

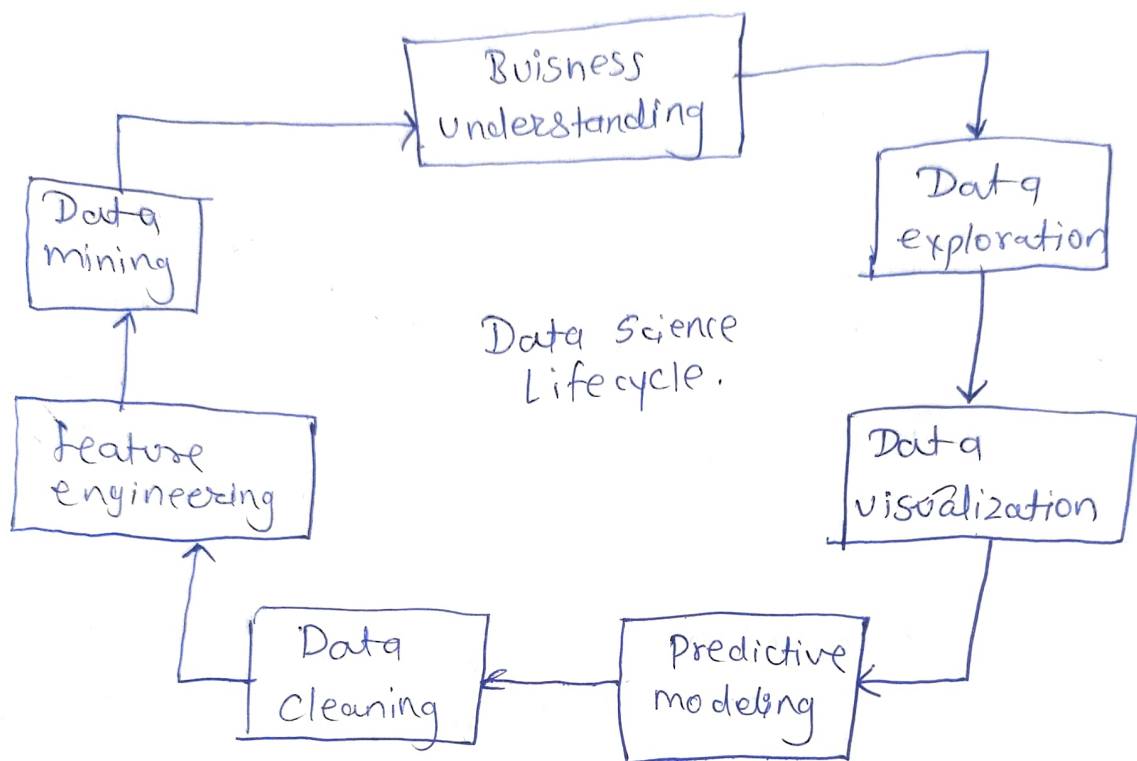
Q1] Explain Data Science & Application of Data Science

- Data is a collection of facts and figures which relay something specific. It can be number, word, measurement or description of something. Data is a raw material.
- Types of data are record data, data matrix, document data, ordered data, graph data etc..
- Data Science is the field to extract the ~~data~~ information from various form of data.
- Data science aims to discover knowledge from the data which can be used in business decision & prediction.
- Data Science can help in business Analysis. From historical data, instead of knowing how many product sold it forecast the future sale of product.

Applications of Data Science. →

- 1] Healthcare
- 2] Gaming
- 3] Image Recognition
- 4] Logistics
- 5] Predict Future market trends
- 6] Recommendation System.

Q] Explain Data Science Life Cycle.
→ Data Science lifecycle is as follow



a] Business Understanding → Understand problem to solve

b] Data exploration → Understand pattern & bias your data.

c] Data visualization → Create and study of the visual representation of data.

d] Predictive modeling → It is the step where the machine learning finally comes into your data.

e] Data cleaning → Detect and correct corruption inaccurate data.

f] Feature engineering → Process of cutting down the feature.

Q] Data mining → Gathering your data from different sources.

Q] Data explosion plays major Role in big data. justify the Statement with proper explanation along with examples. & Factor responsible for data explosion

or

Q] Explain 7 V's of Big Data.

- Data explosion means rapid growth of the data from different resources and stored in computer system is called data explosion.
- Data is generated automatically through mobile devices, computer system, etc...
- The phenomena of exponential multiplication of data that get stored is termed as data explosion.
- Sending email, making phone calls, collecting info each day we create massive amount of data.

* Role of Data Explosion in Big Data *

- 1] Volume : Volume refers to the size of the dataset.
- It could consists of billions of rows & millions of columns.
 - Usually this dataset is stored in multitiered storage.
 - for example imagine a sale on Flipkart, a person views 10 mobile a day of 5 different companies each it would make 50 data points now let's suppose Flipkart have 1 million active users on sale day it would make 50 million data points on single day.

2] Velocity → Velocity refers to the rate of speed at which data is generated & processed. It can be ^{based on} real time or based on historical in nature.

- Examples → In stock market millions of transactions are made in real time.

3] Variety → It refers to the data accumulated from multiple data sources. It can be structured, semi-structured or unstructured.
eg → Ecommerce collect ^{structured (sale orders), semi-struct (review)} ^{unstructured (photos & videos)}

4] Veracity → It refers to the measure of ~~of~~ data quality & usefulness of the data.

- The data should ~~have~~ be true or relevant we could not perform useful analysis if incoming data is false or has error.
eg → fake news.

5] Value → Value measures the usefulness of data which helps in improving business decisions or enhancing AI models.

- eg → Netflix uses data to recommend shows to users.

6] Variability → The changing nature of data, including seasonal trends & inconsistencies.

- eg → Google news.

7] Visualization → The ability to represent complex Big data in visualizing or understandable way using charts, graphs & dashboards.

- eg → COVID-19 uses heat map to show spread of virus.

* Factors responsible for data Explosion *

- 1] Growth of internet
 - 2] Social media & Digital Content.
 - 3] Smart Devices
 - 4] Ecommerce & Online Transaction
 - 5] AI & ML
-

Q] List and explain data processing infrastructure challenges in Big data with suitable example.

→ Big Data requires a strong infrastructure to process large volumes of data efficiently.

• Challenges in big data infrastructure are →

1] Data storage → The increase in volume of data, increases the need for storing data. It requires the medium with higher I/O speed to store the data. Traditional database struggle to handle such large datasets.

eg → Facebook stores petabyte of data which requires hadoop HDFS.

data.
2] Processing Speed → Processing high-velocity data in real time is difficult, especially in applications that requires instant responses.

eg → Stock market platform have millions of transactions per second, so they use Apache spark

3] Data Integration → Big data comes with ~~various~~ various sources like databases, sensors, api's etc. Integrating different format & structure is complex.
eg → E-commerce platform integrates data from review, rating, purchase in real time update.

4] Scalability → As data growth exponentially, traditional IT infrastructure fails to scale effectively.

eg → Netflix uses cloud based AWS for streaming.

5] Data privacy & Security → Protecting data from cyber ~~the~~ threats, unauthorized access ~~data~~ is ~~be~~ crucial. Data must be encrypted & should not shared.

6] Higher Cost → Growth in data increased the demand of data storage ~~att~~ which is directly proportional to Higher Cost.

7] Data Quality & Cleaning → Big Data comes with errors, missing values, inconsistency, which affects decision making.

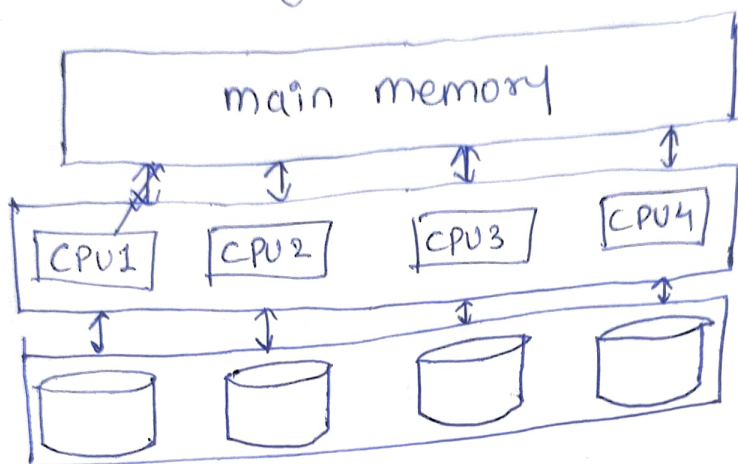
Q] Difference between Big Data & Small data, with processing architecture.

→ feature	Big Data	Small Data.
1] Size.	Very large.	Small.
2] Speed.	Processed in real time.	slowly or manually
3] Sources	social media, transaction etc..	Excel files, small databases
4] Processing method	uses distributed computing (Hadoop, Spark)	Traditional processing (Excel, SQL)
5] Storage	stored in cloud i.e. HDFS or NoSQL databases	stored in local files, spreadsheet, SQL DB
6] Complexity	structured, semi-structured or Unstructured	mostly structured
7] Scalability	Highly Scalable	Limited scalability
8] Example	Netflix recommendation, fraud detection in banking	student attendance, sale tracking

QJ Explain Big Data processing Architecture

or
Explain shared Nothing, shared everything architecture.

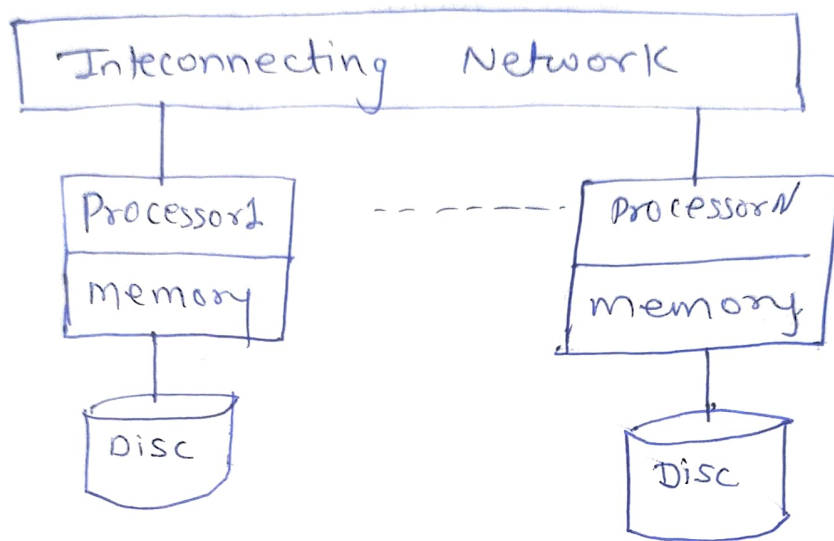
→ shared everything Architecture



- ~~This~~ In shared everything architecture, all server are ~~int~~ connected or all server shares the same memory. storage & has access to the store.
- Main idea behind this system is maximum resource utilization. Disadvantage is performance.
- Scalability is main problem
- Symmetric multiprocessing & Distributed shared memory are two types.
- In Symmetric multiprocessing (SMP) architecture, CPU shares a single pool of memory for read-write operations, sometimes also called as uniform Memory Access (UMA) architecture.

- Distributed shared memory (DSM) addresses the scalability problem by providing multiple pools of memory for processors to use. It is also called as non-uniform memory access (NUMA) architecture.

* Shared Nothing Architecture *



- In shared nothing architecture it consists of multiple nodes having ~~at~~ own OS.
- Each node is connected with other using interconnecting network.
- Each node contains its own memory (M), processor (CPU) & ~~shared~~ storage device.
- Each node is in control of its own OS.
- Data is partitioned horizontally across nodes.

8] Difference Between Data Warehouse & Data mining

→ ^{function} Definition	Data warehouse.	Data mining.
1) Definition	A storage system that collects and organise data.	A process of analysing data to find patterns & trends.
2) Purpose	Storage large amount of data.	Extract useful insights from stored data.
3) process	Data is collected, cleaned & stored from diff resource	Data is analysed using algorithms to find hidden patterns.
4] Technique used	OLAP.	ML, Clustering, Classification
5] Tools.	Snowflake, Google Big Query	Python, R.
6] who use it?	Data engineer & analyst.	Data scientist & AI/ML experts.
7] Example.	A bank store transaction details in data warehouse	A bank used Data mining for fraud detection.

Q] Define the relation between Artificial Intelligence, statistical learning & machine learning

→ AI, statistics & ML are interconnected fields, but they differ in scope & application.

AI → It is the field that focuses on creating machine that can perform tasks that typically requires human intelligence such as reasoning, problem solving decision making

ML → ML is a subset of AI that enables computers to learn from data without being explicitly program

Statistical Learning → It is a model that helps in predictions & relationships using statistical models

Relation

1] AI → Intelligence which include ML model to fulfill requirement

2] ML → Uses Statistical Learning technique to train model from data.

* Statistical Learning → provides model to ML Algo's.

eg → fraud detection system in bank.

1] Statistical Learning → Logistic Regression analyses past fraudulent transaction

2] ML → Random forest improves accuracy

3] AI → make decision based on ML model.