

(DA) Data Analytics

08/01/2023

Data:

- ↳ Data → Information → Knowledge → Actionable insights
- ↳ since the mid-1900s, people have used the word data to mean computer information that is transmitted or stored
- ↳ Data means stored in a computer and the data are various types data is important because the
- ↳ it is interpreted, by a human or machine to derive meaning
- ↳ present heterogeneity since as well as heterogeneous source
- ↳ need of the hour is to understand, manage and take data, for analysis to derive valuable insights

Importance of Data:

- ↳ searching data and analyze data and get concluded to user preference and analyze and send notification on user
- ↳ Boost Company Revenue (ex: amazon)
- ↳ Customer demand

Characteristics of Data:

- ↳ Data → Composition, Condition, Context
 - ↳ Deal with the structure Data analysis use
 - ↳ Deal with 'where' has been generated

Human vs. Machine Readable Data

- ↳ Human-readable refers to the information that only human can interpret, it required a person to interpret it, that information is human-readable
- ↳ Machine readable information is a set of instructions for manipulating data.
ex. CSV, JSON, XML.

Classification of Digital Data

- 1) structured Data (10%)
 - 2) semi-structured Data (10%)
 - 3) unstructured Data (80%)
- } Data (World)

Structured Data

- ↳ stored in table form (pattern makes it easier for any person to sort, read & process data).
- ↳ Any type of person (using ~~SQL~~ SQL) can query predefined form, stored table form, reside in a fixed field within a record field, attributed mapped, used to query and report against predetermined data types.

Structured Data

↳ Relational Database

- ↳
- ↳
- ↳

↳ Insert / update / delete, Security, Indexing, scalability, Transaction processing.

Initial to final state ←

Semi-structured Data

- ↳ schema-less or self-describing structure
- ↳ ex: email, XML, markup language like HTML
- ↳ rows and columns.

Semi-structured Data

- ↳ web data in the form of content
- ↳ XML
- ↳ JSON
- ↳ other Markup language

- Characteristics
- ↳ Inconsistent structure
 - ↳ self-describing (level / values)
 - ↳ other schema information is blended with data values
 - ↳ Data objects may have different attributes not known beforehand

~~relation~~

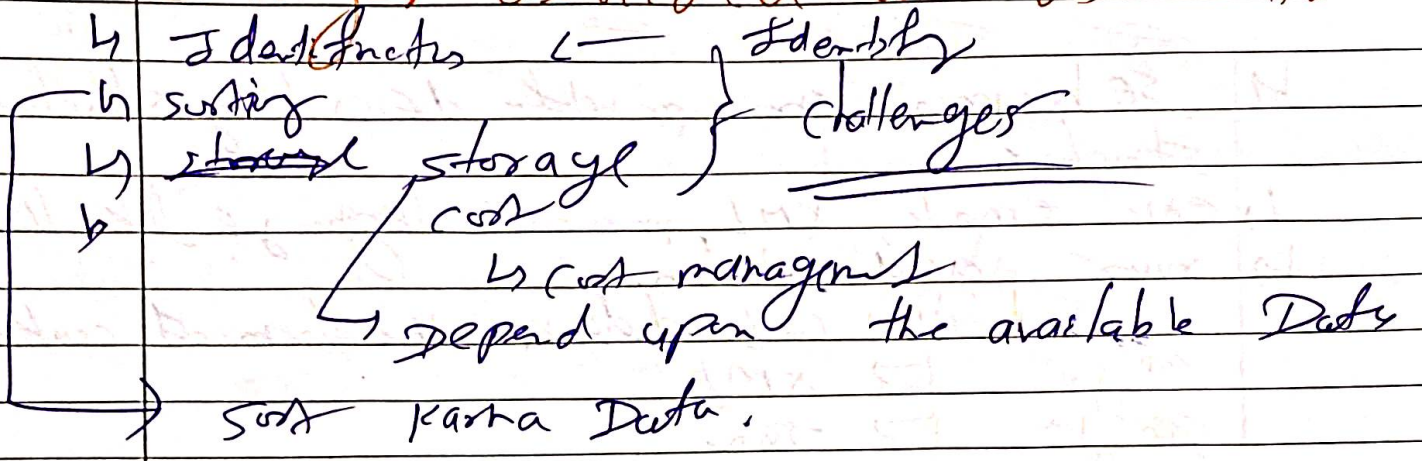
Unstructured Data

- ↳ Data may have no logical, logical or positional
- ↳ 80% of unstructured Data
- ↳ typically consists of metadata i.e. additional information related to data

- unstructured data → text both internal and external to org.
- ↳ Social media data
 - ↳ Body of email
 - ↳ chat, Text messages
 - ↳ media data
 - ↳ images, audio, videos.

02/01/2024

Challenges associated with Unstructured Data



- Dealing →
- ↳ Data Mining (DM)
 - ↳ Natural Language Processing (NLP)
 - ↳ Text Analytics (TA)
 - ↳ Noisy Text Analyser

Data Mining → store Data in one store and analysis karne ke bad ek solution mil rha hai us domain mai.

~~Natural~~ Natural Language Processing → Google Assistant

Text Analyser → Unstructured Data

Big Data → High Volume → velocity → variety

why → structured Data + semi-structured Data + Unstructured Data = Big Data

Elements of Big Data :-

Volume, variety, velocity, veracity

more data → More accurate analysis → Greater confidence

Greater operational efficiencies, cost reduction, time reduction, new product development, and optimized offerings etc. in decision making

Challenges of Traditional system

- ↳ outstructure data → (on structure Data)
- ↳ RDBMS using Relation (No Denormalize table) associated Data → (on structure Data)
- ↳ Use ETL process (extract, transfer and load)
- ↳ Active Parallelism with the help of hardware
- ↳ in dependent support of aggregated summary Data.

Difficult to manage

Condition

- ↳ Data challenges → Volume, velocity, variety
- ↳ process challenges → Data Discovery
- ↳ Management Challenges → Scalability

Capturing Data, Aligning Data, Transferring Data, Storable Data, Modeling data, security, privacy, Governance, Ethical use

Web Data :-

- ↳ Available public Domains
- ↳ document Pdfs, docx, Plan text

↳ Best - from source
↳ Data unstructured and in appropriate.

*) parallel computing vs Distributed computing

- | | |
|---|---|
| ↳ shared memory sys. | Distributed Memory sys. |
| ↳ Multiple processing stage
a single bus and memory used | Autonomous computing nodes
connected via network |
| ↳ processing in order of tasks
Tbps | apps |

↳ Limited Scalability

Better scalability and cheap

Distributed computing in local network (called cluster computing), Distributed computing in wide area network (grid computing)

*) EDW, OLTP, MPP

- ↳ Enterprise Data Warehouse
- ↳ Online transaction processing
- ↳ Massively parallel processing

*) Hadoop

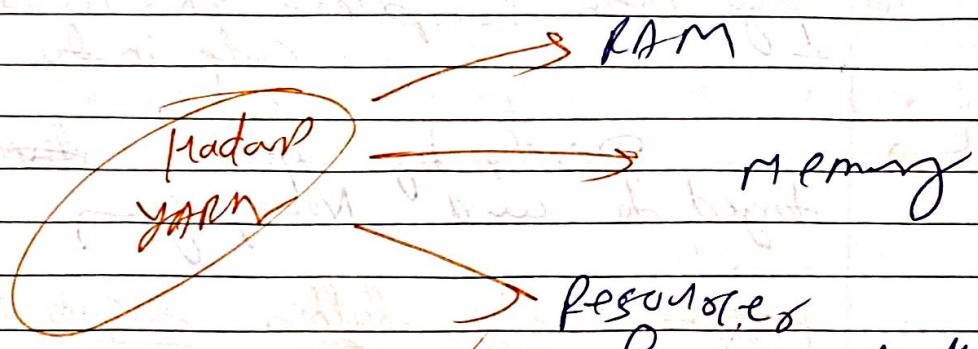
Store Big Data using Distributed

- ↳ store Huge Amount of Data and process
- ↳ Computers (Distributed nodes) / provide work
- ↳ High Computational Power, Fault tolerance, Flexibility, Low cost, scalability

19/10/2024

YARN: (Resource Management Unit)

Jisko use karke Data ko process krna hota hai



Yarn is an acronym for yet another Resource Negotiator. It handles the cluster of nodes and acts as Hadoop's resource management unit. YARN allocates RAM, memory, and other resources of yarn.

↳ Resource Manager (Master)

↳ Node Manager (Slave)

MAP Reduce: (Data processing)

1) The Mapreduce perform the processing of data dataset in a distributed parallel manner.

2) The Map reduce consist of 2 steps
1) Map → to divide the data in a customized way
2) Reduce and combines them into key value pair

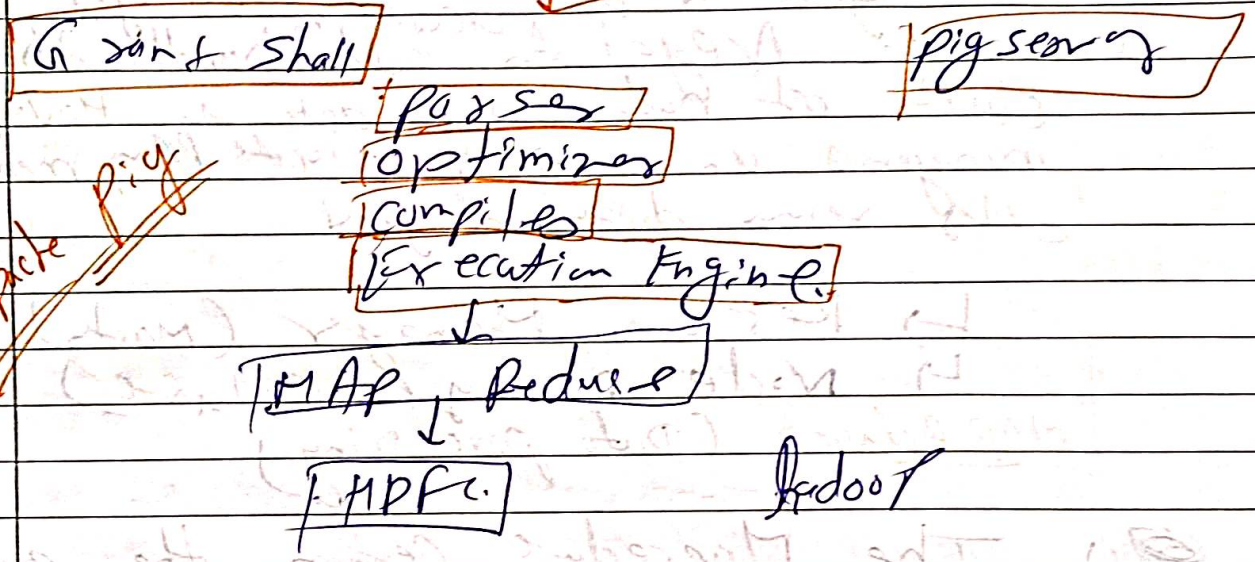
↳ To collect intermediate Results

Apache Pig

- 1) It consists of Pig Latin which is the Java script
- 2) Pig Latin compiler which converts pig latin code into executable code

Developed by ~~apart~~ ~~level~~ yahoo research, targeted to work Non programming

Pig Latin Script



Hive

- ↳ Hive Command line
- ↳ JDBC/ ODBC driver

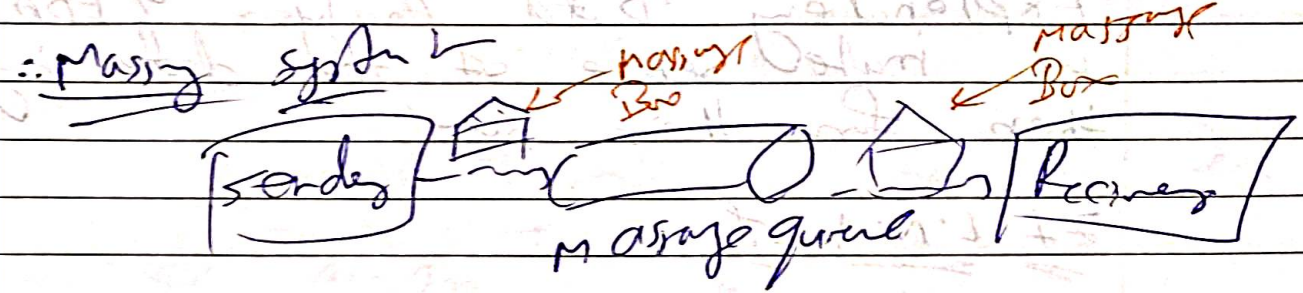
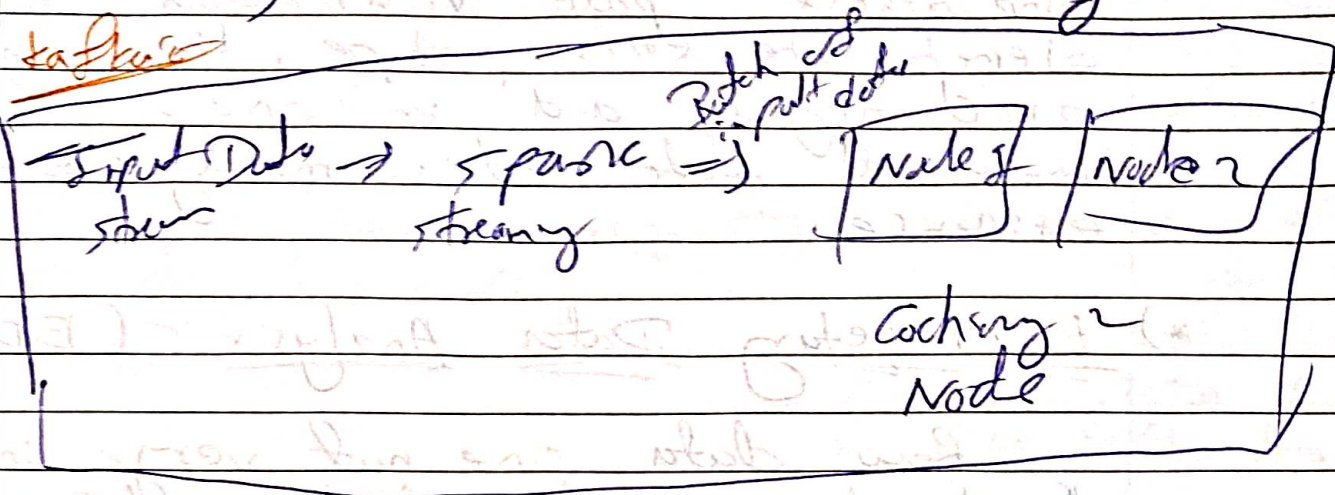
- ↳ Spark open source
- ↳ All Data process
- ↳ processing

MapReduce

MapReduce builder machine learning Algorithm

↳ provide platform

Ambaioir ↳ Decided way cluster
↳ How many applications running



→ Streamio ↳ Real time streamy of Data processing

↳ Ranger security ~~input~~ security implerid in Ranger
↳ Apache Knox proxy services
Authentication services
Client services

↳ Oozie 1) workflow engine
2) condiditry engine

① Data Explacation ↳
Data explanation is the first step.
in Data analysis involving the use of
data visualization tool and statistical
te unferen - data set chendents

and initial pattern.

Importance:- Data visualization tools and elements like color, shape, lines, graphs and angles aid in challenging & recognizing flow or data index among without existence.

2) Exploratory Data Analysis:- (EDA) :-

Raw data are not very informative. Exploratory Data Analysis (EDA) is how we make sense of the data by analyzing them from their raw.

Its Limits of:-

- i) Analyzing and summarizing the new data.
- ii) Discovering important features and patterns in the data and any striking dependencies from those patterns.
- iii) Interpreting our findings in the context of the problem.

3) Sampling:-

→ Sampling is a technique of selecting individual members as a subset of the population to make a statistical conclusion on the basis of evidence from them and estimate characteristics of the whole population.

Sampling Methods:-

Probability Sampling

Non-probability Sampling

Any element can be chosen randomly from the population. It deals with choosing the sample randomly.

Every element will be chosen in the subjective judgment of the person on the basis of their past experience & knowledge rather than non-selection.

Non - Sampling Methods

Every element will be chosen in the subjective judgment of the person on the basis of their past experience & knowledge rather than non-selection.

~~06/01/2019~~

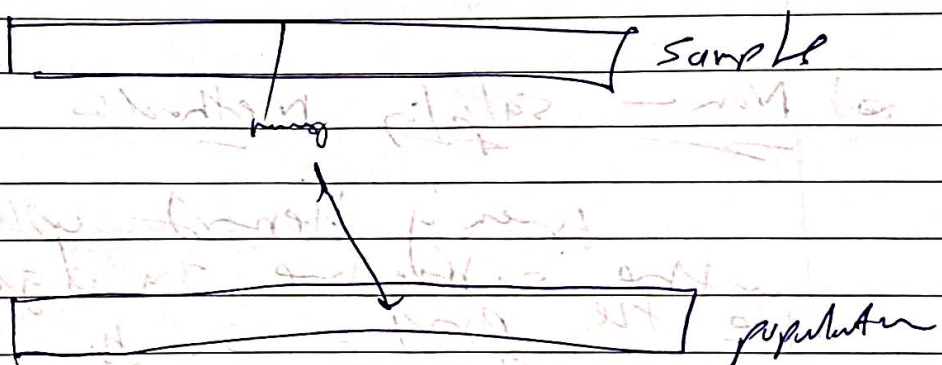
Stats checking

12/02/2024

Estimation:

- ↳ It is often of interest to learn about the characteristics of a large group of elements such as individuals, households, business units.
- ↳ Data from the sample are then used to check estimates characteristics of the larger population.

Point Estimator



Point Estimator Condition

- ↳ Most often, the methods of July the parts of large population as the error are subtle average, age of baby budget, it is impossible to check or count age of every person in the world.

Interval Estimation

- ↳ Can include Interval Estimation can support a student meaning the boiling temperature of a certain 100°C

12/04/2019

Observed the reading (in degree Celsius) for 10's, 100's, 1000's, 10000's, 100000's and 1000000's as a distribution sample of the liquid about is the interval estimation for the population mean at a 95% confidence level

mean is 106.82 , $n=10$

$$\mu = 106.82 \pm 1.96 \cdot (0.984/\sqrt{10}) = 106.82 \pm 0.60$$

Standard error (SE) $\sigma/\sqrt{n} = 0.984/\sqrt{10}$

Margin of error $z^* (\sigma/\sqrt{n}) = 1.96 \cdot (0.984/\sqrt{10})$
 $= 0.60$