

Page No. _____
Date _____
Question. _____

* Clustering -

cluster \rightarrow a no. of similar things that occur together.

technique in which the data points are arranged in similar groups dynamically without any pre-assignment of groups

partitioning a set of data in a set of meaningful subclasses.

- properties of a cluster -

- all data points in a cluster should be similar to each other
- data points from different clusters should be as different as possible.

Types $\begin{cases} \rightarrow \text{Hard clustering} \rightarrow \text{each data point} \rightarrow \text{only one cluster} \\ \rightarrow \text{soft clustering} \rightarrow \text{each data point} \rightarrow \text{separate cluster} \\ \quad \quad \quad \rightarrow \text{probability or likelihood of data.} \end{cases}$

Distance = $P = (p_1, p_2, \dots)$ $Q = (q_1, q_2, \dots)$

$$d = \sum_{i=1}^k (p_i - q_i)^2$$

Centroid is a point whose coordinates are averages of corresponding coordinates of a given set of points

- Appⁿ →
- Customer Segmentation
 - Image Processing
 - Recommendation Engine
 - Marketing
 - Insurance
 - Seismology
 - Land Use
 - Urban planning
 - Healthcare

*

K-Means-

- heuristic method.
- unsupervised learning algo.
- solve clustering problems.
- groups are unlabelled dataset into different clusters.
- $K \rightarrow$ defines no. of predefined clusters. need to be created.
 - $K=2 \rightarrow$ two clusters
 - $K=3 \rightarrow$ three clusters.
- ~~divides~~ each dataset belongs only one group has similar projects.
- minimizes sum of distances between data point & their corresponding clusters.
- distance calculation \rightarrow Euclidean Distance.

Adv \rightarrow • efficient in computation
• easy to implement.

Disadv \rightarrow • Applicable only when mean is defined.
• need to specify K , no. of clusters, in advance
• Trouble with noisy data & outliers
• not suitable to discover clusters with non-convex shapes.

• Hierarchical clustering -

- Hierarchical cluster analysis or HCA -
- method of cluster analysis
- data points are arranged in hierarchy of clusters.

• Dendrogram -

- diagram representing a tree or hierarchy.

Hierarchical clustering Strategies- (Algorithms)

Agglomerative

- bottom-up approach
- each item in its own cluster
- iteratively clustered are merged together.

Steps → initialization, similarity measurement, merge clusters, creating dendrogram, cut the dendrogram.

- no need of specifying no. of clusters
- expensive → large dataset
- handling noise or outliers is difficult

Divisive

- top-down approach
- all items in one cluster.
- large clusters are successively divided.

Steps → initialize, split cluster, recursive division, creating dendrogram, cut the dendrogram.

- expensive → large dataset
- recursively splitting clusters
- overcome limitation of agglomerative
- determining appropriate no. of clusters & interpreting result



* Time Series Analysis -

- attempts to model underlying structure of observations taken over time.
- Applⁿ → Retail Sales
→ spare parts planning
→ stock trading
- characteristics → Trend
→ Seasonality
→ Cyclic
→ Random

* ARIMA → (Auto Regressive Integrated Moving Average).

- forecasting technique that projects future values of series
- short term forecasting → 40 historical data points.
- Parameters of ARIMA model -
 - p (lag order) - no. of lag observations.
 - d (degree of differencing) → no. of times raw observations are differenced.
 - q (order of moving average) → size of moving avg window.

• Text Analysis -

- also called text analytics
- refers to representation, processing & modelling of textual data to derive useful insights.

• CORPUS → large collection of texts used for NLP.

→ Challenges → • High Dimensionality
• Unstructured Data

• Steps in Text Analysis - (Preprocessing)

- 1) Tokenization
- 2) Lowercasing
- 3) Stop word Removal (noise in text)
- 4) Punctuation Removal
- 5) Lemmatization or Stemming
- 6) Removing HTML tags or special characters.
- 7) Removing Numbers
- 8) Removal of frequent words.
- 9) Removal of rare words.

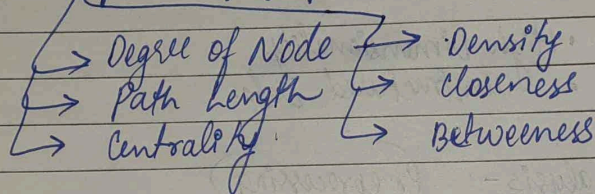
Term Frequency → measures frequency of a word in given document.

$$TF = \frac{\text{count of terms in document (+)}}{\text{total no. of terms in document (d)}}$$

* Social Network Analysis (SNA)

process of investigating social structures through use of networks & graph theory.

- graph theory & mathematical models.
- examines patterns of interactions, connections & dependencies for insights.
- properties of graph



- vertex or Node \rightarrow represent an object in graph.
- connection between \rightarrow nodes \rightarrow edge or a link.

* Need of SNA -

- \rightarrow Identifying Influencers
- \rightarrow Human Resource Management (HRM)
- \rightarrow Contact Tracing
- \rightarrow Identify themes & connections.
- \rightarrow Fraud Detection.

Fundamental Concepts

- \rightarrow Actor
- \rightarrow Relational tie
- \rightarrow Dyad
- \rightarrow Triad
- \rightarrow Subgroup
- \rightarrow Group
- \rightarrow Relation.

* Business Analysis -

involves identifying business needs, analyzing requirements, & finding solutions to organisational problems.

• Business analyst → intermediaries b/w stakeholders & technology terms to ensure that business goals.

- understand business process
- areas for improvement & propose solutions
- documenting business requirements.
- analyzing & modeling business processes.
- facilitating communication & collaboration b/w stakeholders.
- managing change & ensuring smooth implementation of solutions
- defining & validating solution requirements.

1) Leave P-Out Cross Validation

2) Sub-Sampling

- variant of k-fold cross validation → evaluate performance of ML models.
- doesn't divide dataset into k folds
- divides dataset into all possible combinations of leaving p samples & remaining for training
validⁿ set.

- steps →
- value of p (choose)
 - create possible combinations
 - each combinⁿ → train model on training remaining sample & evaluate
 - calculate avg performance
 - more comprehensive evaluation.
 - small datasets or order of sample matters.

2) Sub-Sampling

- down sampling or undersampling
- ml to address class imbalance problems.
- no. of class → significantly uneven.
- Biased model performance.
- Balance class distribution → reduce samples.

- Process steps →
- Identify class → large no. of samples.
 - Randomly select subset of sample.
 - combine selected subset from majority class to minority class.
 - train ml model using balanced dataset.

Q13 - Stemming & Lemmatization.

- techniques used to in NLP to reduce words to their base or root forms.
- normalize words & reduce inflectional or derivational variations.

• Stemming -

- process of removing prefixes & suffixes from words to obtain word's base form.
- resulting stem may not always be a valid word.

original words - fishing, fishes, fished

stemmed words - fish, fish, fish.

In this case, the stem 'fish' is derived by removing the suffixes '-ing', '-es' & '-ed' from original words.

• Lemmatization -

- aims to transform words to their base form or lemma, considering word's meaning & part of speech.
- uses linguistic rules & morphological analysis to achieve this.
- resulting lemma is a valid word represent base form.

original words - walking, walks, walked

stemmed words - walk, walk, walk.

stemming → faster

lemmatization → more accurate results

Q14- Preprocessing Techniques -

essential steps in preparing raw ^{text} data for NLP tasks.

- Transforming & cleaning text to improve its quality & facilitate effective analysis.

1) Tokenization -

↳ process of breaking a text into smaller units called tokens.

- Tokens → words, sentences, or even subword units.
- segmenting text & forms basis for subsequent analysis.

Input Text → "I love to play soccer."

Tokenized Output → ["I", "love", "to", "play", "soccer", "."]

2) Lowercasing -

↳ involves converting all text to lowercase.

It helps in standardizing text & treating words with diff cases (eg. "hello" & "Hello") as same token.

Input Text → "I love NLP"

Lowercased output → "i love nlp"

3) Stop Word Removal -

common words → no significant meaning.

Input Text - "I love to read books & watch movies."

Stop word Removed Output - "love read books watch movies."

4) Punctuation Removal -

5) Lemmatization or Stemming -

6) Removing HTML tags or special characters.

89

Apriori Algorithm

→ association rule mining algorithm → discover frequent items in

Components → support, confidence & Lift a dataset

1) Generating itemsets of size 1 (singletons)

Principle → If an itemset is a frequent itemset then its subsets must also be frequent.

2) Generating Frequent itemsets of size k ($k > 1$): frequent.

3) Iteratively Repeating Step 2-

min support → 3 (occurs at least 3 transactions).

frequent itemsets of size 1 would be {Milk}, {Bread}, {Butter}, & Eggs

→ handles large dataset
→ describe relation b/w items
→ meaningful associations & pattern in transactional data.