

Subject (Write in full) : Text, web & SMA (Regular / KT)

Exam : May 2024 / Nov. 20 _____ Exam Date : 29/ May/ 2024 Q. Paper Code : _____

Department : AI Year : FE/SE/TE/BE/ME/MMS Semester : 8 Scheme : CBSGS/CBCS

Handwritten Solution Prepared by : Amit Aylani

Name of the Subject Cluster : AI ML

Name of the Cluster Mentor / Assessor : Prof. Avinash Shrivastava

1st Assessment :

Question No.	Marks Obtained						Total
	(a)	(b)	(c)	(d)	(e)	(f)	
1							
2							
3							
4							
5							
6							
Total (Out of						marks)	

Signature of Assessor / Cluster Mentor : _____

2nd Assessment :

Question No.	Marks Obtained						Total
	(a)	(b)	(c)	(d)	(e)	(f)	
1							
2							
3							
4							
5							
6							
Total (Out of						marks)	

Signature of Assessor / Cluster Mentor : _____

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
	(Q. 1)	<p>Named entity recognition (NER) is a process in NLP that identifies and classifies key elements in text into predefined categories such as name of people, organization, location, date & many more.</p> <p>for example - In the sentence "Apple Inc was founded by Steve Jobs, Steve Wozniak, a NER System would identify "Apple Inc" as an organization, "Steve Jobs" "Steve Wozniak" & Ronald Wayne" as person, "Cupertino, California" as a location & April 1, 1976 as a date.</p> <p>One effective approach to NER is using Conditional Random fields (CRFs) which are statistical Model that predict the sequence of label (entity categories) for a sequence of word. CRFs consider the context of each word & use feature like word itself - part of speech tags, orthographic feature & neighbouring word to accurately predict entities.</p> <ul style="list-style-type: none"> > Feature extraction > Training the CRF Model > Interface

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
		Advantage :-
		<ul style="list-style-type: none"> ➢ CRFs Consider the Context of each word by taking into account neighbouring words & their feature, making them well suited for Sequential labeling task.
		<ul style="list-style-type: none"> ➢ Unlike Hidden Markov Model, CRFs don't make the independence assumption i.e. "the absence" feature allowing for more flexibility & accuracy in Modeling
		<p>Q-L (b) Probabilistic Document clustering → It is a technique used to group documents into clusters based on probabilistic Model that capture the underlying data distribution. This method assumes that documents are generated by a mixture of probabilistic distribution, it aims to find the parameters of these distributions that best explain the data. Common Model used in probabilistic document clustering include Latent Dirichlet Allocation (LDA) & GMM</p>
		<ul style="list-style-type: none"> ➢ Assignment of Document :

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	
Space for Marks	Question No.	START WRITING HERE	
		<p>1) Distance based clustering - In method like K Mean or hierarchical clustering. Each document is assigned to a single document is assigned to a single cluster based on a distance metric. (e.g. Euclidean distance or cosine similarity). Document are grouped by minimizing the within cluster variance.</p> <p>2) Modeling data -</p> <ul style="list-style-type: none"> > They methods rely on predefined notion of distance or similarity. They do not explicitly model the underlying distribution of the data. > They method use probabilistic model to represent the data. They capture the statistical properties of relationship within the data, providing a generative process of document. <p>3) Handling overlap -</p> <p>Typically assume that clusters are distinct & non-overlapping. Each document belongs to exactly one cluster.</p>	

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
	Q1>	Complexity of Computation - 10
		→ Often involve simpler computation, like distance calculation & centroid update. Then mean method can be faster & easier to implement.
Q=2	(a)	Rule Based classifier:- Rule Based classifier use a set of predefined rules to assign labels to text document. These rules are often constructed based on the presence or absence of specific keyword, phrases or pattern within text. Each rule typically consists of condition & an action.
		Ex - Rule-1 : If the email contain the phrase "Congratulations you have won" then label it as "spam".
		Rule-2 : If the email contain the word "free" more than three times then label it as "spam".
		Rule-3 : If the email contain the phrase "unsubscribe" then label it as "spam".
		Rule-4 : If email contain the word "urgent" AND it is from a known contact

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	
Space for Marks	Question No.	START WRITING HERE	
		<p>then label it as "not spam".</p> <p>> Strength of Rule Based classifier:</p> <ul style="list-style-type: none"> > Simplicity & Interpretability:- The rules are straight forward & easy to understand. Each decision can be traced back to a specific rule making the classification process transparent. > Customization - rule can be tailored to specific needs & updated easily as new pattern emerge. > Efficiency - for small to medium sized dataset, rule based system can be very efficient since they do not require extensive training. <p>Q-2 → Distance based clustering Algo.</p> <p>(b)</p> <p>> K-Means:- K-mean is a popular distance based clustering algo that partitions data into K clusters. It aims to Min. the variance within each cluster by assigning each document to the nearest cluster centroid.</p>	

Total Marks of Question no.		Examiner
		Moderator
		Re-Assessor

Space for Marks	Question No.	START WRITING HERE
		Advantage - <ul style="list-style-type: none"> → Simplicity → effectiveness → Scalability
		Disadvantage - <ol style="list-style-type: none"> 1) fixed no of clusters. 2) Sensitivity to initialization 3) Assumption of spherical clusters
	2)	Hierarchical clustering:- <p style="text-align: right;">Hierarchical clustering</p> <p>build a tree of clusters by either agglomerative bottomup or top down approaches</p>
		Advantage - <ul style="list-style-type: none"> - No need to specify no of clusters - flexibility - Dendrogram visualization
		Disadvantage - <ul style="list-style-type: none"> - Computationally intensive - not suitable for very large dataset - Sensitivity to noise & outliers.
		> Comparison:- (1) Kmean is more efficient & Scalable for large dataset but require the no of clusters to be predefined

Total Marks of Question no.		Examiner
		Moderator
		Re-Assessor

START WRITING HERE

Space for Marks	Question No.

Answers Spherical cluster.

> Hierarchical clustering - It provide more flexibility & does not require the no. of clusters to be specific in advance but it is computationally intensive & less suitable for very large dataset.

Q=3

(a) web spam employs various technique ..

1. Cloaking :- Serving different content to Search engine & user. A web page might show high quality, keyword rich content to search engine while displaying irrelevant or spammy content to user.

2. Keyword stuffing :- Overloading a web page with keyword to manipulate search engine ranking.

3. Hidden Text & link :- Using technique to hide text & links from user while making them visible to search engine.

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
		→ Challenges Associated with detecting & Combating Spam
		1) Evolving technique :- Spammer continually develop new tactic to evade detection
		2) false positive & negative : legitimate Content may be mistakenly classified as spam or spam may go undetected.
		3) Resource Intensiveness :- effective spam detection often require significant computational resources & sophisticated algorithm.
	(b)	Challenges in Multilingual Sentiment Analysis
	1)	Language Variability :- - Different languages have unique syntax, grammar & idiomatic expression that affect how sentiment is conveyed.
	2)	Resource availability - There is often a lack of annotated dataset & linguistic resource (like Sentiment lexicon) for many languages

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	
Space for Marks	Question No.	START WRITING HERE	
		<p>3) Translating Issue:</p> <ul style="list-style-type: none"> - Translating Issue for Sentiment analysis can introduce error & may not preserve the original sentiment context. <p>> Sentiment can be expressed differently across cultur might have a neutral or negative connotation in another.</p>	
		<p>4) Addressing challenges in multilingual Sentiment analysis</p> <ul style="list-style-type: none"> > Multilingual Pre trained Model > Cross Lingual Embedding > Transfer learning > Data Augmentation & crowd sourcing > cultural Adaptation > Translation aware Model 	

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
	Q.4 (a)	<p>Session & visitor analysis in web Mining involve examining the interaction & behaviour of user on a website over a specific period. A session typically represent a single visit by a user, capturing all interaction from time they enter the website until they leave. visitors analysis consider both the session level & broader behaviour of user, including return visit & overall engagement patterns</p> <p><u>Important :-</u></p> <ul style="list-style-type: none"> > Understanding user behaviour :- Analyzing session help in understanding how user interact with website, what page they visit & how long they stay on each page. > Improving User experience :- by studying visitor behaviour journey can identify pain points in the user journey such as drop off rates on specific page or during particular process. > Personalization :- Insight from session & visitor analysis can be used to tailor content & offer to individual user, enhancing personalization.

Total Marks of Question no.	
Moderator	
Re-Assessor	

START WRITING HERE

Space for Marks	Question No.

* Marketing & Conversion optimization
 → understanding the path to
 reach conversion point ; making
 purchase sign up for newsletter

Insigned Grained

- > Traffic sources & Referrals
- > User engagement & Content effectiveness
- > Customer Journey mapping
- > User demographic & Segmentation

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
Q4	(b)	Analysis of Sequence & Navigational pattern in web mining
		> Sequential & Navigation pattern analysis In web mining focus on understanding the order & structure of user action as they navigate through a website. This analysis aims to uncover common path, sequence & pattern in user navigation.
		> Sequence Mining Algorithm.
	1)	Apriori Algorithm:- It used for Market analysis, it can be adapted to identify frequent navigation sequence by finding frequent item sets. & then generating association rules.
	2)	prefixSpan Algo - It focus on finding frequent sequential pattern without generating candidate sequence making it more efficient for large datasets
	3)	GSP - Extends Apriori by considering time constraint & allowing for gaps let "sequences.

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	
Space for Marks	Question No.	START WRITING HERE	
		Application:-	
		<p>These algorithms can be applied to web logs to identify frequent sequence of pages that user visit helping to understand common navigation path & user behavior.</p>	
		<ul style="list-style-type: none"> > Navigational pattern path. 1) Path analysis 2) Clickstream analysis 	
		<p>⇒ Insights gained</p>	
		<ul style="list-style-type: none"> > User profile & interest > Optimization & Navigation structure > Improving user flow > Predictive Analytics. 	

Total Marks of Question no.	Examiner	
	Moderator	
	Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
(Q.5)	(a)	<p>Behaviour Analysis in SM Minin :-</p> <p>behaviour analytics in sm minim involves the study & interpretation of user action & interaction on sm platform. They include tracking metric like, share, comment, retweet, follow click, & time spent on content. By analyzing their behaviour organization can gain insights into user engagement profile & emerging trends.</p> <p>1) User engagement:- Analyzing metric such as like, share, comment & retweet helps measure how user interact with content.</p> <p>2) User preference:- By tracking which type of content gain the most engagement, we can infer user preference.</p> <p>3) Trend identification - Monitoring the volume & velocity of specific keywords hashtag & topic over time.</p> <p>4) Audience & Segmentation:- Analyzing demographic data & behaviour pattern to segment user to distinct group.</p>

Total Marks of Question no.		Examiner	
		Moderator	
		Re-Assessor	

Space for Marks	Question No.	START WRITING HERE
(Q-5)	(a)	Behaviour Analysis in SM Miniy:- behavioural analysis in sm miniy involves the study & interpretation of user reaction & interaction on sm platform. They include tracking metric like, share, comment, retweet, follow click, & time spent on content. By analyzing their behaviour organization can gain insights into user engagement profile & emerging trends. 1) User engagement:- Analyzing metric such as like, share, comment & retweet helps measure how user interact with content. 2) User preference:- By tracking which type of content seem the most engaged user can infer user preference. 3) Trend identification - Monitoring the volume & velocity of specific keywords chartage & topic over time. 4) Audience & Segmentation:- Analyzing demographic data & behaviour pattern to segment user to distinct group.

Total Marks of Question no.		Examiner
		Moderator
		Re-Assessor

Space for Marks	Question No.	START WRITING HERE
		5) Sentiment Analysis - Analyzing the tone & sentiment of user comments review & post
		6) Influence identification - Tracking the engage of each individual user to identify influential figure within a week
	(b)	<p>Collaborative filtering -</p> <p>It is a recommendation approach can be based on the assumption that user with similar behaviour & preference with like similar item</p> <p>> User based Collaborative filtering : Recommended item by finding similar user based on their rating & behaviour</p> <p>> Item based Collaborative filtering :- - It recommend item by finding similar item based on user interaction with them</p>

Total Marks of Question no.		Examiner										
		Moderator										
		Re-Assessor										
Space for Marks	Question No.	START WRITING HERE										
		<p>> Content Based filtering It recommends item based on characteristics of item & user past interaction.</p> <p>> Analyzes the content of item & matches them with user (profile derived from their past behaviour)</p> <p>> Difference b/w Collaborative & Content based filtering</p> <table border="1" style="width: 100%; border-collapse: collapse;"> <tr> <td style="padding: 5px; vertical-align: top;">1) Data dependency</td> <td style="padding: 5px; vertical-align: top;">Dependent on user item interaction & rating</td> <td style="padding: 5px; vertical-align: top;">Depend on item feature & user profile.</td> </tr> <tr> <td style="padding: 5px; vertical-align: top;">2) personalization</td> <td style="padding: 5px; vertical-align: top;">provide recommendation based on similarity.</td> <td style="padding: 5px; vertical-align: top;">provide recommendation based on individual user profile</td> </tr> <tr> <td style="padding: 5px; vertical-align: top;">3) cold start problem</td> <td style="padding: 5px; vertical-align: top;">Suffer from cold start problem for new user & new item</td> <td style="padding: 5px; vertical-align: top;">> Left affected by cold start problem for new item as recommendation</td> </tr> </table>		1) Data dependency	Dependent on user item interaction & rating	Depend on item feature & user profile.	2) personalization	provide recommendation based on similarity.	provide recommendation based on individual user profile	3) cold start problem	Suffer from cold start problem for new user & new item	> Left affected by cold start problem for new item as recommendation
1) Data dependency	Dependent on user item interaction & rating	Depend on item feature & user profile.										
2) personalization	provide recommendation based on similarity.	provide recommendation based on individual user profile										
3) cold start problem	Suffer from cold start problem for new user & new item	> Left affected by cold start problem for new item as recommendation										

