

Space for Marks	Question No.	START WRITING HERE
		<u>NLP MSE-1 Solution</u>
	Q1(a)	Define affixes & its list its types Ans Affixes are a grammatical element that is combined with a word, stem or phrase to produce derived or inflectional form Types of affixes - prefix suffix
		Eg unwillingness here there are three morphemes (affixes) un - prefix willing - stem word ness - suffix
	Q1(b)	What is lexical ambiguity in natural language processing? Give one example. Ans Lexical ambiguity occurs when a word carries a different sense i.e. having more than one meaning & sentence in which it is combined can be interpreted differently depending on its correct sense. Lexical ambiguity can be resolved using part of speech tagging technique
		Eg - The chicken is ready to eat.

Space for  
MarksQuestion  
No.

START WRITING HERE

Q1(c) List the application of NLP

- (a) Text summarization
- 2) Recommendation engine
- 3) Opinion Mining
- 4) Sentiment Analysis
- 5) Quillbot - Paraphrasing tool
- 6) ChatBot / Chat GPT
- 7) Alexa / Siri

Q1(d) For a corpus, the maximum likelihood estimate (MLE) for bigram "Battery life" is 0.27. frequency of "Battery" is 800. After applying laplace smoothing the mle for battery life 0.025. What is the vocabulary size of the corpus.

Ans

$$P_{MLE}(\text{life/Battery}) = \frac{f(\text{Battery life})}{f(\text{Battery})}$$

$$0.27 = \frac{f(\text{Battery life})}{800}$$

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

With Laplace Smoothing

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

$$0.025 = \frac{216 + 1}{800 + V}$$

$$V = 197$$

$$V = 7880$$

Space for Marks	Question No.	START WRITING HERE													
	8(e)	Differential between Stemming & Lemmatization with example													
Ans		<p>1) Stemming is an elementary rule based process for removing inflectional forms from token's output are the stem/root of the word.</p> <p>Eg - ① Laughing, Laughed, Laughs will become Laugh ② Changing → change</p> <p>2) Lemmatization is a systematic step-by-step process for removing inflectional forms of a word. It makes use of vocabulary, word structure, part of speech tags, and grammar relations.</p> <p>Output of Lemmatization is the root word called a lemma.</p> <p>Eg - ① Running, Run, Ran → Run. ② Changing → change.</p>													
	8(f)	Differential Between Inflectional & Derivational Morphology with example													
Ans		<table border="1"> <thead> <tr> <th></th> <th>Inflectional Morphology</th> <th>Derivational Morphology</th> </tr> </thead> <tbody> <tr> <td>①</td> <td>The study of the modification of words to fit into different grammatical contexts.</td> <td>The study of the formation of new words that differ either in syntactic category or in meaning from their base.</td> </tr> <tr> <td>②</td> <td>Create new forms of same word.</td> <td>Creates new words.</td> </tr> <tr> <td>③</td> <td>Run - Running VB            VB</td> <td>Danger → Dangerous Noun            Adj. (JJ)</td> </tr> </tbody> </table>		Inflectional Morphology	Derivational Morphology	①	The study of the modification of words to fit into different grammatical contexts.	The study of the formation of new words that differ either in syntactic category or in meaning from their base.	②	Create new forms of same word.	Creates new words.	③	Run - Running VB            VB	Danger → Dangerous Noun            Adj. (JJ)	
	Inflectional Morphology	Derivational Morphology													
①	The study of the modification of words to fit into different grammatical contexts.	The study of the formation of new words that differ either in syntactic category or in meaning from their base.													
②	Create new forms of same word.	Creates new words.													
③	Run - Running VB            VB	Danger → Dangerous Noun            Adj. (JJ)													

Total Marks of Question no.		Examiner 1	
Space for Marks	Question No.	START WRITING HERE	
		<p>Q1(g) Apply Porter's Algorithm —</p> <p>(1) Running → Run  <span style="margin-left: 10em;">(remove) ing → φ      Rule used  <span style="margin-left: 10em;">(*V k) ing → ε</span></span></p> <p>(2) King → King  <span style="margin-left: 10em;">motoring → motor  <span style="margin-left: 10em;">sing → sing</span></span></p>	
		<p>Q1(h) Define Regular Expression &amp; its usage</p> <p>Ans A regular expression is a language for specifying text search string. RE helps us to match strings or set of strings, using a specialized syntax held in a pattern.</p> <p>Usage in NLP — A RE search corpus function will search through the corpus returning all texts that contain the pattern.</p>	
		<p>Q2(a) Explain the stages of NLP.</p> <p>Ans Stages of NLP —</p>	
		<p>(1) Morphological Analysis — In this phase individual words are analysed and distinct words are identified depending on the morphemes. Stemming &amp; lemmatization are used as part of text preprocessing to identify the tokens (words)</p>	
		<p>(2) Syntactic Analysis — In this phase linear sequence of words are transformed into structure</p>	

Total Marks of Question no.		Examiner 1
-----------------------------	--	------------

Space for Marks	Question No.	START WRITING HERE
		that shows how the words relate to each other. In this phase Part of speech tagging takes place.
3)		<u>Semantic Analysis</u> - In this phase a mapping is made from the input text to an internal representation that reflects the meaning. Here word sense disambiguation using wordnet/babelfox is performed.
4)		<u>Pragmatic Analysis</u> - To interpret what was said to what was actually meant.
5)		<u>Discourse Analysis</u> - Using Habb's algorithm resolves references.
Q2(b)		Explain FST in detail with example. Ans Finite State Transducers (FST) is a sophisticated tool that is used for understanding & transforming language. It FST is used in various applications like auto correcting, search engine. It is computational model used for representing & manipulating finite state machine (FSM) that map input sequence to output sequences. Key component of FST are - state, transition, input symbol, output symbol (annotation), & FSM

Total Marks of Question no.		Examiner 1	
Space for Marks	Question No.	START WRITING HERE	
		<ul style="list-style-type: none"> <li>- One common application of FST in NLP is Morphological analysis (which involves analysing the structure &amp; meaning of words at the morpheme level)</li> <li>- Stemming with FST - Reducing words to their root form often by removing the affixes.</li> <li>- Eg - English stemming with an FST - Let's consider an English stemming FST will have states representing the process of removing common English suffixes.</li> </ul>	
		<p><u>Step 1</u> - Define the FST states &amp; transitions</p> <p>States by defining the states of the FST representing different stages of stemming</p>	
		<ul style="list-style-type: none"> <li>- Define transitions between states based on rules for removing suffixes</li> <li>- E.g transitions</li> <li>- <u>state 0</u> - Initial state</li> </ul>	
		<p>Transition - If the input ends with "ing" remove "ing" &amp; transition to state 1</p>	
		<ul style="list-style-type: none"> <li>- State 1: "ing" suffix removed</li> </ul> <p>Transition: if the input ends with "ly" remove "ly" &amp; transition to state 2</p>	
		<ul style="list-style-type: none"> <li>- State 2: "ly" suffix removed</li> </ul> <p>Final state : Output the stemmed word.</p>	

Space for Marks	Question No.	START WRITING HERE

- Q2(c) Write a note on Language model
- A language model is the core model component of NLP. It's a statistical model that is designed to analyse the pattern of human language & predict the likelihood of a sequence of words or tokens.
- Language model determine the probability of the next word by analysing the text in data. These models interpret the data by feeding it through algorithm.
  - For training a language model, a number of probabilistic approaches are used. These approaches vary on the basis of the purpose for which a language model is created.
  - Types of language model -

- ① Statistical Language model -
- It includes the development of probabilistic models that are able to predict the next word in the sequence, given the word that precedes it.

Total Marks of Question no.		Examiner 1			
Space for Marks	Question No.	START WRITING HERE			
	a)	<p><u>N-gram Model</u> — In this a probability distribution for a sequence of <math>n</math> is created, where '<math>n</math>' can be considered as size of gram (a sequence of words). So if we use <math>N</math>-gram model, so it predicts the next word in the sequence based on <math>N-1</math> words.</p>			
	b)	<p><u>Neural Language Model</u> — It overcomes of <math>N</math>-gram model &amp; are used for complex task such as speech recognition or machine translation.</p>			
<h3>Common Challenges of NLP Language Model</h3>					
<p>① Long term dependency — Capturing relationships between words that are far apart in a sentence. Traditional RNN face vanishing gradient problem.</p>					
<p>② Low Resource Language — Building effective model for indigenous language is challenging due to limited data.</p>					
<p>③ Difficult to handle sarcasm, irony, idioms.</p>					
<p>④ Handling noisy text</p>					
<p>⑤ Contextual Ambiguity</p>					

Space for  
MarksQuestion  
No.

START WRITING HERE

Q5(A) Find the probability of the following sentence  
<s> Michael & Zack played at the playground</s>  
from the following corpus. Assume a trigram language model (use Laplace smoothing)

<s> The school was open </s>

<s> Michael & Zack went to the school</s>

<s> The Playground at the school was huge</s>

<s> Bob, Michael & Zack were friends </s>

also find the perplexity of the sentence  
which has highest probability.

Soln

Trigram Model Probability

Probability Table

$$P(\text{Michael} | \langle s \rangle) = \frac{1}{5}$$

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

$$P(\text{Zack} | \text{Michael} \&) = \frac{2}{2}$$

$$P(\text{Played} | \text{and Zack}) = \frac{1}{3}$$

$$P(\text{at} | \text{Zack played}) = \frac{1}{1}$$

Space for  
MarksQuestion  
No.

START WRITING HERE

$$P(\text{the} | \text{played at}) = \frac{1}{1}$$

$$P(\text{playground} | \text{at the}) = \frac{1}{2}$$

$$P(\langle x \rangle | \text{the playground}) = \frac{1}{2}$$

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

$$\boxed{P(s) = \frac{1}{60}}$$

$$\text{Perplexity} = P(s)^{-1/n} \quad n = \text{tokens} = 9$$

$$\Rightarrow \left(\frac{1}{60}\right)^{-1/9}$$

$$\boxed{\text{Perplexity} = 1.576}$$

Total Marks of Question no.		Examiner 1	
Space for Marks	Question No.	START WRITING HERE	

(Q3(a)) Discuss different Text Preprocessing techniques with example.

Ans Text preprocessing refers to a series of techniques used to clean, transform & prepare raw textual data into a format that is suitable for NLP or ML task. The goal of text preprocessing is to enhance the quality & usability of the text data for subsequent analysis or modelling.

- It involves following steps -

- 1) Lowercasing
- 2) Removing punctuation & special character
- 3) Stop word removal
- 4) Removal of URLs
- 5) Removal of HTML tags
- 6) Stemming & Lemmatization
- 7) Tokenization
- 8) Text normalization

① Lowercasing : `text.lower()`

② Removing punctuation & special character using NLTK & RE library  
`re.sub(punctuation_pattern, "text")`

③ Stop word removal can be done by importing `stopword` from `nltk.corpus import stopwords`  
`remove_stopword(en_text, "english")`

④ Removal of URLs using Regular Expression RE  
`re.compile(r'https?://(s+www\.)|(\s+)')`

Space for  
MarksQuestion  
No.

START WRITING HERE

(5) Stemming & Lemmatization — Stemming removes common suffixes from the end of word tokens, lemmatization ensures the output word is an existing normalized form of the word. Running — run  
Smiling — smile

(6) Tokenization — It is the process of dividing the text into smaller units known as tokens. The word tokenization divides the text into individual words. Sentence Tokenization — This is useful for tasks requiring individual sentence analysis or processing.