

CBCS SCHEME

USN

--	--	--	--	--	--	--	--	--	--

18CS743

Seventh Semester B.E. Degree Examination, Dec.2023/Jan.2024 Natural Language Processing

Time: 3 hrs.

Max. Marks: 100

Note: Answer any FIVE full questions, choosing ONE full question from each module.

Module-1

- 1 a. What is NLP? Explain the challenges of NLP. (08 Marks)
- b. Explain different levels of language processing. (08 Marks)
- c. Differentiate between the rationalist and empiricist approaches to NLP. (04 Marks)

OR

- 2 a. Explain n-gram model. How to handle data sparseness problem in n-gram model. (10 Marks)
- b. Represent the following sentence in GB (S-structure and d-structure) (06 Marks)
- c. With example, explain the important features of Indian Languages. (04 Marks)

Module-2

- 3 a. What is Morphology? Explain two step morphological parser. (06 Marks)
- b. Write and explain minimum edit distance algorithm. Compute minimum edit distance between tutor and tumour. (10 Marks)
- c. What is POS (Part of Speech) tagging? List POS tagging methods. (04 Marks)

OR

- 4 a. Discuss the disadvantages of the logic top-down parser with the help of an appropriate example. (10 Marks)
- b. Write a note on :
 - i) CFG (Context-Free Grammar) for natural language
 - ii) Lexicalization (10 Marks)

Module-3

- 5 a. Explain the shortest path hypothesis with example. (10 Marks)
- b. Explain how to capture relation pattern with a string Kernel. (10 Marks)

OR

- 6 a. With diagram explain learning framework architecture. (10 Marks)
- b. Explain the strategies used in active learning approach. (10 Marks)

Module-4

- 7 a. Explain with neat diagram evolutionary model for KDT (Knowledge Discovery from Text) (10 Marks)
- b. Define :
 - i) Cohesion
 - ii) Interestingness
 - iii) Coherence
 - iv) Coverage
 - v) Plausibility of origin. (10 Marks)

OR

- 8 a. Explain SVM learning method in sequential model estimation. (10 Marks)
b. Write a note on various approaches to analyzing texts. (10 Marks)

Module-5

- 9 a. What are the benefits of eliminating stop words? Give examples in which stop word elimination may be harmful. (06 Marks)
b. What is IR (information system)? Explain design features of IR with neat diagram. (10 Marks)
c. State and explain Zipf's law. (04 Marks)

OR

10 Write a short note on :

- i) POS tagger
- ii) WORDNET
- iii) FRAMENET
- iv) STEMMERS

(20 Marks)



Visvesvaraya Technological University
Belagavi, Karnataka - 590 018.

Scheme & Solutions

Signature of Scrutinizer

Subject Title : Natural Language Processing

Subject Code : 18CS743

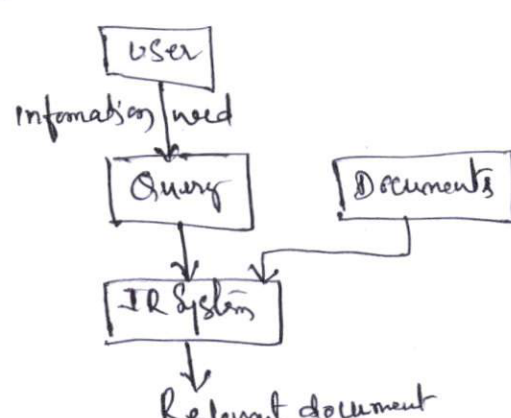
Question Number	Solution	Marks Allocated
1 a.	<p>Definition: concerned with the development of computational models of aspects of human language processing</p> <p>challenges in</p> <ul style="list-style-type: none"> * representation & interpretation. + identifying semantics. * frequency of words * Idioms, metaphor & ellipses * Quantifier Scoping * Ambiguity 	2 Marks 6 Marks
1 b.	<p>Levels of language processing:</p> <ul style="list-style-type: none"> * lexical analysis. * Syntactic analysis. * Semantic analysis. * discourse analysis * pragmatic analysis. (with explanation) 	8 Marks
1 c.	<p><u>Rationalist approach</u>: which assumes the existence of some language faculty in the human brain. (2M)</p> <p><u>Empiricist</u>: do not believe in existence of a language faculty. Instead they believe in the existence of some general organization principles. (2M)</p>	4 Marks
2 a.	<p>Explanation of n-gram (Probability estimation)</p> $P(w_i/b_i) = P(w_i/w_{i-n+1} \dots w_{i-1})$ $P(w_i/w_{i-n+1}, \dots w_{i-1}) = \frac{C(w_{i-n+1}, \dots w_{i-1}, w_i)}{C(w_{i-n+1}, \dots w_{i-1})}$ <p>Techniques to handle data sparseness problem</p> <ul style="list-style-type: none"> * Add-one Smoothing * Good - Turing Smoothing * Caching Technique 	5 Marks 5 Marks

Question Number	Solution	Marks Allocated
2 b	<p>S- structure representation :</p> <pre> S / \ NP INFL VP / \ Mukhsh past / \ BE VP / \ Be V NP e Killed </pre> <p>D- structure</p> <pre> S / \ NP INFL VP / \ V VP / \ V NP Mukhsh. Killed </pre>	3 Marks 3 Marks
2c	<ul style="list-style-type: none"> * Morphologically rich languages . * flexibility allow word groups in any order * Nouns are followed by post positions * Verbs are formed differently in Indo-Aryan & Dravidian languages . <p>Explanation with example</p>	4 Marks
3 a	<p>Morphology: It studies word structure & the formation of words from smaller units (morphemes).</p> <pre> Surface form → [Step 1: Split word into possible morphemes] → [Intermediate form] → [Step 2: Map morphemes to stem & morphological features] → Lexical form </pre> <p>Explanation</p>	2 marks 2 marks 2 marks

Question Number	Solution	Marks Allocated																																																								
3 b	<p>Algorithm: Input : two strings X and Y output : The minimum distance b/w X & Y $m \leftarrow \text{length}(X)$ $n \leftarrow \text{length}(Y)$ for $i=0$ to m do $\text{dist}[i,0] \leftarrow i$ for $j=0$ to n do $\text{dist}[0,j] \leftarrow j$ for $i=0$ to m do for $j=0$ to n do $\text{dist}[i,j] = \min \{ \text{dist}[i-1,j] + \text{insert_cost},$ $\text{dist}[i-1,j-1] + \text{Subst_Cost}(X_i, Y_j),$ $\text{dist}[i,j-1] + \text{delet_cost} \}$</p> <p>Explanation —————</p> <table border="1" data-bbox="368 829 975 1256"> <tr><td></td><td>*</td><td>t</td><td>u</td><td>m</td><td>o</td><td>4</td><td>r</td></tr> <tr><td>x</td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td></tr> <tr><td>t</td><td>1</td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr> <tr><td>u</td><td>2</td><td>1</td><td>0</td><td>1</td><td>2</td><td>3</td><td>4</td></tr> <tr><td>t</td><td>3</td><td>2</td><td>1</td><td>1</td><td>2</td><td>3</td><td>4</td></tr> <tr><td>o</td><td>4</td><td>3</td><td>2</td><td>2</td><td>1</td><td>2</td><td>3</td></tr> <tr><td>r</td><td>5</td><td>4</td><td>3</td><td>3</td><td>2</td><td>2</td><td>2</td></tr> </table> <p>Computing Minimum distance.</p>		*	t	u	m	o	4	r	x	0	1	2	3	4	5	6	t	1	0	1	2	3	4	5	u	2	1	0	1	2	3	4	t	3	2	1	1	2	3	4	o	4	3	2	2	1	2	3	r	5	4	3	3	2	2	2	<p>4 Marks</p> <p>3 Marks</p> <p>3 Marks</p>
	*	t	u	m	o	4	r																																																			
x	0	1	2	3	4	5	6																																																			
t	1	0	1	2	3	4	5																																																			
u	2	1	0	1	2	3	4																																																			
t	3	2	1	1	2	3	4																																																			
o	4	3	2	2	1	2	3																																																			
r	5	4	3	3	2	2	2																																																			
3 c	<p>POS tagging: process of assigning a part of Speech Such as a noun, verb, pronoun, preposition, adverb & adjective to each word in a sentence.</p> <p>POS tagging methods: * Rule-based, * Stochastic (data driven), * Hybrid</p>	<p>2 Marks</p> <p>2 Marks</p>																																																								
4 a	<p>Disadvantages: * left recursion forms stuck in an infinite loop * Structural ambiguity, * Coordination ambiguity, * local ambiguity, * repeated parsing</p> <p>Explanation —————</p> <p>Example —————</p>	<p>3 Marks</p> <p>4 Marks</p> <p>3 Marks</p>																																																								
4 b	<p>Explanation of (i) CFG. _____ (ii) Lexicalization _____</p>	<p>5 Marks</p> <p>5 Marks</p>																																																								

Question Number	Solution	Marks Allocated
5 a	<p>If e_1 & e_2 are two entities mentioned in the same sentence such that they are observed to be in a relation R, then the contribution of the sentence dependency graph to establishing the relation $R(e_1, e_2)$ is almost exclusively concentrated in the <u>shortest path</u> between e_1 & e_2 in the undirected version of the dependency graph.</p> <p>Explanation of the above concept example —</p>	<p>4 marks 3 marks 3 marks</p>
5 b.	<p>Explanation of along with the patterns:</p> <ul style="list-style-type: none"> * [FB] Force-Between . * [B] Between * [BA] Between-After * [M] Modifier 	<p>6+4 marks</p>
6 a	<p>Activity learning</p> <p>Explanation</p>	<p>5 marks 5 marks</p>
6 b.	<p>Strategies:</p> <ol style="list-style-type: none"> (a) Divide the corpus in clusters of sentences with the same target verb. If a cluster has fewer sentences than a given threshold, group sentences with verbs evoking the same frame into the same cluster. (b) within each cluster, group the sentences with the same parse subtree together. (c) Select sentences from the largest groups of the largest clusters and present them to the user for annotation. (d) Bootstrap initialization: Apply the labels assigned by the users to groups of sentences with the same parse subtree. 	

Question Number	Solution	Marks Allocated
	<p>(e) Train all the classifiers of the Committee on the labeled instances; apply each trained classifier to the Unlabelled sentences.</p> <p>(f) Get a pool of instances where the classifiers of the Committee disagree & present to the user the instances belonging to sentences from the next largest clusters not yet manually labelled.</p> <p>(g) Repeat steps d-f a few times until a desired accuracy of classification is achieved.</p>	<p>7 Marks</p>
7 a	<p style="text-align: right;">Explanation</p> <p style="text-align: center;">Preprocessing & training</p> <p style="text-align: center;">Knowledge discovery.</p> <p style="text-align: right;">feg — 5 Marks Explanation — 5 Marks</p>	<p>3 Marks.</p>
7 b	<p>(i) cohesion (H) = $\sum_{r_i, p_i \in H} \frac{Pos(P_i/r_i)}{ H }$ $P_i \rightarrow$ predicate $r_i \rightarrow$ rhetorical role</p> <p>(ii) interestingness (H) = \langle semantic dissimilarity b/w Antecedent and Consequent \rangle.</p> <p>(iii) Coherence (H) = $\frac{\sum_{i=1}^{ H -1} SemSim(P_i(A_i), P_{i+1}(A_{i+1}))}{(H -1)}$ SemSim \rightarrow semantic similarity, (H -1) \rightarrow adjacent pair</p> <p>(iv) Coverage (H) = $\frac{ Rules\ Covered(H) }{ Ruleset }$</p> <p>(v) plausibility (H) = $\begin{cases} S_p & \text{if H wa created from a Swanson's Crossover} \\ 0 & \text{if H is the original population or is a result of another operations.} \end{cases}$</p>	<p>2x5 = 10 Marks</p>

Question Number	Solution	Marks Allocated
8 a	Explanation of how SVM (Support Vector machine) learning Method used in Sequence model estimation.	10 Marks
8 b	Explanation of various approaches to analyzing texts * Bag-of-words-based approach. ——— * Highlevel representation approach ———	5 Marks 5 Marks.
9 a	Benefits of eliminating stopwords ——— drawback of eliminating stopwords with ex ———	2 Marks 4 Marks
9 b	Definition of IR: deals with the organization, storage, retrieval, and evaluation of information relevant to a user's query  <pre> graph TD User[User] -- Information need --> Query[Query] Documents[Documents] --> IRSystem[IR System] Query --> IRSystem IRSystem --> Relevant[Relevant document] </pre> Explanation of design features: indexing, stopword elimination, stemming	2 Marks 2 Marks 6 Marks.
9 c	Zipf's law statement and explanation 2+2	4 Marks.
10	(i) POS tagger — (5 Marks) (ii) WORDNET — (3 Marks) (iii) FRAMENET — (5 Marks) (iv) STEMMER — (5 Marks)	20 Marks