

Assignment - 1

Sec-A

- ① Big Data Integration refers to the process of combining data from different sources into a unified view. It involves extracting data from various Structured, Semi-Structured and Unstructured Sources, transforming it into a common format, and loading it into a central system for analysis.
- ② Value in Big data refers to the insights and actionable information that can be extracted from large datasets. The value comes from analyzing data to uncover trends, patterns and correlations that can drive business decisions, improve efficiency and create competitive advantages.
- ③ Data explosion refers to the rapid and exponential growth in the volume of data being generated from various sources like Social media, sensors, devices and business processes. This increases in data requires advanced storage, processing and analysis techniques to extract meaningful insights.
- ④ A NoSQL model refers to a type of database that is designed to handle large volumes of unstructured & semi-structured data. Unlike traditional relational databases, NoSQL databases use flexible data models such as key-value pairs, document-based, column-based or graph-based models. They are optimized for scalability and fast data access.

1. HDFS (Hadoop Distributed File System) is a distributed storage system used to store large volumes of data across a cluster of machines. It follows a master-slave architecture where:
 - NameNode - Acts as the master and manages the metadata and file system structure.
 - DataNode - Acts as the slave and stores the actual data.HDFS splits large files into smaller blocks and distributes them across the cluster. It ensures fault tolerance by replicating the data across multiple nodes. HDFS supports high-throughput data access and parallel processing making it ideal for handling big data workloads.

2. ◦ Volume: Managing large amounts of data generated at high speed.
- Variety - Handling data in different formats (Structured, Semi-structured, Unstructured).
- Velocity - Processing data at a high rate of incoming flow.
- Scalability - Designing a system that can scale horizontally as data volume increases.
- Veracity: Ensuring the accuracy and consistency of data.

- Security: Protecting sensitive data from unauthorized access and breaches.

Sec - C

- ① Processing in Big Data Infrastructure refers to the methods and frameworks used to analyze and extract insights from large datasets. It includes:
 - Batch Processing: Processing large datasets in bulk (Using Hadoop MapReduce).
 - Real-time Processing: Processing data as it is generated (Using Apache Storm or Spark Streaming).
 - Stream Processing: Handling continuous data streams to provide real-time insights. Efficient processing enables organizations to uncover trends, patterns and actionable insights quickly.

- ② A relational Database management System (RDBMS) is a database system that organizes data into tables (relations) with rows and columns. Key concepts include:
 - Tables - Store data in a structured format.
 - Primary key - A unique identifier for each row.
 - Foreign key - A reference to a primary key in another table to establish relationships.
 - SQL (Structured Query Language).
 - Used to query and manipulate data.

Applications of RDBMS include banking, Customer relationship management (CRM)

and inventory management, where structured data and complex queries are involved.

SQL is widely used in various industries for managing large amounts of data. Some of the key applications of SQL include:

- Relational databases: SQL is the standard language for managing relational databases. It is used for creating, querying, updating, and deleting data stored in tables.
- Business intelligence: SQL is used for extracting, transforming, and loading (ETL) data from various sources into data warehouses for analysis.
- Web development: SQL is often used in conjunction with web technologies like PHP, Python, or Java to interact with databases.
- Cloud computing: SQL is used in cloud environments for managing data stored in cloud databases like Amazon Redshift or Google BigQuery.
- Big data: SQL is used for managing large datasets in distributed systems like Apache Hadoop or Apache Spark.

Another important application of SQL is in the field of data warehousing. Data warehousing involves collecting data from various sources, integrating it, and storing it in a central repository for analytical purposes. SQL is used for defining the schema of data warehouses, performing data extraction, transformation, and loading (ETL) operations, and querying the data for reporting and analysis.

SQL is also used in various other domains such as scientific research, financial modeling, and machine learning. Its ability to handle structured data and perform complex queries makes it a valuable tool for a wide range of applications.