

Branch	Date	Sem.	Roll No. / Exam Seat No.	Subject	Student's Signature	Junior Supervisor's Name and Sign
AII	2/4/24	8		Text web SMA		

Question No.	A	B	C	D	E	F	G	H	Total	Total out of (20 / 30 / 40)
1										
2										
3										
4										

Examiners Signature	Student's Sign (After receiving the assessed answer sheet)
---------------------	---

Q-2	<p>a) Text Mining - It is also known as text data Mining or textual data Mining. It is a process of extracting valuable information from text. It involves the discovery by computer of new previously unknown information by automatically extracting information from different written resources</p> <ul style="list-style-type: none"> <li>→ Knowledge Discovery</li> <li>→ efficient information retrieval</li> <li>→ Sentiment Analysis</li> <li>→ Enhancing customer support</li> </ul> <p>b) Tokenization is technique in Text Mining &amp; NLP that involve dividing text into smaller units called token. These tokens can be words, phrases, symbols or other meaningful element that constitute a piece of text. The process of tokenization helps in preparing text data for further analysis or</p>
-----	---

or processing text data for further analysis or processing. Such as parsing, part of speech tagging etc.

c) Social Media Analysis involve collecting data from SM platform & analyzing it to gain insights into various aspects, such as, public opinion, market trends, brand presence, user behaviour. This analysis can be applied across various domain including market, customer service, public relation.

d) feature selection is crucial step in text clustering as it significantly impacts the performance & outcome of clustering process. Text clustering involve grouping a set of texts into cluster grouping a set of texts into clusters that contain similar text.

- 1) Term frequency inverse document frequency
- 2) feature Pounding
- 3) Latent Semantic Analysis
- 4) word embedding
- 5) chi Square

e) Rule based classifier - It is a type of ML algo that makes classification decision based on set of predefined rules. These rules are often created by human experts who understand the domain well or can be automatically generated through algo.

Characteristics:-

- Transparency
- Simplicity
- Domain expertise

f) Text Mining - It is a process of extracting meaningful information & insight from textual data. It involves techniques from linguistics, statistics, & ML to analyze, understand, & interpret large volumes of text. Common tasks in Text Mining include sentiment analysis, topic modeling, NER, & text classification.

g) NER stands for Named Entity Recognition. It is a subtask of NLP that aims to identify & classify named entities within text into pre-defined categories such as names of person, organization, location, date, numerical expression & more.

h) Hidden Markov Model - They are statistical models used to describe sequences of observation events, where each event in the sequence corresponds to an underlying state that is not directly observed. Hence "hidden". HMMs are widely used in various fields, including speech recognition, NLP, bioinformatics & finance.

O=2

(a) N-grams are contiguous sequences of N items (word, character or token) extracted from given text or sequence. They are widely used in NLP & other fields for various tasks such as language modeling, text generation & sentiment analysis.

Type of N-grams

g) Unigram - ( $N=1$ )

- unigram are single word or token in sequence
- Example - "Apple", "orange", "banana"

### 2) Bigram ( $N=2$ )

- Bigram are sequence of 2 adjacent word or token
- ex - "Apple orange"

### 3) Trigram ( $N=3$ )

- Trigram consist of 3 adjacent word or token
- ex - "Apple orange banana"

### 4) N-gram ( $N \geq 3$ )

- N-gram represent sequence of  $N$  adjacent word or token
- ex - "Apple orange, banana, pear"

Q-2

(b) Bayesian N/W - It is probabilistic graphical Model that represent a set of random variable & their conditional dependencies in a directed acyclic graph (DAG)

### Component of Bayesian Network

1) Edge - Directed edge tell Node indicate the dependencies "tell" variable.  
An edge from node A to node B means B is conditionally dependent on A

2) Node - Node represent random variable event, or attribute in the domain of interest.

## Advantage of BN

- 1 Probabilistic Reasoning
- 2 graphical Representation
- 3 Modularity & Scalability
- 4 Causal Inference
- 5 efficient inference
- 6 handling missing data
- 7 Decision Support

Q-3 Recommendation Algorithm: are technique used  
(a) to suggest item or action to users  
based on their preference . behavior . or  
characteristics.

There are mainly 2 Major type

(i) Collaborative filtering: Recommend to the user to a based on the preference of user with similar tastes & identifier . user with similar item preference & suggest item liked by those user . this is mainly the type of collaborative named . as User based collaborative filtering

(ii) Item bases collaborative filtering. This method focuses on the similarity between item . rather than user .  
The similarity bet" item is calculated based on the rating . or interaction they give from user .

2) Content based filtering → The approaches recommend item by analyzing the content or feature of item & profile of user's preference.

Step -

1) Feature extraction - Analyze items to extract features

2) Profile Building - construct a user profile based on features of item

3) Item scoring - for each item, compute score based on similarity profile

4) Recommend item - Recommend item with the highest score that user hasn't interacted with

(Q-3)

(b) Distance based clustering algorithm partition a set of object into cluster, based on similarity measure often using distance between data points.

1) K-mean clustering:- partitioning approach, divides the dataset into k cluster by iteratively updating cluster centroid & assigning data point to the nearest centroid.  
- require specifying the no. of clusters.

- Assume clusters are spherical & equally sized.

- More Scalable with large no of variables & data points.

- ② Hierarchical - It is approach where build a tree of clusters either by starting with individual data points & merging them.
- Does not require specify no of clusters.