# SYMBIOSIS INTERNATIONAL (DEEMED UNIVERSITY)

(Established under Section 3 of the UGC Act 1956)

Re-accredited by NAAC with "A++" Grade | Awarded Category - I by UGC

No. **075662**

**SIU-24 E**

| | |
|---|---|
| Verified all entries & found correct | |
| Jr. Supervisor's Signature & Date | |

Centre: _SCMHRD_

Seat No. (in figures) | 2 | 0 | 4 | 0 | 0 | 7 |

Seat No. (in words) _Two Zero Four Zero Zero Seven_

Examination : Day & Date : _Tuesday - 22/10/2024_

Programme : _Business Analytics_ Semester : _III_

Course : _Natural Language Processing_

Main Ans. Script 1 + No. of Supplementary Answer Scripts _____1_____ = Total | 2 |

| Q.No. | | | | | | | | | | | | Total | Signature of Examiner |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Max. Marks | | | | | | | | | | | | | |
| Marks Obtained | | | | | | | | | | | | | |

## INSTRUCTIONS

1. Mention your details only in the space provided for in the main answer script & the supplement. If any other details (including seat number, name, contact details, etc.) are written anywhere else in the answer script and or supplement it will be treated as adoption of unfair means and the performance will be treated as null and void for the entire examination.

2. Write answer in legible hand. Answers written in an illegible and undecipherable hand are liable to be marked as zero.

3. An act of copying or of impersonation at an examination is punishable under the Maharashtra Prevention of Malpractices at University, Board and other specified examinations Act 1982.

4. Candidates should write answers in BLUE/BLACK ink only. Use of Pencil and other colors are permitted only in case of diagrams, graphs etc. Answer Scripts written with pencil or ink of other colors will not be evaluated.

5. Write on both sides of paper.

**BEGIN WRITING HERE :**

**Q1.** National Language Processing (NLP) is a branch of computer science that make use of computational linguistics and statistics to understand the human language by the machines.

It has taught the machines to understand, process and generate human language.

Computer Science, Advanced Modelling & Computational Linguistics make NLP.

Advanced Modelling can be Statistical way or by using Artificial Intelligence like Machine learning.

It consist of National language Understanding & National language Generation.

In Natural language understanding, machine tries to break down the input text to understand the structure, meaning & context of the language

Some of Real world applications are:

Text Classification: the input text is classified based on common features.

Eg: SPAM Detection

Sentiment Analysis: The sentiments of the text such as positive, negative or neutral can be extracted by finding the sentiment words (positive,negative) & scoring the sentence, document based on scores.

Benefits:

Helps in understand human language analyze & extract the visible & hidden meanings. Automate & helps in processing & understanding volumes of text.

This process will help in model creation.

This also helps Machine understand Human language.

Challenges:

☆ Ambiguity in language, Same words having different meanings.

☆ Complex languages & huge volumes of data required for accurate model.

☆ Cannot recognize Sarcasm and other subjective aspects of communication

Natural language Generation
Machine generates the human language up
to interact with the humans.

Some Real world Applications

Machine translation: The translation of
text from one language to other

Content Generation: The content is generated
based on given prompt

Benefits:
- Reduces the manual work of
humans, used in text summarization
- GPTs make content generation
useful to interact with humans.
- chatbots to generate a offers
to human queries.

Challenges:
- Depends on the training dataset.
- It may contain baisness in answers/replies
- Complex task requires computational
power.

Q.2. Commercial uses of NLP:
- Text-speech conversion ⎫ Speech
- Speech-to-text conversion ⎰ recognition
- text classification
- Sentiment Analysis
- Machine translation
- word or sentence prediction (generation)
- Text summarization
- Speech recognition
  language or document classification

**Text-speech or Speech-text conversion:**

NLP is used in conversion of 'text' of a language into speech using the inbuilt models of computer science.

In Python, the following inbuilt language models are used for text-speech conversion.

NLTK → Natural Language toolkit, Spacy, Gensim are the other python packages that contains inbuilt models for text Analysis.

The broken down text is converted to digitized form into numbers then to audio signals, then passed to speaker for audio output.

Application in voice recognition, voice enabled automation.

Here the human speech is converted to digital signal then to text then processed. The commands of voice is then executed.

Eg: Siri & Alexa.

Python Package 'Sphinx' can do speech recognition.

## Documents classification:

Documents which are collection of text into separate files can also be classified & categorized using nlp

Example, Content of the document is rated, plagirism check, etc.

In plagiarism checker the words & sentence are checked with Documents found in Internet, if it matches then it flags & classifies as "plagirism detected".

## Text summarization:

Text corpus is summarized to form the idea behind the text. It will help reduce the long text into short & precise format help us to label the group of text.

Q.3 Significance of POS tagging:

Parts of speech tagging used to tag the words to the grammatical categories, as Nouns, pronouns, verbs, adverbs, adjectives, Interjections and conjuction.

It's a disambiguation method used to address one of the major chalenges of NLP - Ambiguity.

The words and its context & role in the sentence is recognized then ambiguity is cleared.

Example:

"I want to fly, like a fly!"

Here without POS tagging the word 'fly' is considered as a same entity though it has 2 meanings an insect fly and action fly.

Now when we do POS tagging using ~~NLP~~ or NLTK or spacy, we get.

I — noun.
want — verb          Here the
to — preposition.    'fly' is
fly — verb           seperated to
like — verb.         verb &
a — Articles         noun.
fly — noun

Thus ambiguity is rectified

POS also helps the data to be applicable for NER (Named entity recognition)

Since the Noun - fly and Verb-fly are seperated, the NER will be more accurate.

Example:

Mr. plant's office is near the biogas plant, where biofue was planted
Here Mr.plant → noun.

biogas plant → noun

planted → Verb

Hence NER since tagged from POS recognizes Mr. plant as a name

POS also categorizes the word = (word category tagging) used to further syntactic Analysis

Q4. N-grams are sequence of n items such as words, letters and base pairs from sentences.

- when N=1, unigrams

Eg: 'NLP', 'is', 'an', 'interesting', 'Subject' → unigram

when N=2, Bigrams

Twin towers is a large building. Twin towers was constructed in 19th Century.

Here Twin towers → Bigrams as they always appear in a sentence together

Also, when we divide the words in a sentence in a pair of two or three or more, n-grams concept is used.

for same example:
      NLP is an interesting subject

('NLP', 'is', 'an', 'interesting', 'subject') unigram

('NLP', 'is'),
('is', 'an'),                    } Bigrams
('Interesting, subject'),
etc.

N-grams come under statistical modelling of language modelling where the next word or sequence of words are predicted using their probability of occurrence that comes from frequency of words.

Limitations:
N-grams cannot be applied for huge dataset as it makes the computations more complex.

Approach:
The next word prediction using n-gram here unigram,

('NLP', 'is', 'a', 'good', 'Subject') unigram.

P('NLP' | 'is', 'a', 'good', 'subject')
probability of occurence of 'NLP' given the sequence of sentence & (is, a good subject).

Here we say, the model then trained for new data it also creates baisness in data as the training happens due to proximity of words

Cannot work on real life test data, unless trained on huge volumes of datasets

**Q5.**

| Consistency parsing | Dependency parsing |
|---|---|
| ☆ It identifies the constituents of the sentences from words. Into noun phrases verb phrases | ☆ It helps to identify the relation (or) Grammatical relation among the words in a sentence. example. subject-verb relation |
| ☆ Used in Natural Language Understanding | ☆ Used in Natural Language generation |
| ☆ The tree (phrase tree) follows top down approach flows from root node to leaves | ☆ The phrase tree follows bottom up approach. The flows from leaves to root node. |
| ☆ It follows tree structure with nodes → root node & leaves. | ☆ It follows graph structure where nodes - words edges - relations. |
| ☆ Extracts complex information & characteristics of words, of a language. Thus Rich Morphological language of rich syntactic structure can be used. | ☆ extracts simple information & characteristics thus, poor morphological & poor syntactical structure languages like 'English' 'Chinese' are used. |
| ☆ Old traditional model Approach | ☆ new Advanced model approach (machine Learning) |

**B6.** Application of Naive Bayes classifier:

It is supervised Machine Learning algorithm of probabilistic classification modelling make use of Bayes theorem to predict the classes of target variables. Used in classification. more suitable in text classification.

Usage:
* Naive bayes used in high dimentional features of data.
* Speed of prediction is high
* Process & predict complex relations
* features should not be conditionally dependant & data must not contain missing values
* The features should be equally weighted or given importance.
* f continous features should be normally distributed & discrete features should be multinomial distributed.

Text classification with an example:

multiple applications include:
* Text classification
* Document Classification
* content rating
* SPAM Detection.
* Sentiments classification.

SPAM Detection:
Here the Naive bayes classifies the new email or SMS as a 'spam' or 'not spam', using training dataset

# SYMBIOSIS INTERNATIONAL (DEEMED UNIVERSITY)

(Established under Section 3 of the UGC Act 1956)

Re-accredited by NAAC with 'A' grade (3.58/4) | Awarded Category - I by UGC

Supplement No. SIU-19 A

04432

Centre: SCMHRD

Seat No. (in figures) 2 | 0 | 4 | 0 | 0 | 7

Seat No. (in words) twozero fourzero zero seven

Programme: BA

Course: NLP

Semester: IV

Supplementary Answer Script No. I / II / III

---

**BEGIN WRITING HERE:**

assuming the features of spam class

i) Unknown resource → yes, no

ii) suspicious email → yes, no

iii) Occurance of words → spam → 'Discount', 'offer'

↳ Discount, offer, lottery,

So, using Bayes theorem, the new email has the following

P( Spam )

features (Unknown, suspicious, account)

test ⇒ ( yes, yes, offer) ⇒ Neutral
data

then

$$P\left(\dfrac{Spam}{\substack{new \\ email}}\right) = \dfrac{P\left(\dfrac{Unknown}{Spam}\right) \times P\left(\dfrac{suspicion}{Spam}\right) \times P\left(\dfrac{Place}{Spam}\right) \times P(Spam)}{P(new email)}$$

P(new email)

Similarly

we also calculate for    $P\left(\dfrac{Not\ spam}{new\ email}\right)$

Finally in this case   $P\left(\dfrac{Spam}{new\ email}\right) > P\left(\dfrac{not\ spam}{new\ email}\right)$

Thus it identifies as SPAM

Similarly based on occurence of positive & negative words, sentiment of document is classified as positive or negative.

Q.1 Steps in text-preprocessing :

1. Tokenization:

The text preprocessing first starts with,
i) Document triage
ii) Text extraction/tokenization

Document triage:

The documents are converted to digitized format that can be processed by the computers software for Analysis & stores them into corpus

Text extraction:

The text data in documents are extracted & stored in machine or software accessible format) into object & variables.

## Sentence tokenization:

The sentences are extracted & stored in a list by breaking down a document or passage using delimiters.

### Importance:

Organizes the text into suitable format for text preprocessing. Breaking down also helps to analyze sentence use or construct helps for further text analysis.

### Word tokenization:

Breakdown of sentence into words and letters.

## 2. Lemmatization & stemming:

Reduces the words into its base form that are meaningful words for lemmatization.

Example:
- fly → fly
- flies → fly
- flown → fly
- flying → fly.

**Stemming:** If it reduces the words removing suffixes to meaningless words knowns stemming.

- Studying → studi

**Importance:** Reduces the complexity in words & different forms of words into their base form, reduces vocabulary & errors.

3. Stopwords removal:

Remove filler words &
insignificant words from the database.

Importance: Removes words that does not
add any meaning, reducing volume
of dataset.

Punctuations can also be removed in
similar way.

4. POS & NER tagging:

The words then categorized
based on parts of speech and entities
like name, place, & location can
be extracted Using Named entity
recognition

Significance of text preprocessing:
Removes Bais & variance improves
performance of model

Q.8. The main components of NLP

(NLU) Natural language Understanding:
It understands the human language
by breaking it down to machine readable
& recognizable structure.

This stage consists of meaning extraction
& understanding role & structure of words
in sentence.

POS tagging & parts of speech tagging
helps to identify words as parts of
speech securing ambiguity

Semantic Analysis used to extract meaning
& context.

Syntactic Analysis is used to extract the grammatical structure.

Natural Language Generation (NLG) where the model generates the language or text from the understanding of input data.

Text planning:
Plans the text - words that should be present in language generation

Sentence planning: The sentence generation sequence & order of words placed

Realization: The structure of sentence, words put into language that conveys meaningful message & generated

6.9. Language Modelling:
It is the probabilistic method of assigning probabilities to words in a sentence based on count or frequencies used to predict the next probability of next word or sequence of words generated

Two Methods
- Statistical Method
- Neural Method

Statistical Method: Make use of probability to predict next word found sequence of words or input string

Use the concept of n-grams are used. N-grams are sequence of n items.

assuming our
n-grams can be bigrams, trigrams etc..

Now, Probability of sentence = P(w)

Probability of words = P($w_1, w_2 ... w_n$)

$$P(w) = P(w_1, w_2, ... w_n)$$

P(W) & P($w_1, w_2, ... w_n$) is Language
modelling.

<u>Using Markov, we</u>

Neural Model :- Neural method
It overcomes the limitation
of statistical method, by using
advanced a complex modelling
done using Machine learning
or deep learning. But it also
needs huge computation power & huge
dataset to then.

6.10. Languages like Chinese and
Arabic does not have
structure (or) order of words
in a sentence.
The order of words in a sentence
does not change the meaning of
sentence.

Out of them these are Orthographic

Language.

EXAMPLES

(Rough ALL)

SIU SIU SIU SIU SIU SIU SIU SIU SIU SIU SIU SIU SIU SIU SIU SIU

12

① text class content generation

→ machine translation

→ sentiment

→ chatbots

→ Spam Detection

→ speech recognition

*Natural Lang...

⑧ NLU | NLG

NLP

⑦ Context things...

→ sentence → word → text

→ SEO

② Logographic alphabetic syllabic