



SYMBIOSIS INTERNATIONAL (DEEMED UNIVERSITY)

(Established under Section 3 of the UGC Act 1956)

संयुक्त विश्वविद्यालय

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

Total Pages 12 (Twelve)

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

No. 075785
SIU-24 E

Centre: SCMHK D

Seat No. (in figures) 2 0 4 0 0 5

Seat No. (in words) Two Zero Four Zero Zero Five

Examination : Day & Date : 22nd Oct '24, Tuesday

Programme : MVA - IT 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.



1. Natural Language Processing :- (NLP)

The two real world applications are as follows:-

Here the NLP helps in understanding the essence or sentiment of the market or public, by ~~using~~^{processing} the text or document.



as an input and giving out the sentiment which can be positive or negative.

e.g. Understanding the sentiment of public or customer towards a new launch of product.

Benefits:-

- This helps in understanding the performance of the product in the market.
- This will help the business in strategizing the actions to be taken in future.

Challenges:-

- NLP ^{may} find difficulty in categorizing the sentiment clearly due to complex feedback or text.
- The processing of the text or document given can take time and also can be computationally intensive.

(5) Speech Recognition.

→ NLP can analyse the real on going speech.

For instance, the analysis of Federal Reserve President speech so that it can help in prediction of impact of the speech on the stock market.



Benefit:-

- Helps in strategizing the investment and trading position on per the speech given.
- It saves time as the real impact of the speech can follow time.

Challenges:-

- If the speech is in different language, then NLP can find difficulty in converting it into local language.
- Phonem. or sign used in speech, is difficult to understand by NLP.

2. The few commercial uses of NLP are as follows:-

(a) NLP is used in understanding the sentiment of the customers about the new launch of product.

This helps in depicting the performance of the market campaign and also to understand the real time feedback.

Accordingly, the organization can strategize the action to be taken in future.

(b) Analyzing the Ancient Symbols or Pictorial

NLP has great contribution in archaeology field. The old script



written on the scriptures are different to analyse.

Now the machine can use suitable approach to analyse and get the maximum information from the scriptures of script.

(c) Machine Translation:-

NLP has been greatly leveraged in translating the speech from one language to another.

For instance, when there is an interaction of people from different culture or background then they can communicate to each other in real time, in their very own mother tongue.

(d) Chatbots

NLP is being leveraged by the organization to build the chatbots. These chatbots are responsible to answer the query of the customer in real time.

This requires machine to understand the human language, break down into text and analyse the question so that it can properly give the feedback in human language.

This saves a lot of time and money for the organization in terms of accuracy and manpower optimization.



3. POS Tagging:-

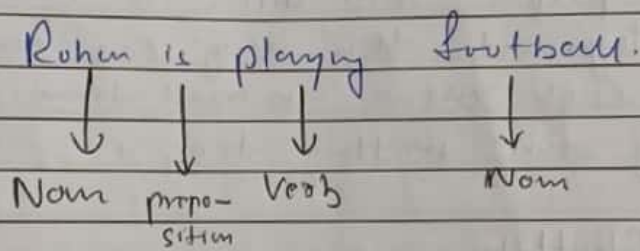
Parts of Speech tagging is the process of analysing the text as per grammar rule.

Synificance:-

This helps in breaking down the ~~words~~ ^{Sentences} to understand the parts of speech of the words belong to.

This bifurcates the words into Noun, Pronoun, Adjective, Adverbs, Preposition, Conjunction and interjection.

Eg.



- This breakdown of the text, helps in understanding the meaning and the action happening.
- If the POS tagging is not done, then machine can confuse with the action, as object get subject relationship.
- POS tagging also helps in for the subsequent steps of NLP.
Once the speech and grammar understanding is clear for the machine, it can now relate or analyse the relationship between the different sentences, and also relate it with the real world.



4. N-gram Models:-

N-gram are the sequence of words in the given text.

For example.

The article is good.

Bigram \rightarrow "The article", " is good"

Unigram \rightarrow 'The', 'article', 'is', 'good'

Language Modelling:-

In language modelling, we try to predict the next word or continue on from the given text.

\rightarrow Now, language modelling, works on the principle of probabilistic modelling, where it try to find the probability of the next word to occur in the text.

\rightarrow Here, N-gram model helps in a way that the sequence of word ~~present~~ present determines the probability of next word or words. This is based on Markov or hidden Markov Model.

For instance.

Playing is good for health.

In this, ~~the~~ after seeing the word like 'is' and 'for' which ~~does not~~ not contribute much.

'Playing' and 'good' determines the probability of



the next word to be 'be' 'hearts'.

Limitation:-

→ There are words whose probability can be zero. And incorporating this word can disturb the overall analysis.

So, we use various smoothing techniques to tackle this problem. The methods like Laplace, Additive, Katz methods etc are used.

→ In N-gram model, the sequence of words can be such which has least contribution in ~~the~~ analysing the sentiment of the document or text.

So this creates issue in building the accuracy of the model.

5. Constituent Parsing.

Dependency Parsing

① In this, the text or document is segmented into constituents or group of words as per grammar rules.

In this parsing, the aim is to find the relationship between the words as per the grammar rule.

② Parse tree formed here is top-down approach.

Parse tree formed by bottom to top approach.



G. Naive Bayes Classifier:-

This works on the principle of Bayes' theorem. It is basically a mathematical model which works on the probability theory.

$$P(A/B) = \frac{P(B/n) \times P(n)}{P(B)} \rightarrow \text{Baye's Theorem}$$

- In text classification, the Naive Bayes Classifier helps in classify the overall sentiment of the text into categorical outcomes.
- For instance, NBC is used in classifying the document or text into sentiments like positive and negative.
- Suppose we have a collection of feedbacks in a raw format from the customers who have used our new product launch.
Now, the raw text can be feeded in the model to classify the overall sentiment of the market segment into positive or negative.
- Another application can be analysing the speech of the financial minister on the "Budget Day" to understand and categorize the sentiment of the budget into →
 - (1) Bullish
 - (2) Bearish
 - (3) Neutral

7. Text Pre-processing

It is one of the crucial step in the NLP Pipeline.

Before the text Pre-processing, the data is collected and stored properly. Then the following steps are taken to pre-process the text before it is exposed to feature extraction and modelling:-

So, Text Pre-processing is all about breaking down the text and ~~analysing~~ ~~data~~ converting it into smallest units.

① Tokenization:- It is process of breaking down the text into smaller units.

It can be done by word tokenization or sentence tokenization.

Two main steps in

(a) Word tokenization:- Dividing the ~~sentences~~ ^{sentences} into words.

(b) Sentence tokenization:- Dividing the sentences or paragraphs into individual sentences.

(c) Stemming:- In this the ~~entire~~ ^{entire} word is broken down into its smallest unit by removing the prefixes or suffix.

Eg: Runny → Run
Playing → Play

i.e. removing 'ing', 'ly', 'es', etc.



(d) Lemmatization:-

In this the words broken down into its root word.

For example,

Runing → run

Run → run

Studying → study.

Note:- Lemmatization uses corpora or dictionary to get the root word and it is more accurate than stemming. But stemming is faster.

(2) POS tagging:- This segments the words as per grammar rule.

For example:-
Rohan is playing football
↓ ↓ ↓ ↓
Noun preposition Verb Noun.

(3) Text Segmentation

(4) Dependency Parsing

Importance of Text Preprocessing:-

(1) Text preprocessing is helpful in a program that it makes 'easy' for the machine to convert the human language into digit form so that it can process the text.

(2) It gives the machine the understanding of feature for modeling.



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

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

Supplement No. SIU-20 A

58220

Centre:

SCHMRD

Seat No. (in figures)

2 0 4 0 0 5

Seat No. (in words)

Two Two Four Two Two Five

Programme:

WBA-BA

Semester:

III

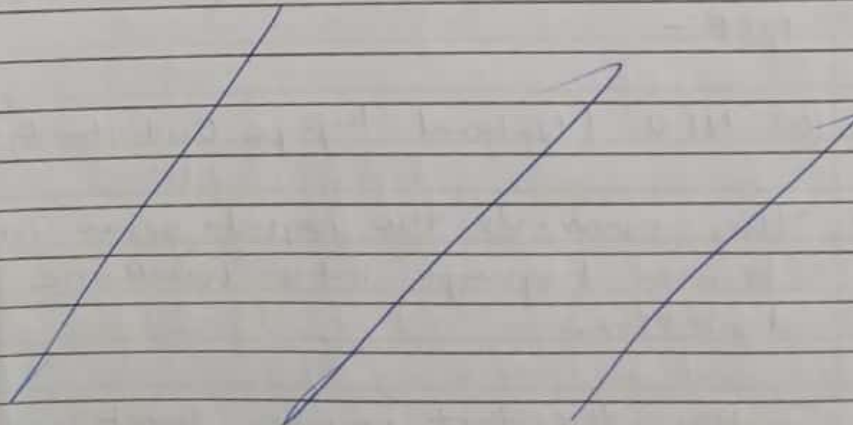
Course:

Natural Language Processing

Supplementary Answer Script No. I/II/III

BEGIN WRITING HERE :

- (3) It enhances the accuracy of the model going forward.
- (4) Feature extraction becomes easy if the preprocessing is done better.



P-7.0



सत्यमेव जयते

SIU SIU

8587

8. NLP (Natural Language Processing)

The purpose is to make machine understand the human language, and so that it can be used for various applications.

There are two basic components of NLP:-

(a) NLU (Natural Language Understanding)

→ This converts the input given in human language into machine digit form.

→ In this the text given as input, are assigned digits and this pattern formation is used by the machine for tasks assigned.

(b) NLG (Natural Language Generation)

→ This converts the machine interpretation into human language. The generated output is received by the user.

For instance,

→ The Use of Chatbots for Question-Answer.

Chatbots work on these two components discussed above.

When the customer put his/her query to the bot, the NLU translates into machine language and NLG generates the final output as the answer to the customer.

9. Language Modeling:-

Language modeling is used to predict the next word or text or sentence as per the given document or text.

The different methods of language modeling are as follows:-

(1) Statistical Models:-

This model works on the principle of probabilistic model. The next word or sentence is predicted on the basis of the previous state of the sentence.

$$P(w_{n+1}) \approx P(w_1 w_2 w_3 w_4 \dots w_n) \approx P(w_n)$$

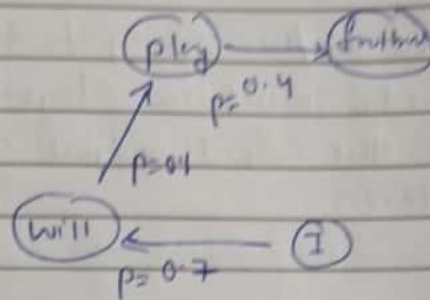


→ In the Markov Model, the future state of the word depends on the current state.

→ In this the probability is generated for the next word to happen.

For instance,

I will play football:



→ There is a transition probability of going from state to another.

→ Models like N-gram Model, Naive Bayes Classifier are the examples of the statistical models.

② Neural Model:-

→ This is an advanced version of language model.

This doesn't consist of any labeling of the input text.

The machine develops the understanding of the text and try to figure out the relationship.

→ Recurrent Neural Networks and Convolutional Neural Networks are the examples of this type.

- This is further and more flooded with enormous amount of input.
- It usesy corpora to train itself and understand the language ~~the~~ so that it can predict.
- Chatbots and Answer-Question applications are few of the examples of this type of model.

10. Text Pre-processing:-

This is used in cleaning the input data by tokenization, pos tagging, ~~and~~ etc in order to pre-process b. for it is fed into the model.

Challenges in pre processing Chinese or Arabic can be:-

- ① The language texture is symbolic in nature. For instance, the symbol used in Chinese depicts the whole word. Also, in arabic, the texture is quite symbolic and logothermic in nature.

This makes difficult for the machine to apply the methods like tokenization, pos tagging.

- ② The texture looks very similar to each other. The words are very similar in look wise. This makes challenging for the machine to



differentiated from each other.

- (3) Arabic is also challenging in sense that it goes backward in reading.

Suggestion to Overcome:-

- (a) Use unilingual corpora, to understand the meaning and analyzing the text in Chinese and Arabic.
- (b) It is also important that the input data is properly collected and stored. Any discrepancy in the 1st step, can create more challenge for the machine to pre process.
- (c) Use dependency parsing, it is more suitable in this case.

Verified and certified & found correct
Dr. Gussone's & Gussone's & Gussone's

2/10/24

NO. 075639
SIU-24 E

Centre: SCMHPD

Seat No. (In figures)	2	0	4	0	0	6
-----------------------	---	---	---	---	---	---

Seat No. (in words) Two Two Four Two Two Six

Examination : Day & Date : Tuesday, 22/02/24

Programme : MBA-GA Semester : III

Course: Natural language Processing

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

[illegible]

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.

2

MLP:

2

i) Chatbots :

(iii) Speech to Text Generation :-

This is used to convert the speech [input] into a text, store it in a document.

* Accurate information response	* Elimination of typing long documents
* No extra/unnecessary response	* Saves a lot of time
CHALLENGES	
* Chatbot may not understand the correct context	* Understanding the text and as it may have semantic errors.
* Eg: Orders, concern by user cannot be properly detected.	* Eg: The dining looks good. When translating, it can type dining @1 dining → Notbook ↳ Malt. product.

Q2 Allp has a lot of commercial uses:

ii) Document Retrieval :

When using packages like "Duffy", queries can be used to retrieve documents from a database through sentence matching, regular expressions

When there are thousands of documents and the naming conventions involve numbers, finding the right document can be a hassle.

In such cases, NLP can be used.

(iii) Spell Check:

When the document contains typos, one had to go through every word to check the spelling, which takes a lot of time. This can be eliminated through NLP.

Eg: The technical competency requires a degree

CORRECTED: The technical competency requires a degree

Eg: The water born disease is contagious.

CORRECTION: The water borne disease is cryptosporidium



(iii) Grammar Check:

NLP understands the language through various corpora. So, it can identify the grammatical errors.

Applications such as Grammarly are developed on this principle.

Eg: I ~~see~~ a bird fly.

CORRECT: I am seeing a bird fly. OR I saw a bird fly.

Eg: I would have ate the cake tomorrow.

CORRECT: I would have eaten the cake by tomorrow.

(iv) Text Summarization:

Companies have a lot of requirements, SOP's to follow. Text Summarization allows to give a summary of the input text.

This helps the person in saving time such as HR's, legal team to save time in reading everything.

Applications such as Merlin are built on this and are used for a web-page summarization.

(v) POS Tagging

Parts of Speech Tagging is a crucial process in the steps of NLP.

It helps to identify each of the tokenised words by its POS form.

Eg Noun, Verb, Adverb, Adjective, etc



Steps in POS Tagging:

- (i) Sentence Tokenization: Word Tokenization - It breaks down paragraph into sentences and sentences into words.
- (ii) Investing language models such as Spacy, NLTK for Pre-processing.
- (iii) Stemming, Lemmatization to reduce the words to lemma (iv) root form with meanings.
- (iv) Apply POS tagging on the words.
- (v) Result Analysis to understand the results.

A simple break up of words give us understanding.

POS tagging provides a semantic network for the words. With POS tagging, we can understand the link between the words as it helps in defining the structure.

This step proves crucial to Named Entity Recognition and even Consistency Parsing.

Eg: I love cakes.

Word Tokenised Result = ["I", "love", "cakes"] → Have no

Eg: I → I → love → cakes
 Subject Verb Object Meaning

Eg: I / have / a flight / tomorrow
 Subject Verb Object complement



Q 5

N-grams are continuous n-words that help in prediction of next word, in language modelling.

N-grams can be of

- * Unigram - One word
- * Bigram - Two words
- * Trigram - Three words
- * N-gram - Many words

When there are multiple words in an N-gram, the complexity increases.

I saw a cat on a wall.

N-gram = [I saw a cat on a]
 Answer can be sofa, wall, table etc.

In this the answers are endless, the model will be confused.

To reduce this, we can use tri-gram, bi-gram

Bigram = [I saw] [saw a], [a, cat], [cat on], [on, a], [a,]

This reduces the complexity.

Limitations:

- (i) Data sparsity, the model will not have enough corpus to predict the outcome. This brings the probability to 0.
- (ii) Lack of corpus words add ambiguity
- (iii) N-gram works on Statistical approach, if probability is 0, the model becomes useless.

(iv) To overcome this, Smoothing techniques such as Laplace Smoothing, Additive Smoothing, Backoff, Interpolation, Christ & Gale, Smoothing are used.



Q 5

Consistency Paradox

Consistency Paradox
 In identifying the words in a

(iii) It is a Top-

(iii) It provides the words like the

(iii) Best suited for Natural language

(iv) Suitable for POS tagging, N

(iv) It is more for analysis, less for generation

(viii) It provides a output

(viii) Good for the generation

help in
distinguishing



Q 5

CONSISTENCY PARSING

DEFINITION

DEPENDENCY PARSING

Consistency Parsing is used to identify the consistent words in a sentence.

Dependency Parsing is used to identify the dependency relations between the words.

(iii)

It is a Top-Down Approach

It is a bottom-up approach

(iii)

It provides the type of words like subject, verb etc.

It identifies the relationship such as subject-verb agreement.

(iv)

Best suited for NLU Natural language understanding

Best suited for NLG Natural language generation

(v)

Suitable for tasks like POS tagging, NER, Stemming

Suitable for Text to Speech, Machine Translation etc.

(vi)

It is more for semantic analysis, like Error Free Grammar.

It is used for less semantic tasks, has dependency based grammar.

(viii)

It provides a parse tree as output.

It provides a dependency tree as output.

(viii)

Good for language with higher semantics, like Japanese, Thai.

Not able to pick up semantics as suitable for English, French.

through corpus
ability to O.

and, the

such as

operation



Q ⑦

Text-Preprocessing is the process of getting a text as input, perform various operations to break down and understand the text.

TEXT PRE-PROCESSING

DOCUMENT TRIAGE

TEXT SUMMARISATION

DOCUMENT TRIAGE:

Document triage is a part of NLP. Natural language understanding.

Here a document is received by the system. It is analysed and converted into a set of digitalised words, which are machine readable.

TEXT SUMMARISATION:

It is a part of NLP. Natural language generation.

This involves various steps.

* Sentence Tokenisation: Breaking the paragraph into various sentences.

* Word Tokenisation: Breaking the sentence into words.

* Stemming: Converting words to root form which may (or) may not have meaning.

Eg: Happiness Intelligence Better
Stemmed: Happy Intelligent Best

* Lemmatization: Converts words into Meaningful words

Eg: Happiness Intelligence Better
Lemma: Happy Intelligent Good

* Identify Stop words: Remove mostly repeated fillers like "a", "an", "the" etc.



of getting a list as
break down and

Summarisation

2. Natural language

system. It
set of digitalised

information.

graph with various

into words.

which may (or)

Better

Better

as

fillers

words that can be removed and still form a meaningful sentence.

* **PostTagging** : Identify Parts of Speech of each word

* **NER** : Named Entity Recognition

This tags each word into Name, Organisation,
Date, Place etc.

Tom went to HSB Bank Head Office in
New Orleans.

Organisation : HSB Bank

Name : Tom

Location : New Orleans

IMPORTANCE :

It is very important to pre-process the text.
Otherwise any model run on the preprocessed text will
yield noisy results.

Or, the model will not be able to predict.

Q ③

NLP has two main components NLU and NLG

i)

Natural language Understanding :

NLU is a way to understand text & learn from various
metadata that is input to the system.

It converts the human language into a machine
understandable form.

NLU is very hard to process than NLG.

This is because, NLU is dependent on context of words, its quality.
It will also take a lot of time to train with the
corpus to get accurate results.



Natural Language Generation:

NLG converts the machine format back into human understandable form.

When the machine processes the input and about to deliver the results, the NLG helps to convert it.

NLG involves tokenization, Stemming, NER, stopwords analysis, etc.

Once the analysis is done, data is finally chunked into various words for the user to understand.

Q 9

LANGUAGE MODELLING:

Language modelling is a way of getting tokenised words / lemmas as input and to check if they are syntactically correct to be accepted as the corpus.

i) STATISTICAL MODELLING:

Statistical Modelling uses probabilistic models to build the corpus, predict the next word.

If a $P = w$, the probability of occurrence of a word,

$$P(P = w) = w_1, w_2, \dots, w_{n-1}$$

Markov and Hidden Markov Models are under this category.



Neural Net



The next is simplified identified model.

HIDDEN MARKOV:

It consists of

NEURAL NETWORK

It is

word

It also

word enters

into a

hidden

happy

sad

happy

happy

happy

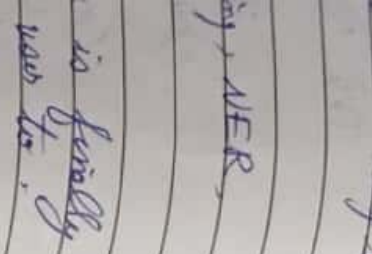
happy

happy

happy

THE UNIVERSITY OF CHICAGO

front back
right and
left help



HIDDEN MARKOV :

NEURAL NET LANGUAGE MODELING:

CNN - Convolution Neural Networks
RNN - Recurrent Neural Networks

word embedding is used to convert a word / ~~language~~ / ~~phrase~~ into a lower order function to easily interpret.

Word Embedding

Neutral	"	"	[1, 1, 1, 0]
Hydrated	"	"	[0, 0, 0, 1]

Engrated " = (iv) [090, 11



Q.10

Challenges of Text Pre-Processing :

languages such as Chinese, Arabic are very different from languages like English.

English has a "space" character in between the words, this information can be generalised to provide a input for word tokenisation.

languages like Chinese, Arabic are complex and do not have such space between them.

Hence it is difficult to analyse the language

Syntactically. Semantically. This makes tokenisation difficult.

RECOMMENDATIONS

(i) Hence more corpus, more understanding of their language is required.

(ii) This also takes a longer time for getting robust model.

(iii) Competency Parsing has the ability to understand the intricacies and complexities as it tries to understand the words by itself using more Fine Grammars.

(iv) Advant of Internet are provides more data which can be interpreted.

(v) Global API's for each language can be used for training if available.

(vi) Go for Neural Network modelling as it does not require Feature Engineering and it can identify on its own.



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 - 1 by UGC

Verified all entries & found correct!
Jr. Supervisor's Signature & Page

Supplement No. SIU-19 A

04430

SC MHRP

Centre: _____

Seat No. (in figures)

2	0	4	0	0	6
---	---	---	---	---	---

Seat No. (in words)

Two Zero Four Zero Zero Six

MBA-B4

Semester: _____

III

Programme: _____

Course: _____

Natural Language Processing

Supplementary Answer Script No. I / II / III

BEGIN WRITING HERE :

Q 6

Naive Bayes Classifier of Text Classification is based on Aberrant Bayes Theorem.

Assumptions :

(i) No Feature Interaction : Does not provide any importance to any particular feature.

(ii) Features of Equal Importance : All the features are treated equally.

(iii) Numerical Variables : Numerical variables are normally distributed.

(iv) Categorical Variables : Categorical variables all follow Multi-nomial Classification.

(v) No Missing Data : There should not be any missing data.

Eg: A cat is white in color

If $P = \text{cat}$ The probability of the next word being

$(P = \text{cat} / \text{is}) = \text{This can be found}$

The word with the highest probability that comes next can be found using this algorithm.