

NUST - Shahzad Khan  
class - TEC(I) T2  
Date - 21/4/25

Roll no - 3025  
sub - DSDBA

\* Theory Assignment + 2 \*

18  
20 Mid

Date: 21/4/25

Q1. What is need of big data analysis? Explain different types of analysis techniques.

=> need of big data analysis

(1) Big Data analysis is essential for

organizations to uncover valuable insights, pattern & trends from large & complex datasets.

(2) Enables them to make informed decisions.

This process involves handling structured, semi-structured & un-structured data, which traditional data processing method cannot efficiently manage.

(3) Types of Data analysis techniques -

i) Descriptive Analysis -

- Focus on summarizing historical data to understand what has happen in past.

- It uses aggregation, data mining & visualization techniques to understand trend pattern & key performance.

- It helps you understand your current situation & make informed decisions.

ii) Predictive Analysis -

- It involves using historical / current data to find patterns & make predictions about future.



Date:

- Accuracy of predictions depends on input variable & type of model used.
- It is useful for forecasting trends outcomes based on historical data.

### ③ Realtime Analysis

- This refers to processing, analyzing data as it's generated.
- Realtime analytics is useful in setting data which are generated quickly.

### ④ Diagnostic analysis

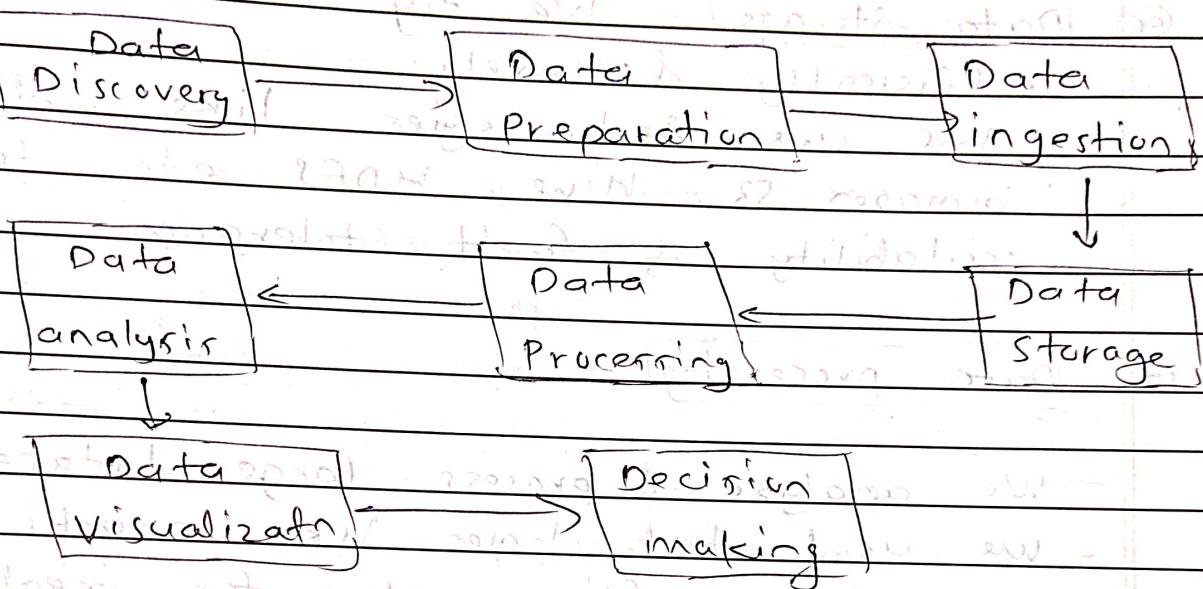
- It typically comes after analysis, taking initial findings involving with certain patterns in data happen.
- It seeks to answer question "why this happen?" by taking more in depth of data to uncover subtle pattern.

Q.2. Explain data analysis life cycle in big data.

### ① Data Discovery

- To understand data sources
- Basically we define what our business

- + Objectives of data mining
- Then we try to explore available data which can be structured or unstructured.
- Assess data volume, velocity, variety, veracity & value (5 v's of big data)



### \* Data Analysis life cycle \*

## ② Data preparation

- In which, we try to make data usable for analysis.
- In which we do data cleaning like handling missing values, null values.
- Data integration from multiple sources are combined, then we perform data transformation & normalization on above data.

## ③ Data Ingestion - In which, we move data into data storage system for processing.



- Here we perform batch / real-time data ingestion tools like apache kafka, flume etc.
- Data is stored in storage-like databases like HDFS or nosql like MongoDB etc.

(4) Data Storage - we try to store data efficiently & securely.

- we use technologies like hadoop, amazon s3, hive, HDFS etc. to ensure scalability & fault tolerance

(5)

Data processing

- We analyze & process large datasets.
- we used technologies like apache spark, mapreduce, flink etc. for realtime & batch processing capabilities.

(6)

Data Analysis

- Extract meaningful insights

- And patterns from processing data.
- Here, we perform different analysis on data like descriptive, predictive, machine learning models to find out different statistical outcome.

(7)

Data Visualization

- It represents insight & decision making

- It uses tools like tableau, powerbi etc. to generate dashboards & interactive reports.

- To guidance business strategy based on insights.

Q. Explain min-max scaling for following data.

$$X = [24, 28, 53, 30, 40, 18, 15, 21]$$

Min-max scaling is a normalization technique used to scale data within fixed range (usually) 0 to 1.

$$\text{Formula - } X_{\text{scaled}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

$$X = [24, 28, 53, 30, 40, 18, 15, 21]$$

$$X_{\text{min}} = 15, X_{\text{max}} = 53$$

X	Calculation formula	Scaled X
24	$\frac{24 - 15}{53 - 15} = \frac{9}{38}$	0.23
28	$\frac{28 - 15}{53 - 15} = \frac{13}{38}$	0.34
53	$\frac{53 - 15}{53 - 15} = \frac{38}{38}$	1.00
30	$\frac{30 - 15}{53 - 15} = \frac{15}{38}$	0.39
40	$\frac{40 - 15}{53 - 15} = \frac{25}{38}$	0.65
18	$\frac{18 - 15}{53 - 15} = \frac{3}{38}$	0.07
15	$\frac{15 - 15}{53 - 15} = \frac{0}{38}$	0.00
21	$\frac{21 - 15}{53 - 15} = \frac{6}{38}$	0.15

$$\text{Scaled value} = X_{\text{scaled}} = [0.23, 0.34, 1.00, 0.39, 0.65, 0.07, 0.00, 0.15]$$



(Q.4.) Data Wrangling - Explain its need & methods.

=>

- Data wrangling is the process of cleaning, transforming & reusing data into format that is more useful for analysis.

### Need for Data Wrangling

- (1) Raw Data is messy (Real world data)
- (2) Different formats (CSV, JSON, XML)
- (3) Improve Data quality (Better insights)
- (4) Data Integration (Wrangling helps in combining multistore to single data)
- (5) Efficiency (Clean data speeds up analysis)

### Methods of Data Wrangling

#### 1. Data Cleaning

- Fixing or removing incorrect, corrupted or missing data.
- Removing duplicates
- Fix spelling mistakes in customer names

#### 2. Data Transformation

- Changing format
- Ex- converting text data to date/time

format, normalizing numeric datasets.

#### (3) Merging

- combining data from multiple datasets into one unified dataset
- Ex - Joining customer & order data based on Customer ID.
- Inner, outer, left/right join are used.

#### (4) Data Reduction

- Removing irrelevant, redundant features
- It helps in improving processing speed, focusing on any imp. features
- techniques - Feature Selection, Dimensionality Reduction (PCA, LDA), Filtering out liars.

