

## NLP Assignment - 01

DATE: 11/09/2023

Q-1 Explain with example different types of ambiguity in natural language?

Ans

1) Lexical ambiguity:-

This occurs when a word has multiple meanings. For instance, the word 'bank' can refer to a financial institute or the side of a river.

2) Syntactic ambiguity:-

This arises when a sentence can be parse in multiple ways due to its structure. For example, 'I saw the man with the telescope' can mean you saw a man with the telescope or you used a telescope to see a man.

3) Semantic ambiguity:-

This occurs when a sentence can be interpreted in more than one way because of the meaning of the words used. For instance, 'She found a bat in the park' could refer to a sport equipment bat or a flying mammal bat.

4) Pragmatic ambiguity:-

This arises when the meaning of the sentence is unclear due to the context or tone. For example, 'Can you pass the salt?' could be a request or a question about someone's ability?

Q.2 Discuss challenges of the NLP and how to overcome them

Ans Natural language processing faces several challenges due to the complexity and nuances of human languages. Here are some major challenges and potential ways to overcome them:

a) Ambiguity:-

Ambiguity in a language can lead to confusion. One way to tackle this is by using context based approaches. Machine learning models that consider surrounding word or sentence can often disambiguate the intended meaning.

b) Variability and Informality:-

Language is incredibly diverse with slang, colloquialisms and variation across cultures. Incorporating more extensive and diverse training data can help models better understand and generate different language styles.

c) Lack of context:-

Understanding context is crucial for accurate comprehension. Models like transformer based architecture leverage attention mechanisms to capture context over longer sequences of the aiding in overcoming challenge.



Q.3

Ans

Tokenization is a fundamental preprocessing step in natural language processing (NLP) that involves breaking down a text into individual words or units called tokens. These tokens serve as the building blocks for the various NLP tasks, such as language modeling, text classification and machine translation.

Tokenization is important for many reasons.

a) Text understanding:-

Tokenization divides text into manageable units allowing machine to understand the structure of sentence, paragraphs and documents.

b) Feature extraction:-

In this task like text classification each token can be treated as a feature, enabling models to learn patterns and relationships within the text.

c) Language modelling:-

Tokenization aids in the building language models by breaking text into smaller units for predicting the next word or generating coherent text.

Code:

```
import re
text = "Tokenization is good"
tokens = re.findall(r'\b\w+\b', text)
print(tokens)
```

Q-4

## Stemming

Stemming is the faster because it chops word without know the context of the word in given sentence.

It is a rule-based approach

Accuracy is less

For example  
studies = 'studi'

## Lemmatization

Lemmatization is slower as compared to stemming but it knows the context of the word before processing.

It is a dictionary-based approach.

Accuracy is more

For example  
studies = 'study'.

## Porter Stemmer algorithm:-

The porter stemmer algorithm is widely used stemming algorithm developed by Martin Porter in 1980. It follows a set of heuristic rules to transform words into their stems. The algorithm consist of several phases, each together targeting specific suffixes. The phases aim to progressively remove common suffixes simple words to their root forms.



Q. 5

Explain Edit distance algorithm

Ans

The edit distance algorithm measures the similarity between two strings by calculating the minimum numbers of single character edit required to transform one string into the other form.

The edit can be insertion, deletion, or substitutions. It is widely in various applications such as spell checking, DNA sequence analysis and plagiarism detection.

Algorithm

function editDistance(str1, str2):

$m = \text{length of str1}$

$n = \text{length of str2}$

    create a dp array of the  $(m+1) \times (n+1)$

    for  $i$  from 0 to  $m$ :

$dp[i][0] = i$

    for  $j$  from 0 to  $n$ :

$dp[0][j] = j$

    for  $i$  from 1 to  $m$ :

        for  $j$  from 1 to  $n$ :

            if  $str1[i-1] == str2[j-1]$

$dp[i][j] = dp[i-1][j-1]$

            else

$dp[i][j] = \min(dp[i-1][j] + 1, dp[i][j-1] + 1, dp[i-1][j-1] + 1)$

return  $dp[m][n]$

Forexample:-

cat      and      hat

∴ edit distance is 1.  
substitute the 'c' with 'h'.

Q.6  
Ans

Use of finite automata and finite state transducers.

1) Tokenization.

FST can be used for segmenting text into tokens. For example in the english you can use an FST to split a sentence into individual words or a sub word.

2) Part of speech tagging:-

FST can help assign parts of speech to words in a sentence. By constructing an FST that recognizes word sequences and assigns pos tags, you can develop a pos tagging system.

3) Morphological analysis

FSTs are extensively used in morphological analysis to break down words into their constituent morphemes.

4) Spell checking and correction:

It can be used to develop a spell checking and correction system.



## Q.7 Importance of Morphological analysis.

- 1) Language Understanding  
Understanding the structure of words is fundamental to understand language.
- 2) NLP tasks:  
It is integral to various NLP tasks, including machine translation, text summarization, information retrieval and sentiment analysis.
- 3) Resource-Sparse Language  
In languages with rich morphological complexity morphological analysis is vital for accurate processing.
- 4) Search engines  
For search engines morphological analysis enables searching for variations of words.
- 5) Sentiment analysis:-  
In sentiment analysis, morphological analysis helps identifying negation or words forms that change the polarity of a sentence.

### Applications

1. Information Retrieval
2. Machine translation
3. Sentiment analysis
4. Speech recognition.

Q.8

N-grams.

Ans

Ngrams are contiguous sequences of  $N$ -items words, letters or symbols from a given sample of text or speech. They are a fundamental concept in the NLP. N-grams are used to analyze and model the relationships between elements in a sequence. They are classified as:

- 1) Unigrams (1-grams):  
These are single items in a sequence / often individual words.
- 2) Bigrams (2-grams):  
These consist of pairs of adjacent items in a sequence of words.
- 3) Trigrams (3-grams):  
Trigrams are composed of the consecutive items.
- 4) N-grams ( $N > 3$ ):  
N-grams can be any sequence of the  $N$  items.

Example:-

This is an group of words.

Unigrams :- 'This', 'is', 'an', 'group', 'of', 'words'.

Bigrams :- 'This is', 'an group', 'of words'.

Trigrams :- 'This is an', 'group of words'.



Q.9

Penn Treebank

Ans

Penn Treebank Pos tag set is a widely used tagging scheme for English that assign tags to words based on their grammatical roles in a sentence.

Following are the some of the tags.

- 1) CC (Coordinating conjunction):  
connect words phrases or clauses.
- 2) CD (Cardinal number)  
Represents numbers like 'one', 'two',
- 3) DT (Determiner)  
modifies the noun and specifies their references such as the 'this', 'some'.
- 4) (JJ) adjective  
Describes noun like, 'happy', 'blue', 'tall'.
- 5) NN (Noun, singular)  
Represent singular or mass noun like 'cat'.
- 6) RB (adverb)  
modifies verb, adjective or other adverb.
- 7) To (to)  
Part of the infinitive form of the verbs like 'to' 'in' 'to go'.

## Q.1a Rule based part of speech (POS) tagging

- 1) Tokenization  
The first step is to tokenize the input text into individual tokens.
- 2) Linguistic Rules  
Develop a set of linguistic rules that describes the characteristic and patterns.
- 3) Ambiguity Handling  
Implementing rule based tagging may lead to ambiguity rules should be designed to handle such cases.
- 4) Exception Handling :-  
Rules should account for exceptions and irregularities in the language.
- 5) Iterative improvement  
This involve analyzing.