



SJBIT

II Jai Sri Gurudev II

Sri Adichunchanagiri Shikshana Trust®

SJB Institute of Technology(Affiliated to Visvesvaraya Technological University, Belagavi and
Approved by AICTE and Accredited by NAAC with 'A' Grade, CGPA-3.22 - New Delhi)

#67, BGS Health & Education City, Dr. Vishnuvardhan Road, Kengeri, Bengaluru - 560060.

Website : www.sjbit.edu.in

INTERNAL ASSESSMENT BOOKStudent Name : Abhishek Kumar RaiSemester & Section : VII A USN : 1JB18IS003Subject: Natural Language Processing Subject Code: 18CS743 Branch: ISEName of Faculty in charge : Chetan R

Internal Assessment Test - I					Internal Assessment Test - II					Internal Assessment Test - III				
Date : <u>19-11-2021</u>					Date : <u>28/12/2021</u>					Date : <u>31/01/22</u>				
Max. Marks : <u>50</u>					Max. Marks : <u>50</u>					Max. Marks : <u>50</u>				
Q No.	PART - A				Q No.	PART - A				Q No.	PART - A			
	A	B	C	Total		A	B	C	Total		A	B	C	Total
1	8	6	6	20	1					1	6	6	8	20
2					2	6	8	2	16	2				
PART - B					PART - B					PART - B				
3					3					3	9	8	4	21
4	4	3	3	10	4	2	7	7	16	4				
I Test IA Marks Total				30	II Test IA Marks Total				32	III Test IA Marks Total				41
Quiz 1/Assignment etc.,				7	Quiz 2/Assignment etc.,				4	Quiz 3/Assignment etc.,				10
Student Signature : <u>Abhishek Kumar Rai</u>					Student Signature : <u>Abhishek Kumar Rai</u>					Student Signature : <u>Abhishek Kumar Rai</u>				
Signature of Invigilator <u>[Signature]</u>					Signature of Invigilator <u>[Signature]</u>					Signature of Invigilator <u>[Signature]</u>				
Signature of Faculty in charge <u>[Signature]</u>					Signature of Faculty in charge <u>[Signature]</u>					Signature of Faculty in charge <u>[Signature]</u>				
Avg. IA Marks for <u>30</u> (A) <u>21</u>					Assignment / Quiz etc. <u>7</u> (B) <u>7</u>					Total IA Marks for <u>40</u> (A+B) <u>28</u>				

HOD

Head of the Department

Principal

Dept. of Information Science & Engineering
S.J.B. Institute of Technology
Kengeri, Bangalore-560 060.

Department of Information Science and Engineering.

Dept. Vision:

We envision our department as a catalyst for developing educated engaged and employable individuals whose collective energy will be driving force for prosperity and the quality of life in our diverse world.

Dept. Mission:

Our mission is to provide quality technical education in the field of information technology and to strive excellence in the education by developing and sharpening the intellectual and human potential for good industry and community.

About Anti - Ragging

SJBIT has zero tolerance policy for ragging. The Institute views ragging, is an uncivilized, and inhuman practice. We do not subscribe to the view that one could wait till something happens in order to initiate stringent action. Any rigorous action in such cases may damage a young career. So we repose faith in averting such eventualities. For this, the Institute has proactive policy.

Punishments for Ragging

1. Cancellation of admission.
2. Suspension from attending classes.
3. Withholding/withdrawing scholarship/fellowship and other benefits.
4. Debarring from appearing in any test/examination or other evaluation process.
5. Withholding results.
6. Debarring from representing the University in any regional, national or international meet, tournament, youth festival etc.
7. Suspensions/expulsion from the hostel.
8. Rustication from the college and University for period varying from 1 to 4 years.
9. Expulsion from the college and consequent debarring from admission to any other college.
10. Rigorous imprisonment of three years and/or a fine of upto Rs.25,000.
11. Collective punishment: When the persons committing or abetting the crime or ragging are not identified, the institution has resort to collective punishment as a deterrent to ensure community pressure on the potential raggars.

Quiz - I

1. d) ✓
2. a) ✓
3. d) ✓
4. a) ✓
5. d) ✓
6. a) ✓
7. b) ✓
8. c) ✓
9. d) ✓
10. b) ✓

(7)

Internals - I

Part - A

- Q.1.A There are 5 different levels of NLP those are:→
- ①. Lexical :- At this level text is processed word by word. Potential meaningful words are recognised.
ex:- Akash is playing cricket

Akash
noun

is playing
verb

cricket

Correct
 - ②. Syntactic :- The text is processed line by line at this level.
ex:- Akash is apple eating

Akash is apple eating

wrong

all the words are correct but sentence all together does not make sense.
 - ③. Semantic :- Once is sentence is processed. Now the context and the meaning of sentence.
ex:- Manu weds Manu

Manu weds Manu

wrong (no meaning).
 - ④. Discourse :- It is used when processing a document the large text is broken down in smaller phrase, sentence and the processed.

- ⑤. Pragmatic :- It is the highest level of processing here the model is aware of outer world information it's knowledge is just limited to the text.

ex:- Tamy weeds Mamy

└─ access the database and check if Tamy and Mamy are really married.

Q.1.B Components of transformational grammar are: →

- ①. Phrase grammar rule :- This divides the text into sub-parts such as noun, verb etc.

$S \rightarrow NP + VP$

$VP \rightarrow V + NP$

$NP \rightarrow Noun + Det$

$V \rightarrow Aux + V$

$Aux \rightarrow$ will, is, be

$Verb \rightarrow$ eating, running

$Det \rightarrow$ a, an, the

ex:- He is running.
└─┬─┬─┬─
Noun Aux Verb

②. Transformational Rule :→ Different types of transformation needs different rules

Ex. - For ~~Active~~ to ~~Active~~ passive

$NP_1 - A - V - NP_2 \longrightarrow NP_2 - A - be - en - V - by - NP_1$

Ex. → The Police will catch the snatcher.

$\underbrace{\text{The Police}}_{NP_1} \quad \underbrace{\text{will}}_A \quad \underbrace{\text{catch}}_{Verb} \quad \underbrace{\text{the snatcher}}_{NP_2}$

6. The snatcher ~~will catch~~ will be caught by the policeman.

$\underbrace{\text{The snatcher}}_{NP_2} \quad \underbrace{\text{will}}_A \quad \underbrace{\text{be caught by}}_{\text{caught.}} \quad \underbrace{\text{the policeman}}_{NP_1}$

③. Morphophonemic rule : - It looks for the words with similar phonemes, sound to distinguish ~~to~~ between words

ex. → bat — bad

Q.1.c. There are two types of language models: →

①. Grammatical based model.

②. Statistical based model

A model is a description of an complex entity or process. NLP is a complex entity.

So the ~~mod~~ description of NLP is a language model.

①. Grammatical: → Grammar is used to create the language model.

6' All rules are hand-coded in it. All known rules of grammar

②. Statistical based model: → A large quantity of text is given to be model to train on.

Model uses the previously learnt data and compare and process the data.

Part - B

Q.4.A. Information Retrieval :→ To extract the required text from the document based the query asked is called Information retrieval.

Issues of Information Retrieval :→

- ① Wrong Question Structure :→ Sometimes user does ask a proper meaningful question which leads to problem.
- ② Match Query in the document :→ Sometimes there are more than one phrases matches the query.
- ③ Unstructured document text :→ Sometimes text does not contain data in proper structure which can be extracted.
- ④ Ambiguity of data :→ Sometimes processing text & create confusion due to ambiguity on phrase and words. ex:- match — Cricket — wedding

Q.4.B. Finite State Automata :-

- Automata taken from greek word means self-work
- It is a pre-defined sequence of processes

A finite Automata can be defined by 5 tuples, —

$$[Q, \Sigma, \delta, q_0, F]$$

$Q \rightarrow$ Set of all States

$\Sigma \rightarrow$ Set of all Symbols

$\delta \rightarrow$ transition function

$q_0 \rightarrow$ Initial State.

$F \rightarrow$ Final State.

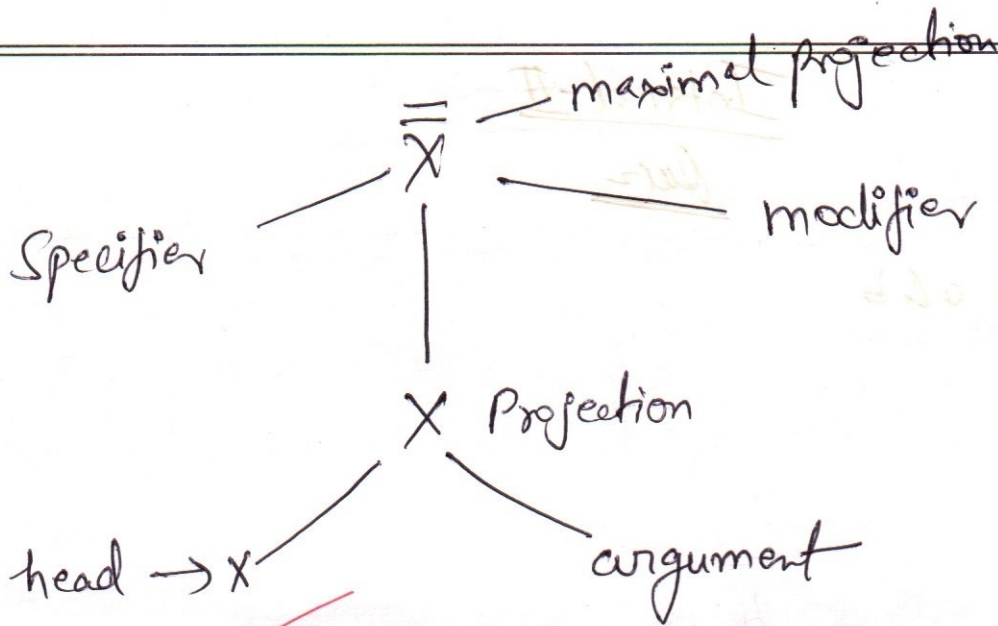
It is of two types.

→ DFA: → Deterministic Finite Automata

→ NFA :- Non-deterministic finite automata,

It is used to find the words and analysis the text by parsing.

Q.4.C.



- \bar{X} theory is a generated theory of a language.
- In this theory Sentence Structure and Phrase Structure are rather than different are considered together as maximal projection.

Sentence Structure + Phrase Structure → Maximal projection

30
50

Internal-II

Quiz

1. ~~a~~. Both a & b
2. ~~b~~.
3. ~~a~~.
4. ~~a~~.
5. ~~b~~.
6. ~~b~~.
7. ~~b~~.
8. ~~a~~.
9. ~~b~~.
10. ~~ca~~.

14

Part-A

Q.2.A. Minimum edit distance algorithm \rightarrow This algorithm is a machine learning algorithm used to compute the minimum number of operations which are needed to be performed to ~~calculate~~ convert one string to another ~~string~~ by using three known methods Insert, Delete and update.

Conversion \rightarrow Given: - Word to convert: - Peaceful
to word: - peaceful

We use the following table matrix to calculate the distance between given words.

	#	P	A	C	E	F	U	L
#	0	1	2	3	4	5	6	7
P	1	0	2	3	4	5	6	7
E	2	1	3	4	5	6	7	8
A	3	2	1	4	5	6	7	8
C	4	3	2	1	5	6	7	8
E	5	4	3	2	1	6	7	8
F	6	5	4	3	2	1	7	8
U	7	6	5	4	3	2	1	8
L	8	7	6	5	4	3	2	1

\rightarrow final output
no. of operation

	#	P	A	C	E	F	U	L
#	0	1	2	3	4	5	6	7
P	1	0	2	3	4	5	6	7
E	2	2	3	4	3	5	6	7
A	3	2	1	3	4	5	6	7
C	4	5	2	1	3	4	5	6
E	5	4	3	2	1	6	7	2
F	6	5	4	3	2	1	7	1
U	7	6	5	4	3	2	1	8
L	8	7	6	5	4	3	2	1

Output
①

The distance between peaceful and paceful is $\$ = 1$.

It means only one operation is to be performed to make peaceful to paceful.

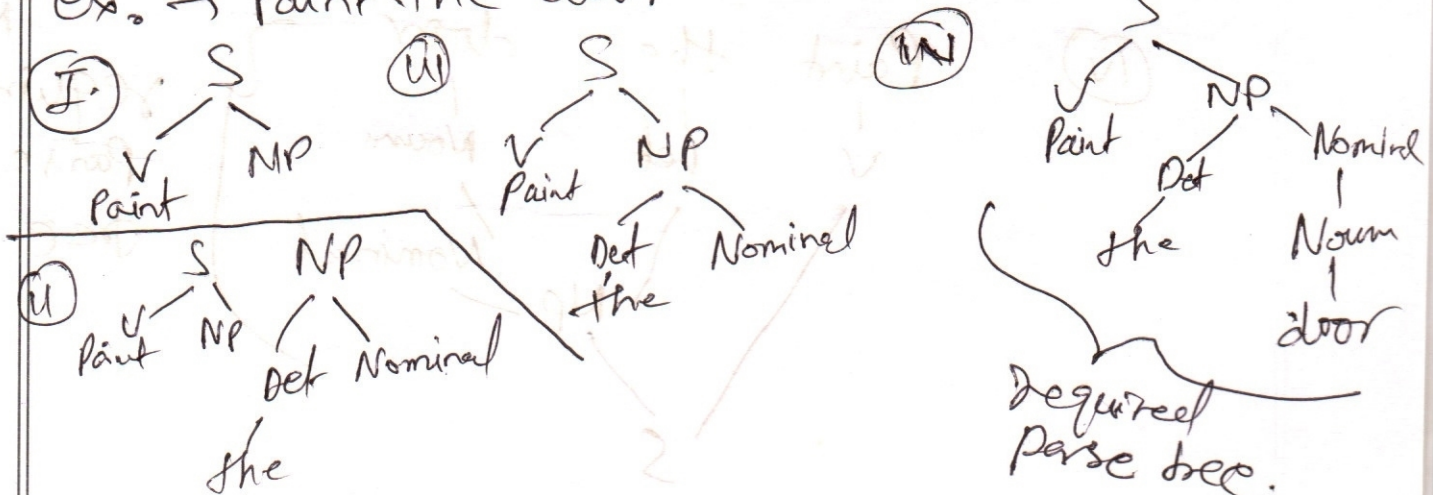
p a c e f u l
 ↑ insert e ——— 1 operation.
 p e a c e f u l

total operation = 1

Q.2.B. Top-down parsing :-

- It is a machine learning algorithm used to find Grammar from starting grammar S.
- It starts from the starting Grammar S of the given input.
- All the sub-tree of S are obtained.
- Use every grammar on S to create sub-tree and expand.
- All the non-terminating nodes are expanded.
- A point is reached when all the leave nodes only contain parts of figure of speech.
- Any non-matching leaves from the input sentence is neglected.
- Successfully obtain the parsed tree.

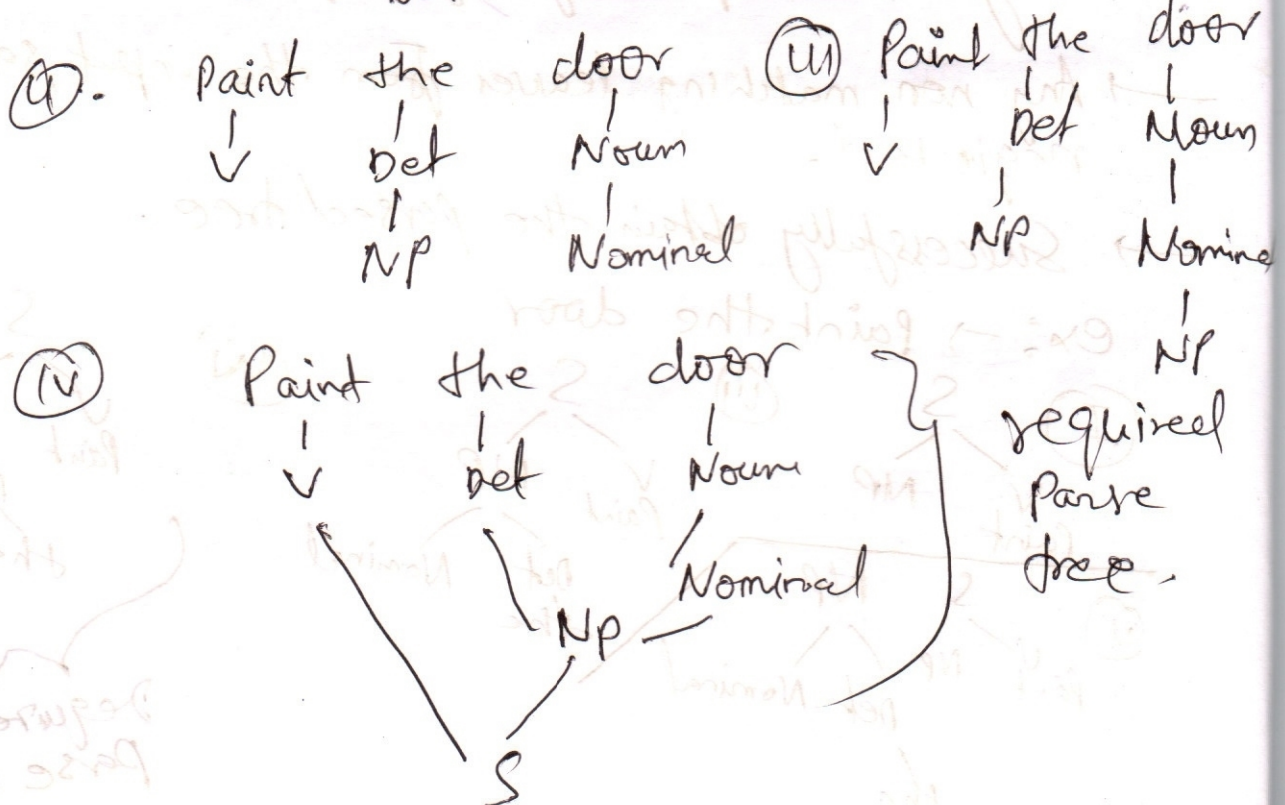
ex: → Paint the door



Bottom-up parsing:-

- It is a machine learning algorithm used to find grammar from sentence.
- It starts from the user input and goes up to until it find the starting grammar S.
- User input it taken to start and further grammar is calculate based on currently present sub tree.
- Once the root node is found check if it is equal to S then accept otherwise reject.

Ex:- Paint the door



Top-down :-

Advantages :- \rightarrow Do not waste time looking for grammar which does not have S as start node.

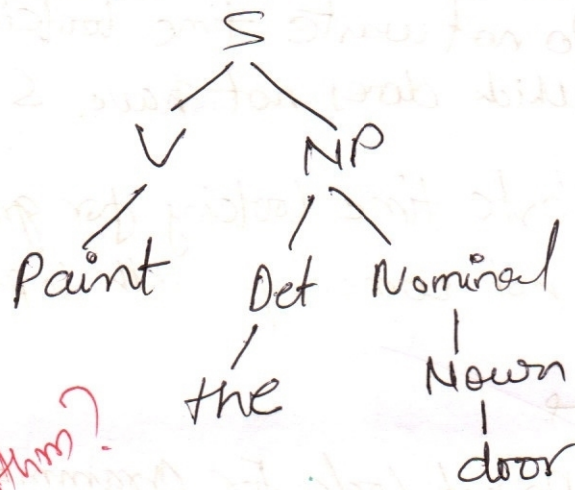
Disadvantages :- \rightarrow Waste time looking for grammar which does not come up with input sentence.

Bottom-down :-

Advantage :- Do not look for grammar which will not come up with input sentence.

Disadvantage :- It waste time looking for grammar which does not have S as starting node.

Q.2.C Top-down depth first algorithm :-



Algorithm?

- This algorithm works from top-downwards the input sentence
- It starts at the top S as first grammar as root node.
- All the sub-tree of S are evaluated for all the grammar and node is expanded
- All the evaluation is done towards the input sentence grammar.
- The leaf node contains the user input sentence grammar
- The algorithm runs ~~can~~ recursively till all the leaf node contains parts of speech grammar.

Part-B

Q.4.c. Active learning approach :-

It is found in research that ~~lab~~-acquiring label from human annotation works better than other approaches.

Steps:-

- ① → All the ~~then~~ sentence cluster are grouped on the basis of same parse sub-tree.
- ② → For each cluster sentences are grouped by same target Verb.
- ③ → Take the largest group from largest cluster and present it to user.
- ④ → Now ~~can~~ allocated the allocated label to all the ~~sub-trees with~~ sentences with same sub-tree.
- ⑤ → Train with the classification.
- ⑥ → Find out any ~~missing~~ miss-match and present to the user.
- ⑦ → Repeat the 4-6 steps until desired accuracy is obtained.

0.4.B

i). The shortest path hypothesis :- This hypothesis is used to calculate the shortest dependency path among the relations in the sentence.

In processing of text by machine learning an error can occur

Ex:- ambiguity

To overcome we find out the different relations in the sentence by shortest dependency path.

This shows how words in a sentence ~~are~~ and the meaning of those words are dependent on other words in the sentence.

The words in the sentence are processed and the dependency through out the sentence is found out

The dependency is used to solve the ~~at~~ ambiguity and to better understand the context of the sentence.

(ii). Learning Dependency path :- To solve the problem to understand the context of the ~~sent~~ sentence it is very important to understand the dependency of the words in the sentence with each other.

7. Dependency path is a way to find the relation b/w the words of a sentence by mapping these relation to the other words in the sentence. The relation mapping can be plotted and referred to further analyses the context of the sentence.

Q.4.A Functional overview of Infact system:-

This system refers to the analysis of language processing and understanding of context of sentence and ~~not~~ relation of words and their dependency to solve any ambiguity in the processing of the text. It uses tagging to solve the problem.

There are three main tagging.

1. Rule-based tagging is used to convert sentences in the form of tuple of sentence or list of sentences.

Stochastic tagging uses frequency and sequence for the analysing and predicting the context of the sentence. When the most frequently used tag is used for tagging it is frequency based.

$$\frac{32}{50}$$

Internals - 1st - Quiz

1. ~~a~~ Zipf's law
2. ~~a~~ Information Retrieval
3. ~~b~~ Search Engine optimization
4. ~~b~~ False
5. ~~d~~ Inverted Index
6. ~~a~~ Gerard Salton
7. ~~b~~ helped, helps → help
8. ~~b~~ Stop word removal.
9. ~~a~~ the terms
10. ~~b~~ recall

10

Part - A

Q.1-A (SVM), Support Vector machine

- It is used for binary classification.
- SVM score is the signed distance from the hyperplane.
- For multi classes problem ~~a~~ series of SVM are trained.
- Every SVM score is mapped with probability.
- Class membership probability is calculated as a function of SVM score trained as ~~for~~ sigmoid function.
- The fit parameters are the slope of sigmoid function.
- SVM score trained as a ~~cost~~ factor both positive and negative ~~cost~~ factors.
- The ratio of positive and negative cost factors gives the decision of right class.

Q1. B. Latent Semantic Analysis feedback system:-

- It is a semantic analysis model.
- It uses similar word from Corpus Document is represent terms.
- It uses Statistical technique for analysing the text Document.
- It uses powerful Mathematical transformation to convert vector space into latent semantic space.
- 6. → It counts the frequency of word ~~are~~, how they occur in the document.
- It does not take in account word order.
- It works good for larger document does not work good for short document.
- It efficiency ~~increases~~ with increase in size of object.
- Whole meaning of the terms is called 'global-focused'.
- Immediate preceeder meaning is called 'local-focused'.

Q.1.C (i) Cohesion :-

→ It measures the relationship between ~~the~~ rhetorical role and a preceding action

$$\rightarrow H = \sum_{r_i, p_i \in H} \frac{P(p_i / r_i)}{|H|}$$

$H \rightarrow$ Hypothesis

$r \rightarrow$ rhetorical role

$p \rightarrow$ action

(ii). Coverage :-

- It measures how good the hypothesis for the mark
- It can use semi-structured data from LSA.
- It define and contain rules for represent.

$$\rightarrow H = \frac{|\text{rules Covered}|}{|\text{Rule Set}|}$$

(iii) Interestingness :→

→ It measures the characteristic of hypothesis by antecedent and consequent.

→
$$H = | \text{Semantic difference between antecedent and consequent} |$$

→ It measures in the unit of KDD approach.

8.

(iv) Plausibility of Origin :→

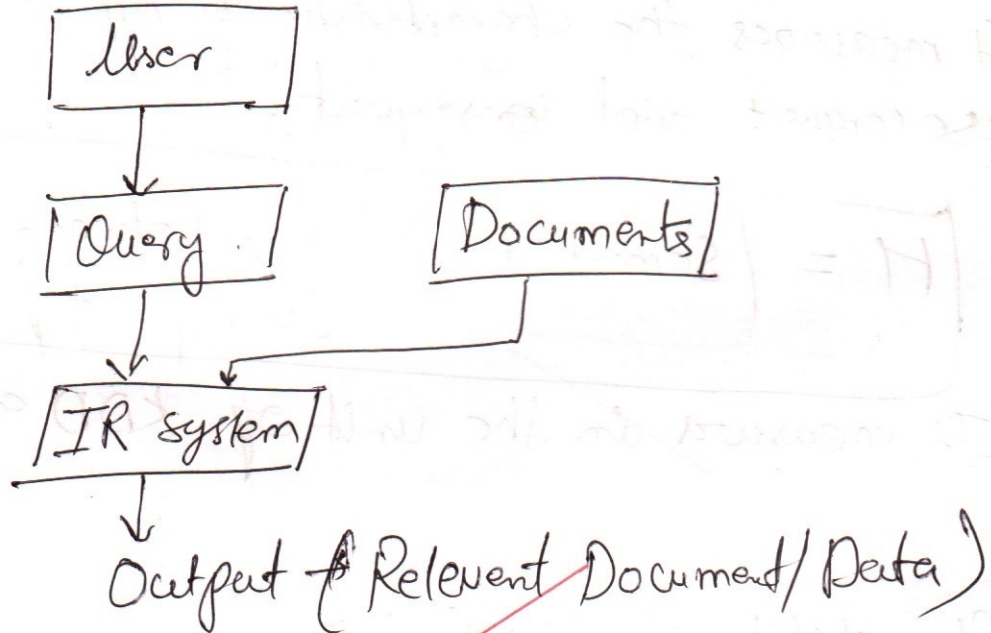
→ It measures how plausible is the hypothesis calculated by sensor's evidence.

→ If a better hypothesis is found it is used or the original hypothesis is given to next generation.

→
$$H = \begin{cases} S_p & \text{if output of Sensor's Crossover} \\ 0 & \text{if output of original} \end{cases}$$

Part-B

Q.2.A.



- User execute the query according to the data required.
- This query is send to the IR-system.
- IR system fetches the relevant information from the documents based on query.
- The relevant information is given as output to user.

Features of IR

Design features:-

①. Inverted Index:-→

- It is a key-value pair which stores the word as key and word frequency as value.
- It is like a hash map and works very fast.
- It is used in search operations.

②. Stop word Elimination:-→

- It is a technique to remove unwanted words from the documents.
- These words contain less or useless meaning.
- ex:- → is, an, the

③. Stemming:-→

- It is a technique to optimize words.
- It reduces the redundancy of the words.
- It helps in memory optimization.
- ex:- → changed → hang, laughing → laugh.
- It reduces the word to its root or stem word.

② Boolean model :-

- It is based on set theory and boolean expression.
- It is most commonly used and easy to use.
- ~~Queries~~ are set theory and ~~Documents~~ are boolean expression.

0.3.B. Boolean Model For Classical Information Model :-

- It is most commonly used Classical model.
- It is based on set theory and boolean expression
- Documents are set theory and queries are boolean expressions.

Ex: →

Documents: → {D1, D2, D3, D4}

D1: → Anuj is a tall boy and Akash is shy.

D2: → Anuj play football and Akash play cricket.

D3: → Cricket is a game of 11 players. Anuj like watching cricket.

D4: → Anuj and Akash are brothers.

Queries: →

①. ~~Anuj and Akash~~ Anuj and Akash.

②. Anuj and Akash ~~but~~ not cricket.
and

Solution: →

①. ~~Anuj and Akash~~
Anuj

Solution: \rightarrow

	D1	D2	D3	D4
Anuj	1	1	1	1
Akash	1	1	0	1
Cricket	0	1	1	0

①. Anuj and Akash: — $Anuj \wedge Akash$.

Anuj — $\{D1, D2, D3, D4\}$

Akash — $\{D1, D2, D4\}$

$Anuj \wedge Akash = \{D1, D2, D4\}$ Ans.

②. Anuj and Akash ~~is~~ ^{and} not Cricket: $\rightarrow (Anuj \wedge Akash) \wedge \neg Cricket$

Anuj — $\{D1, D2, D3, D4\}$

Akash: — $\{D1, D2, D4\}$

$\neg Cricket$: — $\{D1, D4\}$

$Anuj \wedge Akash \wedge (\neg Cricket) = \{D1, D4\}$ Ans.

Q.3.C. FRAME NET :->

(i).

- It is a lexical dataset.
- It contains content which is both machine and human readable.
- It is based on ~~frame~~ frame-syntactic.
- The frame means meaning and lexical means ~~mean~~ one meaning for one word.

(ii). Stemmer :->

- It is a technique used to reduce word to their root or ~~stem~~ stem word/form.
- It helps increase in accuracy.
- It helps reduce memory usage.
- It is fast and ~~can~~ easy to use.
- It is used for stop-word elimination to eliminate less-meaning words.
- Ex:- helped, helps → help.

41/50 ~~41~~



॥ Jai Sri Gurudev ॥
Sri Adichunchanagiri Shikshana Trust

SJB Institute of Technology

College Vision

To become a recognized technical education center with global perspective.

College Mission

To provide learning opportunities that fosters students ethical values, intelligent development in science and technology and social responsibility so that they become sensible and contributing members of the society.

INSTRUCTIONS TO STUDENTS

1. Fill-in the details on front cover page.
2. No sheets be either removed or added after writing a test.
3. No loose sheets are permitted to be used for answering the tests.
4. The question numbers should be mentioned in the margin only.
5. Minimum of 60% marks should be scored in each subject and in all three tests.
6. The candidate shall write answers on both sides of pages of the answer book. Answers must be written using black pen (ball pen or ink pen). If there is a change in pen, the same shall be attested by the Room Superintendent on the facing sheet of the blue book at the top.
7. No candidate shall be permitted to go to toilet during the period of test.
8. After completion of the test, students should handover the test Book immediately to the concerned Faculty.
9. For both theory and Practical classes, 85% Attendance is compulsory.
10. All three tests are compulsory.
11. No candidate shall be permitted to go to toilet during test.
12. Any candidate appearing for the test is liable to be charged with committing malpractice in the following cases.
 - a. Bringing any written material / portions of a book.
 - b. Communicating with any candidate.
 - c. Bringing mobile phone to examination hall.



|| Jai Sri Gurudev ||
Sri AdichunchanagiriShikshana Trust®
SJB Institute of Technology



(A Constituent of BGS & SJB Group of Institutions and Hospitals)

BGS Health & Education City, Dr. Vishnuvardhan Road, Kengeri, Bengaluru-5

Affiliated to Visvesvaraya Technological University, Belagavi. Approved by AICTE, New Delhi. Accredited by NAAC, New Delhi with 'A' Grade. Recognized by UGC, New Delhi with 2(f) and 12(B). Certified by ISO 9001-2015

DEPARTMENT OF INFORMATION SCIENCE AND ENGINEERING

INTERNAL ASSESSMENT PAPER (7th Sem)

Internal Test: 1	Internal Quiz: 1	Academic Year: EVEN /2021-22
Subject : Natural Language Processing	Sub-Code: 18CS743	Sem: 7
Date:19/11/2021	Time:2:00 to 3:45pm	Dur:1.45 min
Internal Test max marks: 50	Internal Quiz max marks: 10	
Staff-Incharge:Chetan R		

I. Quiz (Answer all multiple-choice question in first sheet of your answer book)

Question No	Multiple choice question	BT level	CO Mapped
1	What is the field of Natural Language Processing (NLP)? a. Computer Science b. Artificial Intelligence c. Linguistics d. All of the mentioned.	L2	CO1, PO1
2	NLP is concerned with the interactions between computer and human languages. a. True b. False	L2	CO1, PO1
3	What is the main challenge of NLP? a. Handling Ambiguity of sentences. b. Handling POS-Tagging c. Handling Tokenization d. All of the mentioned	L2	CO1, PO1
4	Modern NLP algorithms are based on machine learning, especially statistical machine learning. a. True b. False	L2	CO1, PO1
5	Choose from the following areas where NLP can be Useful. a. Information retrieval b. Automatic Text summarization c. Automatic Question-Answering Systems d. All of the above.	L2	CO1, PO1
6	What is machine translation? a. Converts human language to machine language. b. Converts one human language to other. c. Converts any human language to English. d. Converts Machine Language to human language.	L2	CO1, PO1
7	Information Retrieval and Information Extraction are the two same thing. a. True b. False	L2	CO1, PO1

8	How many steps are there in NLP? a. 3 b. 4 c. 5 d. 6	L2	CO1, PO1
9	"I saw bats" contains which type of ambiguity? a. Syntactic b. Semantic c. Lexical d. Anaphoric	L2	CO1, PO1
10	"Sita loves her mother and Gita does too" contain which type of ambiguity? a. Syntactic b. Semantic c. Lexical d. Anaphoric	L2	CO1, PO1

II. Internal Test (Answer any two full questions choosing one from each part)
(Each full question carries 25 marks)

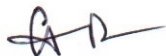
Main Ques	SubQ ues	Full Question	Marks	BT Levels	CO -PO Outcom e
Part A					
1.	A)	Explain different levels of NLP with example.	10	L2	CO1, PO1
	B)	Explain the components of transformational grammar with example.	9	L2	CO1, PO1
	C)	Write a note of different types of language models.	6	L2	CO1, PO1
OR					
2.	A)	Explain the applications of NLP?	8	L2	CO1, PO1
	B)	Explain the different smoothing techniques to handle the data parseness problem in n-gram.	9	L2	CO1, PO1
	C)	Explain Lexical Functional Grammar.	8	L2	CO1, PO1
Part B					
3.	A)	Write the c-structure and f-structure for the sentence: "She saw stars". Given the rules S → NP VP VP → V {NP} {NP} PP* {S'} PP → P NP NP → Det N {PP} S' → Comp S	10	L2	CO1, PO1
	B)	What is Morphological parsing. Explain 2 levels of morphological level with examples.	9	L2	CO2, PO1
	C)	What is NLP? Explain 2 major approaches to NLP.	6	L2	CO1, PO1
OR					
4.	A)	What is Information Retrieval and explain the issues in Information Retrieval.	9	L2	CO1, PO1
	B)	Explain Finite State Automata for parsing word level analysis.	9	L2	CO2, PO1
	C)	Explain \bar{X} theory with example.	7	L2	CO1, PO1

Comments:

— Accepted —



Signature of Faculty



Scrutinizer



HOD

**SCHEMES & SOLUTIONS**

Internal Test: 1	Internal Quiz: 1	Academic Year: ODD / 2021-22
Sub: Natural Language Processing	Sub-Code: 18CS743	Sem: 7
Date: 19/11/2021	Time: 2:00 to 3:45 pm	Dur: 1:45
Internal Test max marks: 50	Internal Quiz max marks: 10	
Staff-Incharge: Chetan R.		

Comments:

Accepted

Signature of Faculty:

Signature of Scrutinizer:

Signature of HOD:

Sl. No.		Marks Alloted
	MCQs. Answer. (1 mark each) 1. d 2. a 3. a 4. a 5. d 6. b 7. b 8. c 9. c 10. b	
1a.	<p><u>Part-A</u></p> <p>Different Levels → Lexical, Syntactic, Semantic, discourse and pragmatic. Explain each level (2 marks each) 2x5=10</p>	10marks.
1b.	<p>Components of transformational grammar.</p> <p>(1) Phrase structure grammar (2) Transformational rules. (3) Morphophonemic rules.</p> <p>3x2=6 marks Example 3marks. } 9marks.</p>	9marks.
1c.	Types of language models. - Grammar-based and Statistical 3marks each.	6marks.
	(OR)	
2a.	<p>Applications of NLP:</p> <p>Machine Translation Speech Recognition Speech Synthesis Natural Language Interfaces to Database</p> <p>Information Retrieval Information Extraction Question Answering. Text Summarization</p> <p>1mark each</p>	8marks.

Q. No.		Marks Alloted
2b.	Add-one Smoothing. Good-Turing Smoothing. Cache Technique. 3 marks each.	9 marks.
2c.	Lexical Functional Grammar C-structure and F-structure } → 2 marks each. CFG rules. — 3 marks. Consistency, Completeness and Coherence } 3 marks.	8 marks.
<u>Part B</u>		
3a.	C-structure → 5 marks. f-structure → 5 marks.	10 marks.
3b.	Definition — 2 marks. Surface level and lexical level → 3 marks each example — 1 marks.	9 marks.
3c.	NLP definition — 2 marks. Explain two approaches :- Rationalist } 2 marks each. Empiricist.	6 marks.

Q. No.		Marks Alloted
4a.	<p>Information Retrieval Definition - 2 marks.</p> <p>issues in IR explanation - 7 marks.</p>	9 marks.
4b.	<p>Finite state Automata — DFA } Explain with NFA } example 4½ each.</p>	9 marks.
4c.	<p>Explain X theory with example.</p> <p>with NP — Noun Phrase VP — Verb Phrase AP — Adjective Phrase PP — Preposition Phrase Sentence structure. projection Maximal projection</p> <p>1 mark each.</p>	7 marks.



|| Jai Sri Gurudev ||
Sri AdichunchanagiriShikshana Trust ®

SJB Institute of Technology

(A Constituent of BGS & SJB Group of Institutions and Hospitals)

BGS Health & Education City, Dr. Vishnuvardhan Road, Kengeri, Bengaluru-5

Affiliated to Visvesvaraya Technological University, Belagavi. Approved by AICTE, New Delhi. Accredited by NAAC, New Delhi with 'A' Grade. Recognized by UGC, New Delhi with 2(f) and 12(B). Certified by ISO 9001-2015



DEPARTMENT OF INFORMATION SCIENCE AND ENGINEERING

INTERNAL ASSESSMENT PAPER (7th Sem)

Internal Test: 2	Internal Quiz: 2	Academic Year: ODD /2021-22
Subject : Natural Language Processing	Sub-Code: 18CS743	Sem: 7
Date:28/12/2021	Time:2:00 to 3:45pm	Dur:1.45 min
Internal Test max marks: 50		Internal Quiz max marks: 10
Staff-Incharge:Chetan R		

I. Quiz (Answer all multiple-choice question in first sheet of your answer book)

Question No	Multiple choice question	BT level	CO Mapped
1	In linguistic morphology, _____ is the process for reducing inflected words to their root form. a. Stemming b. Rooting c. Text-Proofing d. Both a & b.	L2	CO1, PO1
2	Coreference Resolution is: a. Anaphora Resolution b. Given a sentence or larger chunk of text, determine which words ("mentions") refer to the same objects ("entities") c. Both a & b d. None of the above	L2	CO1, PO1
3	Solve the equation according to the sentence "I am planning to visit New York to attend International File Fare Festival." A= (# of words with Noun as the part of speech tag) B= (# of words with verb as the part of speech tag) C=(# of words with frequency count greater than one) What are the correct values of A, B and C? a. 5,5,2 b. 5,5,3 c. 7,5,1 d. 7,4,1	L2	CO1, PO1
4	Which of the following will be POS tagger output when the input sentence is "They refuse to permit" a. [('They', 'PRP', ('refuse', 'VBP'), ('to', 'TO'), ('permit', 'VB')] b. [('They', 'NN', ('refuse', 'VBP'), ('to', 'TO'), ('permit', 'VB')] c. [('They', 'PRP', ('refuse', 'NN'), ('to', 'TO'), ('permit', 'VB')] d. [('They', 'PRP', ('refuse', 'VBP'), ('to', 'PRP'), ('permit', 'VB')]	L2	CO4, PO1

5	In CFG, terminals mainly correspond to While pre-terminals mainly correspond to a. Characters in the language, POS tags b. Words in the language, POS categories c. Words in the language, word relations d. Lexemes, POS Tags	L2	CO1, PO1
6	HMM is used in Phase of NLP. a. Syntactic b. Semantic c. Lexical d. Pragmatics	L2	CO1, PO1
7	Which of the following belongs to the open class group? a. Noun b. Prepositions c. Determiners d. Conjunctions	L2	CO4, PO1
8 is a group of words that may behave as a single unit or phrase. a. Constituency b. Grammatical Relation c. Sub-categorization d. Dependencies	L2	CO4, PO1
9 tagger uses probabilistic and statistical information to assign tags to words. a. Rule based b. Stochastic tagger c. Statistical Tagger d. POS tagger	L2	CO4, PO1
10	"Buy books for children" which type of ambiguity exists in the above sentence? a. Semantic b. Syntactic c. Lexical d. Pragmatic	L2	CO1, PO1

II. Internal Test (Answer any two full questions choosing one from each part)
(Each full question carries 25 marks)

Main Ques	SubQ ues	Full Question	Marks	BT Levels	CO -PO Outcom e
Part A					
1.	A)	Construct the parse tree for the sentence: "The girl plucked the flower with a long stick". Discuss the ambiguity arises from the parse tree.	10	L3	CO1, PO1
	B)	Explain the CYK parser algorithm.	9	L2	CO1, PO1
	C)	Discuss the disadvantages of probabilistic CFG.	6	L2	CO1, PO1
OR					
2.	A)	Explain the minimum edit distance algorithm and compute the distance between "peaceful" and "paceful".	12	L2	CO1, PO1
	B)	Explain the top down and bottom up parsing with suitable example and mention the advantages and disadvantages.	8	L2	CO1, PO1

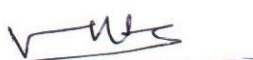
	C)	With example, explain the basic top down depth first algorithm.	5	L2	CO1, PO1
Part B					
3.	A)	With a neat diagram, explain the learning framework architecture.	9	L2	CO1, PO1
	B)	Explain the following: i) Knowledge Roles ii) Domain Knowledge	8	L2	CO2, PO1
	C)	Explain the shortest dependency path hypothesis. Show various shortest dependency path among the relations in the sentence “Jellisc created an atmosphere of terror in the <u>camp</u> by killing abusing and threatening the <u>detainees</u> ”	8	L2	CO1, PO1
OR					
4.	A)	With a neat diagram, explain the functional overview of InFact system.	9	L2	CO1, PO1
	B)	Write the short notes on: i) The shortest path hypothesis. ii) Learning Dependency path.	9	L2	CO2, PO1
	C)	Explain the strategies used in active learning approach for acquiring the labels using the committee based classification scheme.	7	L2	CO1, PO1

Comments:

Accepted


Signature of Faculty


Scrutinizer


HOD

**SCHEMES & SOLUTIONS**

Internal Test: 11	Internal Quiz: 11	Academic Year: ODD / 2021-22
Sub: NLP	Sub-Code: 18CS743	Sem: 7
Date: 28/12/2021	Time: 2:00 to 3:45 pm.	Dur: 1:45
Internal Test max marks: 50	Internal Quiz max marks: 10	
Staff-Incharge: Chetan R.		

Comments:

- Accepted -

Signature of Faculty:

Signature of Scrutinizer:

Signature of HOD:

Q. No.	Marks Alloted
<p><u>Quiz:-</u></p> <ol style="list-style-type: none"> a b d a b a a a b b <p>1a.</p> <p style="text-align: center;"><u>Part A</u></p> <pre> graph TD S --> NP1[NP] S --> VP1[VP] NP1 --> TheGirl[The girl] VP1 --> V1[V] VP1 --> NP2[NP] V1 --> plucked[plucked] NP2 --> Det[Det] NP2 --> Nom[Nom] NP2 --> PP[PP] Det --> the[the] Nom --> N[N] N --> flower[flower] PP --> P[P] PP --> NP3[NP] P --> with[with] NP3 --> along[along] NP3 --> stick[stick] </pre> <p>Attachment Ambiguity : There are two ways of generating prepositional phrase, with a long stick.</p>	<p>10</p> <p>10</p>

1b. CYK parsing algorithm

Let $w = w_1 w_2 w_3 \dots w_i \dots w_j \dots w_n$ and $w_{ij} = w_i \dots w_{j-1}$

// initialization step
for $i=1$ to n do
for all rules $A \rightarrow w_i$ do
 $\text{chart}[i,1] = \{A\}$

// Recursive step

for $j=2$ to n do
for $i=1$ to $n-j+1$ do

begin

$\text{chart}[i,j] = \phi$

 for $k=1$ to $j-1$ do

$\text{chart}[i,j] := \text{chart}[i,j] \cup \{A \mid A \rightarrow BC \text{ is a production and}$
 $B \in \text{chart}[i,k] \text{ and } C \in \text{chart}[i+k, j-k]\}$

 end

if $S \in \text{chart}[1,n]$ then accept else reject.

Algorithm : 5 marks

Explanation : 4 marks

9 marks

1c. Disadvantages :-

1. Independence assumption
2. Lack of sensitivity to lexical information.

3 marks each.

6 marks

2a.

	#	p	a	c	e	f	u	l
#	0	1	2	3	4	5	6	7
p	1	0	1	2	3	4	5	6
e	2	1	2	3	4	5	6	7
a	3	2	1	2	3	4	5	6
c	4	3	2	1	2	3	4	5
e	5	4	3	2	1	2	3	4
f	6	5	4	3	2	1	2	3
u	7	6	5	4	3	2	1	2
l	8	7	6	5	4	3	2	1

Minimum Edit

distance is 1

Algorithm - 6 marks

problem - 6 marks
with explanation

12 marks

2b.

Top down parsing Explanation with example (4 marks)

Bottom up parsing Explanation with example (4 marks)

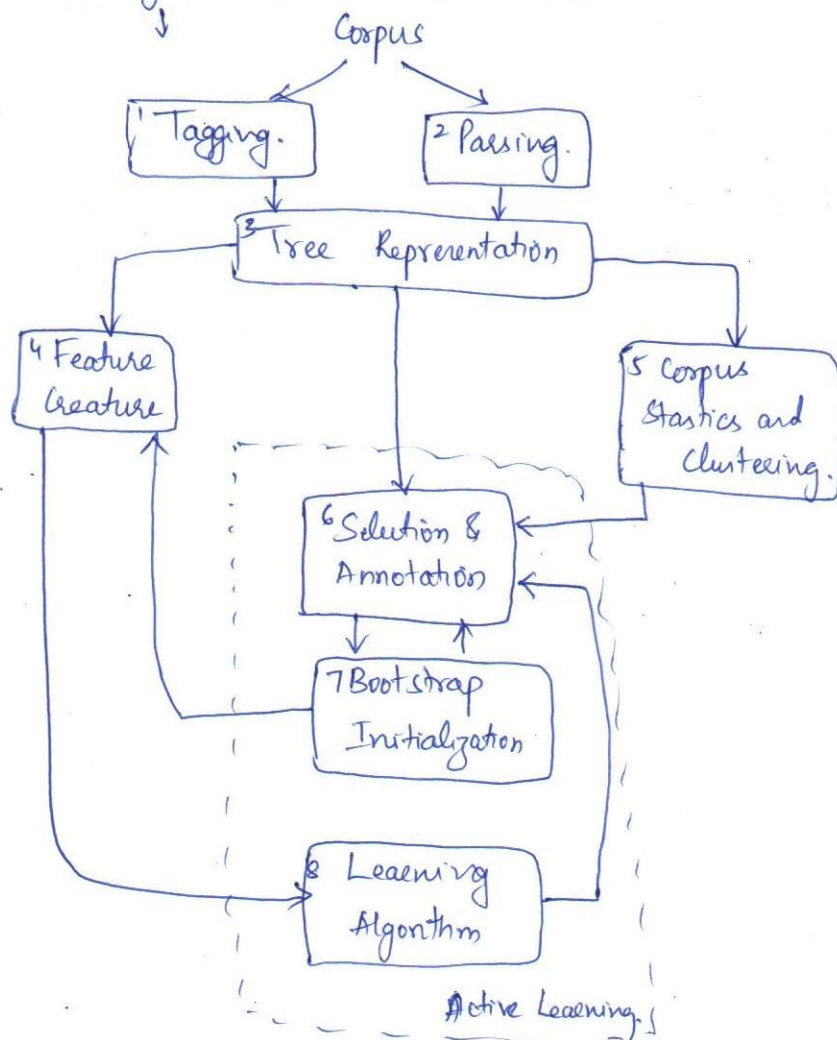
8 marks

- 2c. Basic top down depth first algorithm :-
3 marks explanation 2 marks.

Part B

- 3a. Learning framework architecture

Diagram (4 marks.) Explanation 5 marks.



- 3b. Domain Knowledge :- Domain of interest, Focus is on textual diagnostic reports; Service providers (4 marks)
and customers.
Domain Concepts and Domain specific terms. 8 marks.

Knowledge Roles :- Introduced in Common KADS.

Abstract names that refer to the role a domain concept. Each role has associated with (4 marks)
Text + phrase.

Q. No.		Marks Alloted
3c.	<p><u>Shortest path hypothesis</u></p> <p>If e_1 and e_2 are two entities mentioned in the same sentence such that they are observed to be in a relationship R, then the contribution of the sentence dependency graph to establishing the relationship $R(e_1, e_2)$ is almost exclusively concentrated in the shortest path between e_1 and e_2.</p> <p>"Jelenc created an atmosphere of terror at the <u>camp</u> by killing abusing and threatening the <u>detainees</u>."</p> <p>detainees \rightarrow killing \leftarrow Jelenc \rightarrow created \leftarrow at \leftarrow camp detainees \rightarrow abusing \leftarrow Jelenc \rightarrow created \leftarrow at \leftarrow camp detainees \rightarrow threatening \leftarrow Jelenc \rightarrow created \leftarrow at \leftarrow camp detainees \rightarrow killing \rightarrow by \rightarrow created \leftarrow at \leftarrow camp detainees \rightarrow abusing \rightarrow by \rightarrow created \leftarrow at \leftarrow camp detainees \rightarrow threatening \rightarrow by \rightarrow created \leftarrow at \leftarrow camp</p>	8 marks
4a.	<p>Infact system - diagram \rightarrow 4 marks. Explanation 5 marks.</p>	9 marks
4b.	<p>The shortest path hypothesis \rightarrow 4 marks Learning Dependency path. \rightarrow 5 marks.</p>	
4c.	<p><u>Active Learning:-</u></p> <ol style="list-style-type: none"> Divide the corpus in clusters. within each cluster, group the sentences. Select sentences from largest groups. Bootstrap initialization Train all the classifiers of the committee. Get a pool of instances where the classifiers of the committee disagree and present to the user. Repeat steps d)-f) a few times until a desired accuracy of classification is achieved. 	8 marks



|| Jai Sri Gurudev ||
Sri AdichunchanagiriShikshana Trust ®

SJB Institute of Technology

(A Constituent of BGS & SJB Group of Institutions and Hospitals)

BGS Health & Education City, Dr. Vishnuvardhan Road, Kengeri, Bengaluru-5

Affiliated to Visvesvaraya Technological University, Belagavi. Approved by AICTE, New Delhi. Accredited by NAAC, New Delhi with 'A' Grade. Recognized by UGC, New Delhi with 2(f) and 12(B). Certified by ISO 9001-2015



DEPARTMENT OF INFORMATION SCIENCE AND ENGINEERING

INTERNAL ASSESSMENT PAPER (7th Sem)

Internal Test: 3	Internal Quiz: 3	Academic Year: ODD /2021-22
Subject : Natural Language Processing	Sub-Code: 18CS743	Sem: 7
Date:31/01/2022	Time:1:00 to 2:45pm	Dur:1.45 min
Internal Test max marks: 50		Internal Quiz max marks: 10
Staff-Incharge:Chetan R		

I. Quiz (Answer all multiple-choice question in first sheet of your answer book)

Question No	Multiple choice question	BT level	CO Mapped
1	The formula used to estimate the vocabulary size of a collection is known as: a) Zipf's law b) Power law c) Heap's law d) Compression ratio	L2	CO4, PO1
2	IR Stands for _____. a) Information Retrieval b) Information Retired c) Inform Retrieval d) Information Ready	L2	CO4, PO1
3	SEO stands for _____. a) Search English Optimization b) Search Engine Optimization c) Search Engine Operator d) Search Engine Operation	L2	CO4, PO1
4	The purpose of the inverse document frequency is to increase the weight of terms with high collection frequency. a) True b) False	L2	CO4, PO1
5	A data structure that maps terms back to the parts of a document in which they occur is called an a) Postings list b) Incidence Matrix c) Dictionary d) Inverted Index	L2	CO4, PO1
6	The first large information retrieval research group was formed by _____ at cornell in 1960. a) Gerard Salton b) Ratan Tata c) Ramesh Bush d) Think Roy	L2	CO4, PO1
7	Pick the stemming actions a. was, am, are, is → beb. helped, helps → help c. troubled, troubling, trouble → trouble d. friend, friendship, friends, friendships → friend	L2	CO4, PO1
8	The process of removing most common words (and, or, the, etc.) by an information retrieval system before indexing is known as a) Lemmatization b) Stop word removal c) Inverted indexing d) Normalization	L2	CO4, PO1
9	An inverted index arranges data in a sorted order as per a) the documents b) the frequency of each document c) the frequency of each term d) the terms	L2	CO4, PO1
10	Which of the following is a non-decreasing function of the number of documents retrieved? a) precision b) recall c) accuracy d) mean average precision	L2	CO4, PO1


II. Internal Test (Answer any two full questions choosing one from each part)
(Each full question carries 25 marks)

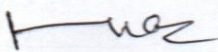
Main Ques	SubQues	Full Question	Marks	BT Levels	CO -PO Outcom e
Part A					
1.	A)	Explain SVM Learning method in Sequence Model Estimation.	8	L3	CO4, PO1
	B)	Explain Latent Semantic Analysis feedback system.	9	L2	CO4, PO1
	C)	Define the following: i) Cohesion ii) Coverage iii) Interestingness iv) Plausibility of Origin	8	L2	CO4, PO1
OR					
2.	A)	Write a note on various approaches to analyzing texts.	10	L2	CO4, PO1
	B)	With a neat diagram explain the evolutionary model for KDT.	10	L2	CO4, PO1
	C)	Define: i) coh-metrix ii) LSI	5	L2	CO4, PO1
Part B					
3.	A)	Explain design features of IR with a neat diagram	9	L2	CO4, PO1
	B)	With an example explain Boolean Model for Classical Information Model	8	L2	CO4, PO1
	C)	Write Short note on: i) FRAME NET ii) Stemmer	8	L2	CO4, PO1
OR					
4.	A)	Explain WORD NET. List the applications	9	L2	CO4, PO1
	B)	Explain non-classical model of IR	6	L2	CO4, PO1
	C)	How stemming effects the performance of IR Systems? Stop words elimination may be harmful, Justify	10	L2	CO4, PO1

Comments:

Accepted.


Signature of Faculty


Scrutinizer


HOD



Sri AdichunchanagiriShikshana Trust (R)

SJB Institute of Technology

(Affiliated to Visvesvaraya Technological University, Belagavi & Approved by AICTE, New Delhi.)

Department of Information Science & Engineering



SCHEMES & SOLUTIONS

Internal Test: <u>III</u>	Internal Quiz: <u>III</u>	Academic Year: ODD / 2021-22	
Sub: <u>Natural Language Processing</u>	Sub-Code: <u>18CS743</u>	Sem: <u>7th A and B</u>	
Date: <u>31/01/2022</u>	Time: <u>1:00pm to 2:45pm.</u>	Dur: <u>1:45 minutes.</u>	
Internal Test max marks: 50		Internal Quiz max marks: 10	
Staff-Incharge: <u>Chetan R.</u>			

Comments:

- Accepted.

Signature of Faculty:

Signature of Scrutinizer:

Signature of HOD:

Q. No.		Marks Alloted
Ans.	<p>1. a 6. a</p> <p>2. a 7. b</p> <p>3. b 8. b</p> <p>4. b 9. d</p> <p>5. d 10. b.</p>	
	<u>Part A</u>	
1a.	<p>SVM Learning method. Explanation. → 4 marks.</p> <p>Sequence Model Estimation. → 4 marks.</p>	8.
1b.	<p>Latent Semantic Analysis Feedback System.</p> <p>LSA → Explanation 4 marks.</p> <p>Feedback System - 5 marks.</p>	9.
1c.	<p>Define. i) cohesion ii) Interestingness each (2marks)</p> <p>ii) Coverage iv) Plausibility of Origin</p>	8.

2a. Various approaches to analyzing texts.

a) Traditional approach

b) Function property.

c). Narrative Theory.

d). LSA.

e) Homogeneous.

(each - 2 marks)

10

2b. Evolutionary model for KDT

Diagram → 6 marks.

Explanation → 4 marks.

10.

2c. Define :- Coh-matrix

LSI

each (2.5 marks)

5

3a)

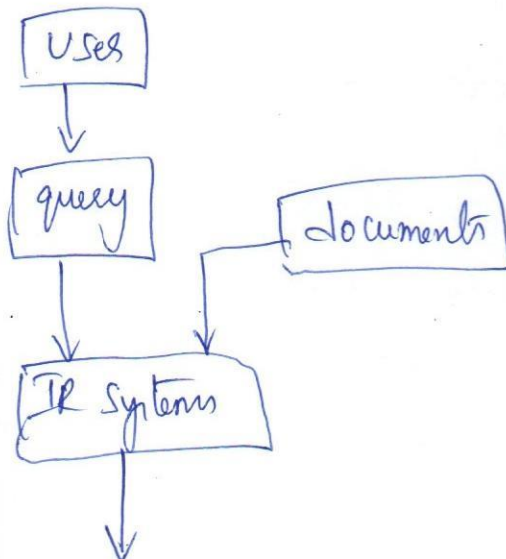


Diagram 3 marks.

Explanation 6 marks.

9.

Q. No.		Marks Alloted
3b.	<p>Boolean Model. \rightarrow AND, OR, NOT. \rightarrow 4 marks.</p> <p>Equation. $Q = \bigwedge (\bigvee \theta_i)$ \rightarrow 2 marks.</p> <p>(i) $R_i \{d_j / \theta_i \in d_j\}$ where $\theta_i = \{t_i, \neg t_i\}$ $\neg t_i \in d_j$</p> <p>(ii). $\cap R_i$</p> <p>Example - 2 marks</p>	8.
3c.	<p>(i) FRAMENET</p> <p>(ii) Stemmer.</p> <p>Explanation. each (4 marks.)</p>	8.
4a.	<p>WORD NET \rightarrow Explanation. 5 marks.</p> <p><u>Applications</u></p> <p>1. WORDNET \rightarrow Automatic text classification Automatic text Summarization. Machine Translation. NLTP Text word disambiguation.</p>	9.

Q. No.		Marks Alloted
4b).	<p>Non-clanical model of IR.</p> <ol style="list-style-type: none"> 1. Information Logic Model 2. Situation theory Model. 2 marks each , 3. Interaction model. 	6mks.
4c.	<p>Stemming affects the performance of IR system</p> <p>Explanation with example - 5 marks.</p> <p>Stop words elimination may be harmful - 5 marks</p> <p>Justification</p>	10mks.



||Jai Sri Gurudev||
Sri AdichunchanagiriShikshana Trust (R)

SJB Institute of Technology

Affiliated to Visvesvaraya Technological University, Belagavi & Approved by AICTE, New Delhi
Accredited by NAAC, New Delhi with 'A' Grade, Recognized by UGC, New Delhi with 2(f) and 12(B)
Certified by ISO 9001- 2015

BGS Health & Education City, Dr. Vishnuvardhan Road, Kengeri, Bengaluru – 560 060
Website: www.sjbit.edu.in, Email : principal@sjbit.edu.in, Mob: +91-6366041109



Department of Information Science and Engineering

Course Name: NATURAL LANGUAGE PROCESSING

Course Code: 18CS743

MODULE 1:

Q. No.	Question	CO Mapped	PO Mapped	Blooms Level
1.	What is NLP? Explain two major approaches to NLP.	CO1	PO1, PO2	L2
2.	Explain the components of transformational grammar.	CO1	PO1, PO2	L2
3.	Explain different levels of NLP with example.	CO1	PO1, PO2	L2
4.	Explain different smoothing techniques to handle the data sparseness problem in n-grain model.	CO1	PO1, PO2	L2

MODULE 2:

Q. No.	Question	CO Mapped	PO Mapped	Blooms Level
1.	What is Morphological Parsing? Explain the two step of Morphological parser.	CO2	PO1, PO2, PO3	L2
2.	Explain spelling correction algorithm.	CO2	PO1, PO2, PO3	L2
3.	With example explain basic top down depth first algorithm	CO2	PO1, PO2, PO3	L2
4.	Explain CYK algorithm.	CO2	PO1, PO2, PO3	L2

MODULE 3:

Q. No.	Question	CO Mapped	PO Mapped	Blooms Level
1.	With neat diagram explain functional overview of InFact System.	CO2	PO1, PO2, PO3	L2
2.	Write a short note on: i) The shortest path hypothesis. ii) Learning with dependency path.	CO1	PO1, PO2	L2
3.	With neat diagram explain the learning framework architecture.	CO1	PO1, PO2	L2
4.	Explain the following i) Domain Knowledge ii) Knowledge roles.	CO1	PO1, PO2	L2

MODULE 4:

Q. No.	Question	CO Mapped	PO Mapped	Blooms Level
1.	Explain SVM learning method in Sequence Model estimation.	CO3	PO1, PO2, PO3	L2
2.	Explain Latent Semantic Analysis feedback system.	CO3	PO1, PO2, PO3	L2
3.	Define the following: i) Cohesion ii) Interestingness iii) Coverage iv) Plausibility of origin.	CO3	PO1, PO2, PO3	L2

MODULE 5:

Q. No.	Question	CO Mapped	PO Mapped	Blooms Level
1.	State and explain Zipf's Law.	CO4	PO1, PO2	L2
2.	Explain Non-classical model of IR	CO4	PO1, PO2	L2
3.	With example explain Boolean model for classical information retrieval.	CO4	PO1, PO2	L2



In-Charge

[CHETAN R]



HOD

Dr. Mohan H S

[Professor & Head, ISE]