

Branch	Test Date	Semester	Div.	Roll No.	Student's Signature		
CSE		VII					
IA Test No.	Subject Natural Language Processing (NLP)						
Junior Supervisor's full signature with date :	Question No.	1	2	3	Total 20	Examiners Signature	Student's Sign After receiving the assessed answer sheet
	Marks obtained						

Q1(a) List the phases of NLP

Ans ① Morphological phase

④ Pragmatic phase

② Syntax analysis phase

⑤ Discourse phase

③ Semantic analysis phase

Q1(b) What is referential ambiguity in natural language processing? Give one example

Ans Either a text is mentioned as an entity (something, someone) and then refers to it again, possibly in a different sentence, using another word. Pronoun causes ambiguity when it is not clear which noun it is referring to.

Eg - The boy told his father about the theft. He was very upset. He is referentially ambiguous because it can refer to both the boy and the father.

Q1(c) List the applications of NLP

Ans ① Auto correction ② Online chatbot ③ Auto completion  
 ④ Sentiment Analysis ⑤ Spam Mail filtering  
 ⑥ Recommendation engine ⑦ Information Retrieval.

Q1(d) Preprocessing steps  
Ans Removing HTML tags, Removal of white space, removal of accented character removal, expand contraction, removal of special character, convert all words into lower case, removal of punctuation mark, removal of words containing digits, removal stop words.

### Q1(e) Stemming

- 1) It is an elementary rule based process for removing inflectional word forms from a token & output are the stem.
- 2) Laughing → laugh  
Studies → studie

### Lemmatization

It is a systematic step by step process for removing inflectional forms of a word.

- 2) Running → Run  
Studies → Study

### Q1(f) Inflection Morphology

- 1) It is a morphological process that adapts existing words so that they function efficiently
- 2) Mostly added as suffix
- 3) Cannot change part of speech
- 4) They are more regular
- 5) Eg - cats → cat

### Derivational Morphology

- 1) It is concerned with the way morphemes are connected to existing lexical form as affix.

- 2) Both as prefix, suffix
- 3) Can change P.O.S.
- 4) They are less regular
- 5) Danger → Dangerous  
Legal → Legalize

### Q1(g) Apply porter algorithm on - Writing, King

Ans As per Porter's Algorithm

(\* vowel \*) ing → φ

(\* vowel \*) ed → φ

∴ Writing → write

King → King

If the word contains vowel then suffix -ing will be removed.

Q1(b) Define Minimum edit distance used in Spelling correction  
Ans Minimum edit distance is the minimum number of insertion, deletion, transposition & substitution reqd to transform the misspelled word into a valid word.

Eg INTENTION EXECUTION

INTENTION  
EXECUTION  
↑↑↑↑↑↑  
I S S D I

insertion (I) = 2

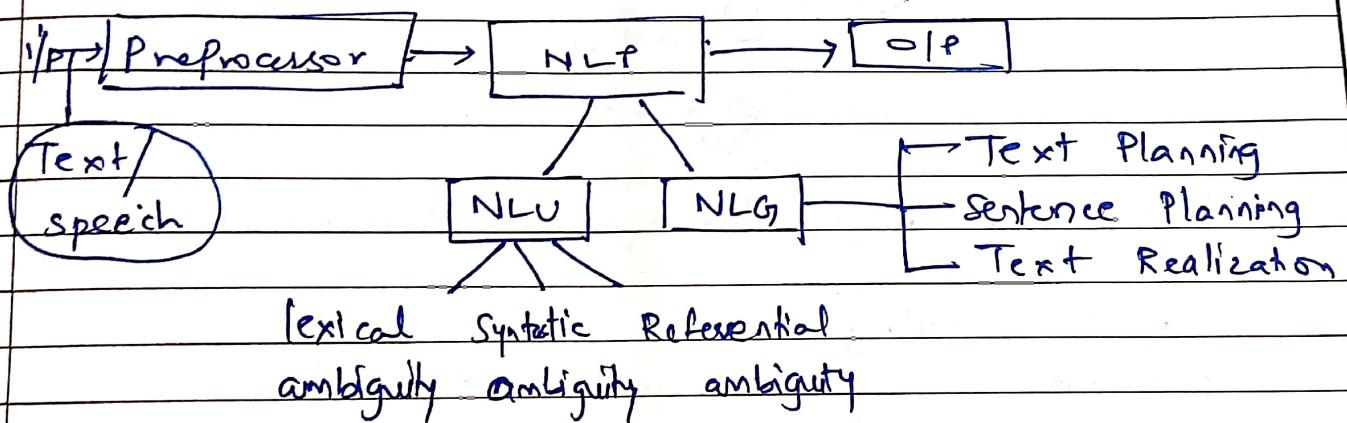
Substitution (S) = 2

Deletion (D) = 1

Minimum edit distance (5)

Q2(a) Explain Natural language generator composed of NLP generic block diagram

Ans Generic NLP system



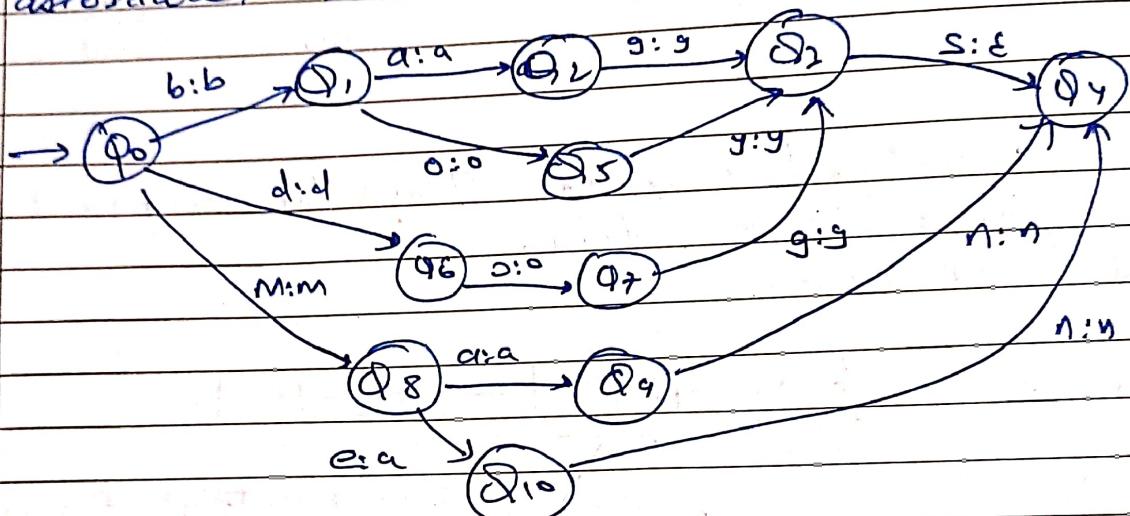
### NLG - Natural Language Generator

It is the process of producing meaningful phrases & sentences in the form of natural language from some internal representation. It involves -

- 1) Text Planning - It includes retrieving the relevant content from the knowledge base.
- 2) Sentence planning - It includes choosing required words, forming meaningful phrase & setting the tone of sentence.

3) Text realization - It is mapping sentence plan into sentence structure.

- Q2(b) Explains FST in detail  
Ans Finite state Transducer are finite state machine to check a string. It is used to build morphological analyzer. It translates strings from one language to another language.  
- Like finite state automata, but here each edge is associated with two strings.



Take men as i/p & recognize Man + PL + N.

### 2-level morphology

cat → cats (regular word + N + PL)

fox → foxes (irregular word + N + PL)

Given i/p cats we get o/p cat + N + PL

2 level morphology [lexical level

[surface level]

lexical level - morphemes & features

surface level - Actual spelling.

lexical       $\{ \underline{\text{c}} \underline{\text{l}} \underline{\text{a}} \underline{\text{t}} \underline{\text{l}} \underline{\text{+N}} \underline{\text{+PL}} \} \}$

surface       $\{ \underline{\text{c}} \underline{\text{l}} \underline{\text{a}} \underline{\text{t}} \underline{\text{l}} \underline{\text{s}} \} \}$

Q2(c) for a corpus, MLE for bigram "battery life" is 0.27  
 frequency of "battery" is 800. After applying  
 Laplace smoothing the MLE for "battery life"  
 becomes 0.025. What is the vocabulary size  
 of the corpus?

Sol:

$$P_{MLE}(\text{life} \mid \text{battery}) = f(\text{battery life}) / f(\text{battery})$$

$$0.27 = f(\text{battery life}) / 800$$

$$f(\text{battery life}) = 800 \times 0.27 = 216$$

with Laplace smoothing

$$P_{MLE}(\text{life} \mid \text{battery}) = \frac{f(\text{battery life}) + 1}{f(\text{battery}) + V}$$

$$0.025 = \frac{217}{800 + V}$$

$$V = \frac{197}{0.025} = 7880$$

$$\boxed{\text{Size of vocabulary} = 7880}$$

Q2(a) find the probability based on given corpus.

① Michael & Zack played at the playground = S<sub>1</sub>

$$P(\text{Michael} \mid S_1) = \frac{1}{5} \quad P(\text{played} \mid \text{Zack}) = \frac{1}{3}$$

$$P(\text{and} \mid \text{Michael}) = \frac{2}{2} \quad P(\text{at} \mid \text{played}) = \frac{1}{1}$$

$$P(\text{Zack} \mid \text{and}) = \frac{3}{3} \quad P(\text{the} \mid \text{at}) = \frac{2}{2}$$

$$P(\text{playground} \mid \text{the}) = \frac{2}{5} \quad P(\text{the} \mid \text{playground}) = \frac{1}{2}$$

$$P(S_1) = \frac{1}{8} \times 1 \times 1 \times \frac{2}{5} \times \frac{1}{3} \times \frac{1}{2} = \frac{1}{720}$$

(2)  $S_2 = \text{Bob went to the school}$

$$P(S_2) = ?$$

$\therefore$  one probability is 0

$\therefore$  without laplace smoothing

$$P(\text{Bob} | S_2) = \frac{2}{5}$$

$$P(S_2) = 0$$

$$P(\text{went} | \text{Bob}) = \frac{0}{2}$$

(3) The school was huge =  $S_3$

$$P(\text{The} | \langle s \rangle) = \frac{2}{5}$$

$$\therefore P(S_3) = \frac{2}{5} \times \frac{2}{5} \times \frac{1}{3} \times \frac{1}{2}$$

$$P(\text{school} | \text{The}) = \frac{3}{5}$$

$$P(\text{was} | \text{school}) = \frac{2}{3}$$

$$= \frac{2}{25}$$

$$P(\text{huge} | \text{was}) = \frac{1}{2}$$

$$P(\langle s \rangle | \text{huge}) = \frac{1}{1}$$

(4) Zack went to the playground =  $S_4$

$$P(\text{Zack} | \langle s \rangle) = \frac{1}{8}$$

$$P(\text{went} | \text{Zack}) = \dots \quad \therefore P(S_4) = 0$$

$\therefore$  highest probability is for sentence  
"The school was huge"

$$\text{Perplexity} = (P(s))^{-1/N}$$

$$= \left(\frac{2}{25}\right)^{-1/6} = \left(\frac{25}{2}\right)^{1/6} = 1.495$$

Q3(b) Write a note on Language model & evaluation metrics for language model

Ans Language modelling is the way of determining the probability of any sequence of words. Language modelling is used in wide variety of application such as speech recognition, spam filtering etc.

- There are two types of language modelling -

i) Statistical Language Modelling

- Development of probabilistic model that are able to predict the next word in the sequence given the word that precedes.  
Eg - N-gram model

ii) Neural Language Modelling

- Neural network methods are achieving better results than classical methods both in stand-alone language model & when it is incorporated into large model on challenging task like speech recognition & machine translation.

Eg - word embedding.

iii) N-gram Language Model

- It can be defined as the contiguous sequence of items from a given sample of text or speech.

- It predicts the probability of a given N-gram within any sequence of words in the language.

## Metrics for language Model

1) Entropy - It is a measure of the amount of information conveyed by Claude Shannon.

$$H(p) = \sum_n p(n) \cdot (-\log(p(n)))$$

$H(p)$  is always greater than equal to 0.

2) Cross Entropy - It measures the ability of the trained model to represent test data.

$$H(p) = \sum_{i=1}^n \frac{1}{n} (-\log_2(p(w_i | u; \theta)))$$

The cross entropy is always greater than or equal to entropy i.e. the model uncertainty can be no less than the true uncertainty.

③ Perplexity - It is a measure of how good a probability distribution predicts a sample. It can be understood as a measure of uncertainty. The perplexity can be calculated by cross entropy to the exponent of 2.

Lower the perplexity = better the model