

USN										
-----	--	--	--	--	--	--	--	--	--	--

Sub:	Natural Language Processing				Sub Code:	18CS743	Branch:	
Date:	27/12/22	Duration:	90 mins	Max Marks:	50	Version/ Sem / Sec:	C/VII/A, B, C	

Answer any FIVE FULL Questions

MARKS

1.	<p>Explain LSA feedback system.</p> <ul style="list-style-type: none"> ➤ LSA uses statistical computations to extract and represent the meaning of words. ➤ Meanings are represented in terms of their similarity to other words. ➤ LSA begins by finding: <ul style="list-style-type: none"> ➤ the frequency of terms used and ➤ the number of co-occurrences in each document throughout the corpus and ➤ then uses a mathematical transformation to find deeper meanings and relations among words. ➤ LSA provides benefit in finding similarity between 2 documents. ➤ The method, does not take into account word order; ➤ Very short documents may not be able to receive the full benefit of LSA. ➤ To construct an LSA corpus matrix, a collection of documents are selected. ➤ A term-document frequency (TDF) matrix X is created for those terms that appear in 2 or more documents. ➤ The row entities correspond to the words or terms (hence the W) and the column entities correspond to the documents (hence the D). ➤ The matrix is then analyzed using Singular Value Decomposition (SVD) ➤ TDF matrix X is decomposed into the product of 3 other matrices: <ul style="list-style-type: none"> ➤ vectors of derived orthogonal factor values of the original row entities W, ➤ vectors of derived orthogonal factor values of the original column entities D, and ➤ scaling values (which is a diagonal matrix) S. ➤ The product of these 3 matrices is the original TDF matrix. $\{X\} = \{W\}\{S\}\{D\}$	[10]
----	--	------

	<ul style="list-style-type: none"> ➤ The similarity of terms is computed by taking the cosine of the corresponding term vectors. ➤ A term vector is the row entity of that term in the matrix W. ✓ The similarity between two documents (i.e., the cosine between the two document vectors) is computed as $Sim(D1, D2) = \frac{\sum_{i=1}^d (D1_i \times D2_i)}{\sum_{i=1}^d (D1_i)^2 \times \sum_{i=1}^d (D2_i)^2} \quad (6.3)$ ➤ The 4 benchmarks include: <ol style="list-style-type: none"> 1) the words in the title of the passage (“title”) 2) the words in the sentence (“current sentence”) 3) words that appear in prior sentences in the text that are causally related to the sentence (“prior text”) and 4) words that were used by two or more subjects who explained the sentence during experiments (“world knowledge”). 	
2.	<p>How to evaluate self-explanations in iSTART</p> <p>iSTART (Interactive Strategy Trainer for Active Reading and Thinking) is a webbased, automated tutor designed to help students become better readers via multimedia technologies.</p> <ul style="list-style-type: none"> ➤ It provides students with a program of self-explanation and reading strategy training, called Self-Explanation Reading Training, or SERT. ➤ The reading strategies include: <ul style="list-style-type: none"> ➤ Comprehension monitoring: being aware of one’s understanding of the text ➤ paraphrasing: or restating the text in different words ➤ elaboration: using prior knowledge or experiences to understand the text ➤ predictions: predicting what the text will say next and ➤ bridging, understanding the relation between separate sentences of the text. ➤ The overall process is called “self-explanation” because the reader is encouraged to explain difficult text to him- or herself. ➤ iSTART consists of 3 modules: Introduction, Demonstration, and Practice 	[10]

	<ul style="list-style-type: none"> ➤ The system evaluates each self-explanation and then provides appropriate feedback to the student. ➤ If the explanation is irrelevant or too short, the student is required to add more information. ➤ iSTART was initially proposed as using Latent Semantic Analysis (LSA). ➤ iSTART used simple word matching algorithms. ➤ Combination of word-matching and LSA provided better results than either separately. ➤ The current evaluation system predicts the score that a human gives on a 4-point scale, where: <ul style="list-style-type: none"> ➤ 0-represents an evaluation of the explanation as irrelevant or too short; ➤ 1-minimally acceptable; ➤ 2-better but including primarily the local textual context; and ➤ 3-oriented to a more global comprehension. ➤ To improve the effectiveness of algorithms, incorporate Topic Models (TM) in conjunction with LSA. ➤ The algorithms are constrained by 2 major requirements: <ul style="list-style-type: none"> ➤ speedy response times and ➤ speedy introduction of new texts. 	
3.	<p>(a) Explain in detail the high-level representation approaches in text mining (5)</p> <ul style="list-style-type: none"> ➤ Another main stream in KDT involves using more structured or higher-level representations to perform deeper analysis so to discover more sophisticated novel / interesting knowledge. Although in general, the different approaches have been concerned with either performing exploratory analysis for hypothesis formation, or finding new connections/relations between previously analysed natural language knowledge, it has also involved using term-level knowledge for other purposes than just statistical analysis. ➤ Some early research by Swanson used an augmented low-level representation (the words in the titles) and exploratory data analysis to discover hidden connections leading to very promising and interesting results in terms of answering questions for which the answer was not currently known. He showed how chains of causal implication within the medical literature can lead to hypotheses for causes of rare diseases, some of which have received scientific supporting evidence. ➤ Other approaches using Information Extraction (IE) which inherited some of Swanson’s ideas to derive new patterns from a combination of text fragments, have also been successful. Essentially, IE is a Natural-Language (NL) technology which analyses an input NL document in a shallow way by using defined patterns along with mechanisms to resolve implicit discourse-level information to match important information from the texts. As a result, an IE task produces an intermediate representation called “templates” in which information relevant has 	[10]

been recognised, for example: names, events, entities, etc., or high-level linguistic entities: noun phrases, etc.

- This technique is meant to be useful as an automated or semi-automated aid for lexicographers and builders of domain-dependent knowledge bases. Also, it does not require an additional knowledge base or specific interpretation procedures in order to propose new instances of WordNet relations [9]. Once the basic relations are obtained, they are used to find common links with other “similar” concepts in WordNet [9] and so to discover new semantic links [18].
- However, there are tasks which need to be performed by hand such as deciding on a lexical relation that is of interest (i.e., hyponym) and a list of word pairs from WordNet this relation is known to hold between. One of the main advantages of this method is its low cost for augmenting the structure of WordNet and its simplicity of relations.
- However, it also has some drawbacks including its dependence on the structure of a general-purpose ontology which prevents it from reasoning about specific terminology/concepts, the restricted set of defined semantic relations, its dependence on WordNet’s terms (i.e., only terms present in WordNet can be related and any novel domain-specific term will be missed), the kind of inference enabled etc

(b) Explain SVM learning method in Sequence Model Estimation. (5)

- The learning method adopted here for estimating the sequence model is a Support Vector Machine[10] (SVM). It is commonly known that Support Vector Machines are well suited for text applications given a small number of training examples.
- This is an important aspect for the commercial use of the system, since the process of gathering, preparing, and cleaning up training examples is time consuming and expensive. Support Vector Machines solve a binary classification problem.
- The SVM score associated with an instance of the considered events is its signed distance to the separating hyperplane in units of the SVM margin. In order to solve multiclass problems, a series of Support Vector Machines have to be trained, e.g., in the case of a one-vs-all training schema, the number of SVMs trained is given by the number of classes.
- The scores between these different machines are not directly comparable and the scores must be calibrated such that at least for a given classification instance the scores are on an equal scale. In this application, the scores not only must be comparable between classes for a given classification instance (page), but also between different classification instances (pages), i.e., the SVM scores must be mapped to probabilities.
- Platt uses SVM scores that are calibrated to class membership probabilities by adopting the interpretation of the score being proportional to the logarithmic ratio of class membership probability. He determines the class membership probability as a function of the SVM score by fitting a sigmoid function to the empirically observed class membership probabilities as a function of the SVM score.
- The fit parameters are the slope of the sigmoid function and/or a translational offset. The latter parameter, given the interpretation of the SVM scores discussed above, is the logarithmic ratio of the class prior probabilities. The method used here fixes the translational offset and only fits the slope parameter.

	<ul style="list-style-type: none"> ➤ In addition, the Support Vector Machines are trained using cost factors for the positive as well as for the negative class and optimize the two costs independently. ➤ Empirical studies performed by the authors showed that cost factor optimization in conjunction with fitting the slope parameter of the mapping function from SVM scores to probabilities yields superior probability estimates than fitting the slope and the translational offset without cost factor optimization, fitting the slope and the translational offset with cost factor optimization, and fitting the slope only 	
4	<p>Explain classical and non-classical information retrieval models with suitable examples</p> <p>classical information retrieval models</p> <ul style="list-style-type: none"> ➤ These are based on mathematical knowledge that is easily recognized and well understood. ➤ They are simple, efficient and easy to implement. ➤ The 3 classical IR models are: <ol style="list-style-type: none"> 1. Boolean model 2. Probabilistic model 3. Vector model ➤ Non-classical IR models are based on principles other than similarity, probability, Boolean operations, etc., on which classical IR models are based. ➤ Non-classical IR models are: <ol style="list-style-type: none"> 1. Information logic model 2. Situation theory model 3. Interaction model. 	[10]

5	<p>Write short notes on:</p> <p>(i) Word Net</p> <ul style="list-style-type: none"> ● WordNet is a large lexical database for the English language. ● Inspired by psycholinguistic theories, it was developed. ● WordNet consists of 3 databases 	[10]	CO4	L2
---	---	------	-----	----

- One for nouns
- One for verbs
- One for both adjectives and adverbs
- ✓ Information is organized into sets of synonymous words called synsets, each representing 1 base concept.
- ✓ The synsets are linked to each other by means of lexical and semantic relations.
- ✓ Lexical relations occur between word-forms (senses).
- ✓ semantic relations occur between word meanings.
- ✓ These relations include synonymy, hypernymy / hyponymy, antonymy, meronymy / holonymy, troponymy, etc.
- ✓ If a word appears in more than 1 synset and in more than 1 part-of-speech.
- ✓ the meaning of a word is called sense.
- ✓ WordNet lists all senses of a word.
- ✓ Each sense belonging to a different synset.
- ✓ WordNet's sense-entries consist of a set synonyms and a gloss.
- ✓ A gloss consists of a dictionary-style definition and examples demonstrating the use of a synset in a sentence, as shown in the figure below.
- ✓ The figure shows the entries for the word 'read'. Read has 1 sense as a noun and 11 senses as a verb.
- ✓ Glosses help differentiate meanings.

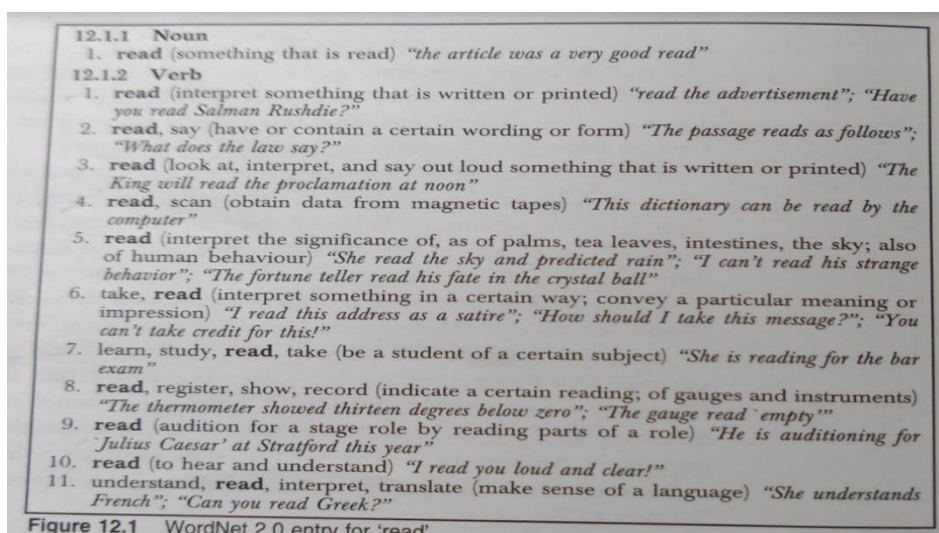


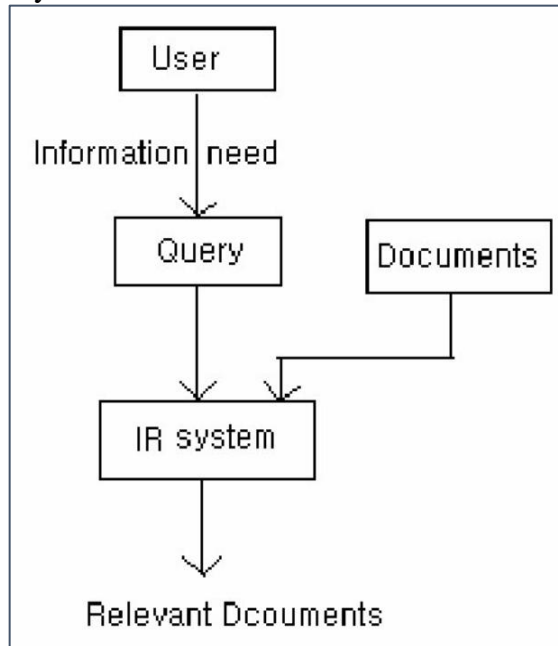
Figure 12.1 WordNet 2.0 entry for 'read'

(ii) Frame Net

- ✓ FrameNet is a large database of semantically annotated English sentences.

	<ul style="list-style-type: none"> ✓ It is based on principles of frame semantics. ✓ It defines a tagset of semantic roles called the frame element. ✓ Sentences from the British National Corpus are tagged with these frame elements. ✓ The basic philosophy involved is that each word evokes a particular situation with particular participants. ✓ FrameNet aims at capturing these situations through case-frame representation of words. ✓ The word that invokes a frame is called target word or predicate, and the participant entities are defined using semantic roles, which are called frame elements. ➤ Each frame contains a main lexical item as predicate and associated frame-specific semantic roles, such as AUTHORITIES, TIME, AND SUSPECT in the ARREST frame, called frame elements. ➤ Example: The sentence below is annotated with semantic roles AUTHORITIES AND SUSPECT ➤ [Authorities The police] nabbed [suspect the snatcher] ➤ The COMMUNICATION frame has the semantic roles ADDRESSEE, COMMUNICATOR, TOPIC, and MEDIUM. ➤ A JUDGEMENT frame contains roles such as a JUDGE, EVALUEE, and REASON. ➤ Example: ➤ [judge She] [Evaluee blames the police] [Reason for failing to provide enough protection] ➤ A frame may inherit roles from another frame. Eg., a STATEMENT frame may inherit from a COMMUNICATION frame, it contains roles such as SPEAKER, ADDRESSEE, and MESSAGE. ➤ Example: ➤ [Speaker She] told [Addressee me] [Message ‘I’ll return by 7:00 pm today’] 			
6	<p>(a) Explain design feature of IR with a neat diagram. (5)</p> <ul style="list-style-type: none"> ➤ The process of IR begins with the user’s information need. ➤ Based on the need, the user formulates a query. ➤ The IR system returns documents that seem relevant to the query. 	[10]	CO4	L3

- The retrieval is performed by matching the query representation with document representation.
- The actual text of the document is not used in the retrieval process.
- Instead documents in a collection are frequently represented through a set of index terms or keywords.



- Representation of keywords provides a logical view of the document.
- The process of transforming document text, to some representation of it, is known as indexing.
- There are different types of index structures.
- The one commonly used is inverted index.
- An inverted index is a list of keywords, with each keyword carrying pointers to the documents containing that keywords.

(b) Define precision and recall. Explain the trade-off between them in evaluation of IR systems. (5)

0.25	1.0
0.4	0.67
0.55	0.8
0.8	0.6
1.0	0.5

Precision:

- ✓ Precision is defined as the proportion of relevant documents in a retrieved set.

✓ It is the probability that a relevant document is retrieved.

✓ It measures the accuracy of a system.

Recall:

✓ Recall is the proportion of relevant documents that are actually been retrieved.

✓ Recall measures the exhaustiveness of the system.

✓

✓ Interpolated Average Precision = 0.745

