utmost importance.
- Recency
- Range
- Robustness
- Reliability.

Data integrity: refers to the maintenance of consistency & accuracy in data over its usual life.

Data Pre processing:

- Step before data ingestion :
- A must before running a machine learning algorithm.
- Prior screening is of data quality is needed.
why is it needed.
 - Dropping out of wrong, inconsistent & outlier values.
 - Filtering unreliable, irrelevant & redundant info.
 - Data cleaning, editing, deduction or wrangling.
 - Data validation, transformation or timestamping.
 - ETL processing.

Data cleaning :

- Process of removing or correcting incomplete, incorrect, inaccurate or irrelevant parts of the data after detecting them.

e.g.: collecting the grade outliers or mistakenly entered values means cleaning & correcting the data.

Data cleaning tools :

It is done before mining of data.

Incomplete or irrelevant data may result

info misleading decisions.

- It helps in refining & structuring data into usable data.

eg.: OpenRefine & DataCleaner.

Data enrichment:

Data enrichment refers to operations or processes which refine, enhance or improve the raw data.

Data Editing :- process of reviewing & adjusting the acquired datasets.

- i) Interactive
- ii) selective
- iii) automatic
- iv) aggregating

Data reduction :-

Enables the transformation into an ordered, correct and simplified form. Only meaningful data will be ingested in the data base.

Data Wrangling: refers to the process of transforming and mapping the data.

eg: mapping enables data into another format, which makes it valuable for analytics and data visualization.

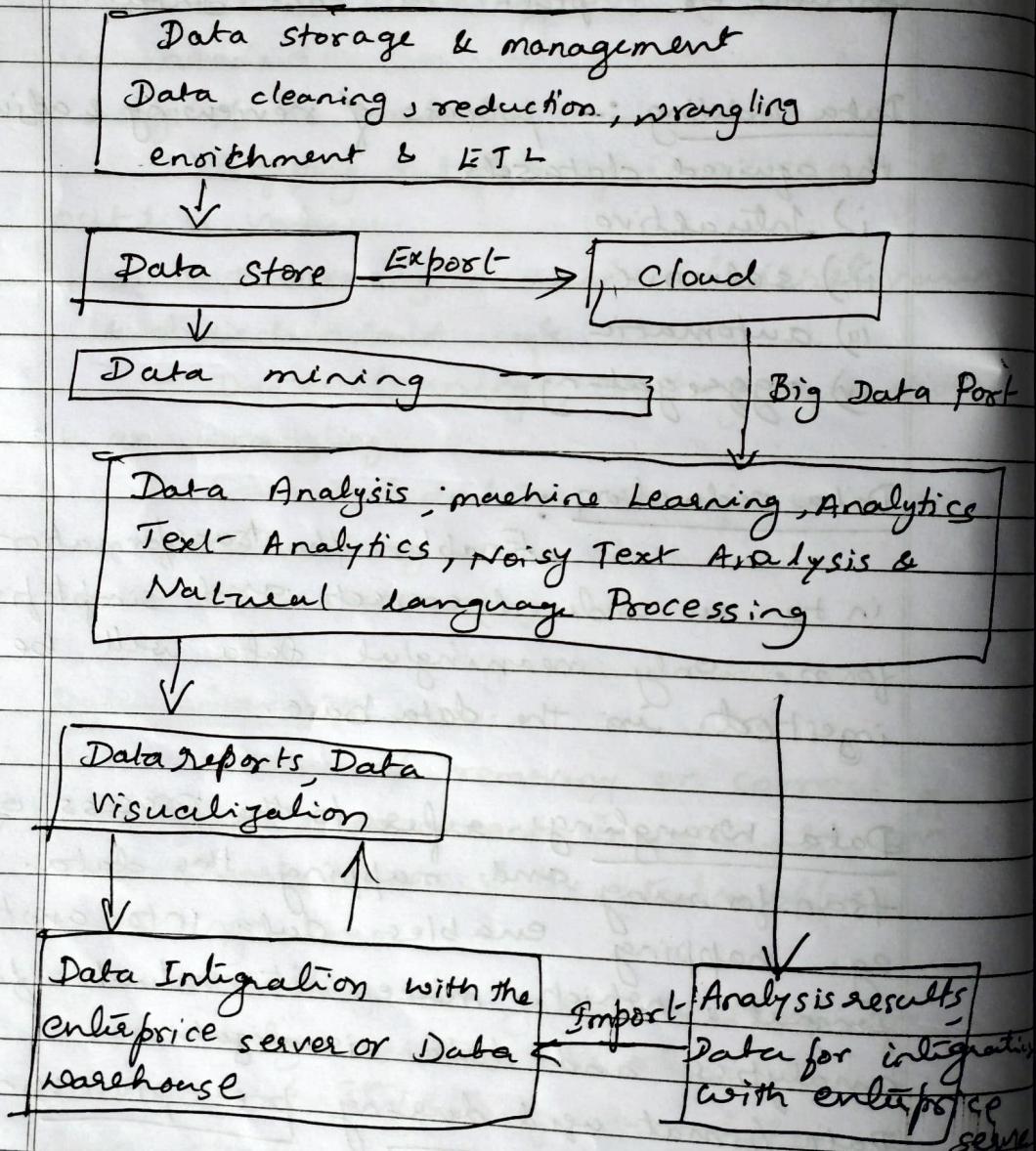
Data format used during pre processing

examples of formats for data transfer from
a) Data storage b) Analytics c) Service
or d) cloud can be

- i) Comma Separated Values (CSV)
- ii) JSON - Java script object notation
- iii) Tag length values (TLV)

CSV format :

Data store export to cloud



Data Pre processing analysis, visualization
data store export.

PAGE NO. _____
DATE: _____

Cloud services: cloud offers various services. These services can be accessed through a cloud client, such as a web browser or web services.

Data Storage and Analysis:

The data store in traditional database management systems and big data are discussed here.

Big Data storage

Big data NoSQL or Not only SQL

- No SQL databases are considered as semi-structured data. Big data stores uses NoSQL.

NoSQL stands for No SQL or Not-only SQL. The store donot integrate with applications using SQL.

NoSQL is also used in cloud data store.

i) ~~day consis~~: It is a class of non relational data storage systems, and the flexible data models & multiple schema.

ii) class consisting of unstructured key / value or big hash table.

iii) class consisting of ordered keys & semi structured data storage systems.

iv) class consisting of JSON (MongoDB)

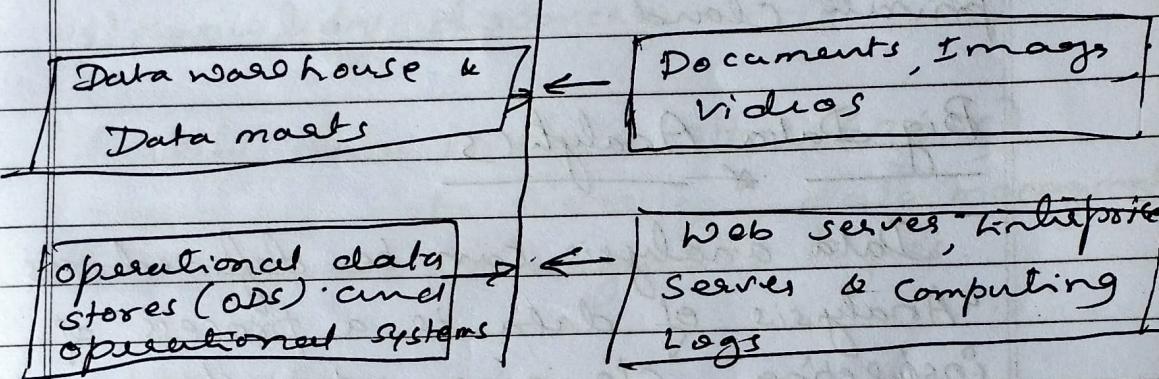
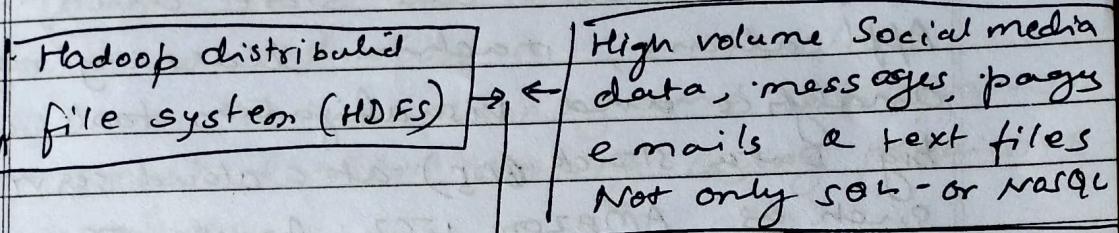
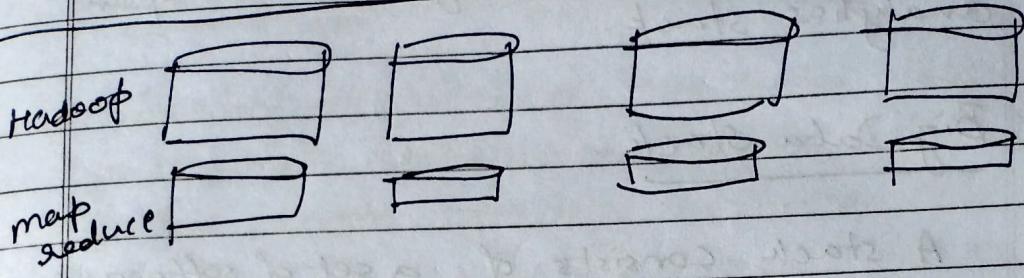
- i) class consisting of name / value in the text (couched DB)
- ii) may not use fixed table schema.
- iii) Do not use JOINs.
- iv) Data written at one node can replicate at multiple nodes, therefore data storage is fault tolerant.
- v) may relax ACID principles during the data store transactions.
- vi) Data store can be partitioned and follows CAP theorem

coexistence of Big Data, NoSQL and Traditional data stores.

Big Data Platform

- Big data platform should provision tools and services for.
 1. Storage, processing and analytics.
 2. Developing, deploying, operating & managing big data environment.
 3. Reducing the complexity of multiple data sources & integration of applications into cohesive solution.
 4. custom development, integrating & integration with other systems &
 5. the traditional as well as Big data techniques.

Hadoop



Hadoop based Big Data Environment.

Big data platform consists of Big data storage, servers & data management & business intelligence s/w.

Storage can deploy Hadoop distributed file system (HDFS), NoSQL data stores, such as HBase, MongoDB, cassandra. HDFS system is an open source storage system.

Mesos:

Mesos v0.9 is a resource management platform which enables sharing of cluster of nodes by multiple frameworks &

which has compatibility with an open analytics stack.

Big Data Stack

A stack consists of a set of software components and data store units. Applications, machine learning algorithms, analytics and visualization tools use big data stack (BDS) at a cloud service such as Amazon Web Services or private cloud.

Big Data Analytics

Data analysis can be defined as "analysis of data is a process of inspecting, cleaning, transforming to model data with the goal of discovering useful information, suggesting conclusions and supporting decision making."

Data analytics can be defined as the statistical and mathematical data analysis that clusters, segments, ranks and predicts future possibilities.

Phases in Analytics

Descriptive analytics

Predictive analytics

Prescriptive analytics

Cognitive analytics

Berkeley Data analytics stack (BDAS)

The importance of Big Data lies in the fact that what one does with it rather than how big or large it is.

- 1) Cost reduction
- 2) time reduction.
- 3) new product planning and development.
- 4) Smart decision making using predictive analysis.
- 5) Knowledge discovery.

Berkeley data analytics stack (BDAS) consists of data processing, data management and resource management layers.

1. Applications .
2. Data Processing combines batch, streaming and interactive computations
3. Resource management software : component provides for sharing the infrastructure across various frameworks.

1.7 Big Data Analytics applications and case studies :-

- Big data in marketing & sales.
- Big data analytics in detection of marketing frauds.
- Big data risks.
- Big data credit risk management.
- Big Data and algorithm trading.

Big data & healthcare.

Big data in medicine.

Big Data in advertising
