

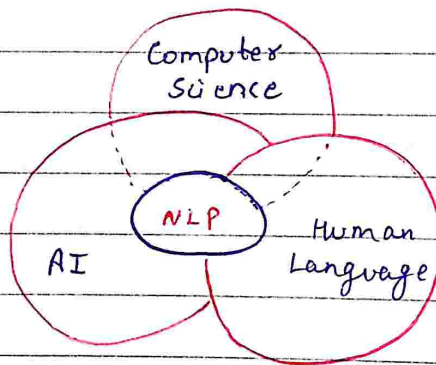
4/4/2024

NLP

IIMT

INTRODUCTION

Natural Language processing is a subfield of linguistics, computer science, and artificial intelligence concerned with the interactions between computers and human language, in particular how to program computers to process and analyze large amounts of natural language data.



Real World Applications

- Contextual Advertisements
- Emails clients - spam filtering, smart reply
- Social Media - removing adult content, opinion mining
- Search Engines
- Chatbots

Common NLP Tasks

- Text/ Document Classification
- Sentiment Analysis
- Information Retrieval
- Parts of speech Tagging
- Language Detection & Machine Translation
- Conversational Agents
- Knowledge Graph and QA Systems
- Text Summarization
- Topic Modelling
- Text Generation
- Spell checking & Grammar Correction
- Text Parsing
- Speech to Text

Approaches to NLP

- Heuristic Approaches
- ML Based Methods
- DL Based Methods

★ Heuristic Approaches

- जुगाड लगाना
- Ex → Regular Expressions
- Wordnet

Challenges in NLP

- Ambiguity
- Contentual Words
- Colloquialisms & slang
- Synonyms
- Irony, Sarcasm & tonal diff
- Spelling Errors
- Creativity
- Diversity

NLP Pipeline

IIMT

8/4/24

↳ WHAT IS NLP Pipeline

- NLP is a set of steps followed to build an end to end NLP software
- NLP software consists of the following steps:

- | | |
|---|---|
| <ul style="list-style-type: none">• Data Acquisition• Text Preparation<ul style="list-style-type: none">◦ Text cleanup◦ Basic Preprocessing◦ Advance Preprocessing• Feature Engineering | <ul style="list-style-type: none">• Modelling<ul style="list-style-type: none">◦ Model Building◦ Evaluation• Deployment<ul style="list-style-type: none">◦ Deployment◦ Monitoring◦ Model update |
|---|---|

↳ Points to Remember

- It is not universal
- Deep learning pipelines are slightly different
- Pipeline is non-linear

Data Acquisition

Data Augmentation

- + Synonym
- + Bigram Flip
- + Back translate
- + Add Noise

TEXT Preparation

→ Clean Tag/HTML Tag

→ Clean Emoji

→ Spelling Checker

Basic Preprocessing

Basic

↓
Tokenization

↓
Syntax Word
Sentence

Optional

+ Stop word Removal

+ Stemming

+ Removing digits, gender

+ Lower casing

+ Lang. detection

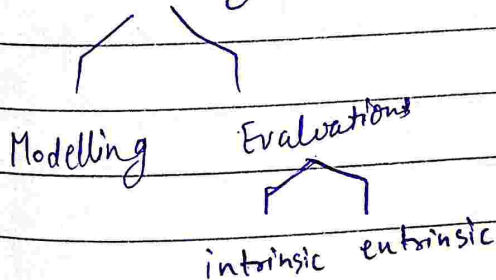
Advance Preprocessing

- + POS tagging
- + Parsing
- + Coreference resolution

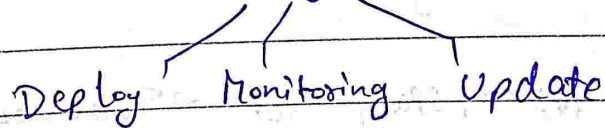
Feature Engineering

- Feature selection
- OHE, Bow, Word2Vec 'a

Modelling



Deployment



TEXT

PREPROCESSING

regen101.com

IIMT

6/4/24

Lowercasing

Remove HTML Tag

Remove URLs

Remove Punctuation

Chat word Treatment

Spelling correction

Removing stop words

Handling emoji's

Tokenization

Stemming

Lemmatization

most common techniques

Punctuation: '!"#\$%&'()*+,-./:;<=>?@[\\]^_`{|}~'

Spelling correction: 'tentblob'

Removing stop words: If you will apply POS then don't apply stop words
'nltk.corpus' + some code

Handling emoji → emoji.demojize

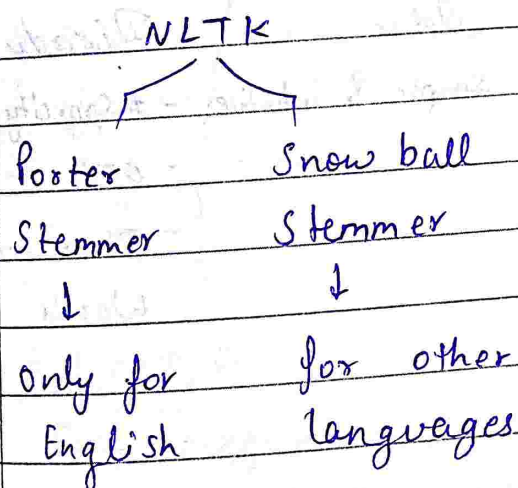
NLTK, spaCy for → Tokenization

Stemming

In grammar, inflection is the modification of a word to express different grammatical categories such as tense, case, voice, aspect, person, number, gender and mood.

en walk, walking, walked, walks

"Stemming is the process of reducing inflection in words to their root forms such as mapping a group of words to the same stem even if the stem itself is not a valid word in the language."



Problem in Stemming

⇒ Non english / language word or apply story → stori

Solⁿ ⇒ Lemitization

6/4/24

Text Representation

IIMT

Introduction

Similar terms: Text vectorization, representation, feature extraction

OHE {One Hot Encoding}

Disadvantages

- sparsity
- No fixed size
- OOV {out of Vocabulary}
- No capturing of semantic

Common terms

- Corpus: Concatenation of all words in dataset (C)
- Vocabulary: Unique words in corpus (V)
- Document: string of each datapoint (D) of each feature
- Word: group of characters (W)

BoW {Bag of Words}

	Data	Label
D ₁	people watch campus	1
D ₂	campus watch campus	1
D ₃	people write comment	0
D ₄	campus write comment	0

Adv

- Simple & intuitive

Disadv

- #Sparsity
- OOV
- Disordering of Words

V-dimensional vector

	People	watch	campus	write	comment
D ₁	1	1	1	0	0
D ₂	0	1	2	0	0
D ₃	1	0	0	1	1
D ₄	0	0	1	1	1

with Binary
without Binary

Why log in IDF? bcz if 1 Lakh total docs & only a word present only in one document then without log $IDF = 100000$ if TF की Dominate कर देगा {Mean TF की कीमत बढ़ा देता है}

N-grams

"Bag of N-grams"

NOTE → हमारे Prev Data पर हम quad-gram भी बना सकते

bi-, tri-, n-, >grams

As per previous data
here we are using bi-grams

	people watch	watch campus	campus watch	people write	write people comment	campus write
--	--------------	--------------	--------------	--------------	----------------------	--------------

D ₁	1	1	0	0	0	0
D ₂	0	1	1	0	0	0
D ₃	0	0	0	1	1	0
D ₄	0	0	0	0	1	1

Benefits

- Able to capture semantic of the sentence
- easy implement

Disadvantages

- Slows down the Algo
- ↑ se dim
- OOV

Why?
If not, if a word present in all documents then $\log_e(\frac{n}{n})$ will 0
final → 0

TF-IDF

Term Frequency

Inverse Dot Frequency

In sklearn → $TF \times (IDF + 1)$

$$0 < TF < 1$$

$$TF \times IDF$$

Advantage

used in Information Retrieval Systems

Search Engines

Disadvantages

Some are prev

$$TF(t, d) = \frac{\text{No. of occurrences of term } t \text{ in document } d}{\text{Total no. of terms in document } d}$$

$$IDF(t) = \log_e \frac{(\text{Total no. of documents in the corpus})}{(\text{No. of documents with term } t \text{ in them})}$$

Word2Vec

IIMT

8/6/24

- A technique of text vectorization which able to understand semantic meaning it actually a Deep Neural Network.

Types

CBOW

Continuous

Bag of Words

↳ for small data

Skip-gram

↳ for large data

CBOW

fake problem

↳ solve \rightarrow vector

dummy problem

↳ by product

Skip-gram

\rightarrow reverse of CBOW neural network

Watch campus for data science

POS TAGGING

spacy
IIMT

What?

In simple words, we can say that POS tagging is a task of labelling each word in a sentence with its appropriate parts of speech.

In traditional grammar, a part of speech or part-of-speech is a category of words that have similar grammatical properties.

Example \Rightarrow Why not tell someone?

Adverb Adverb verb noun Punctuation marks,
Sentence closer

Applications

1. Named Entity Recognition
2. Question Answering System
3. Word sentence disambiguation
4. Chatbots

common example → Will Will Google Campus

How POS Tagging works?

S- Nitish loves campus - E

S- can Nitish google campus - E

S- will Ankita google campus - E

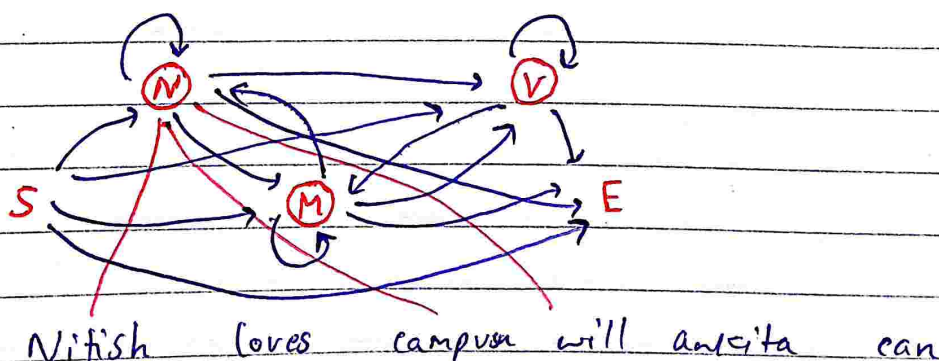
S- Ankita loves will - E

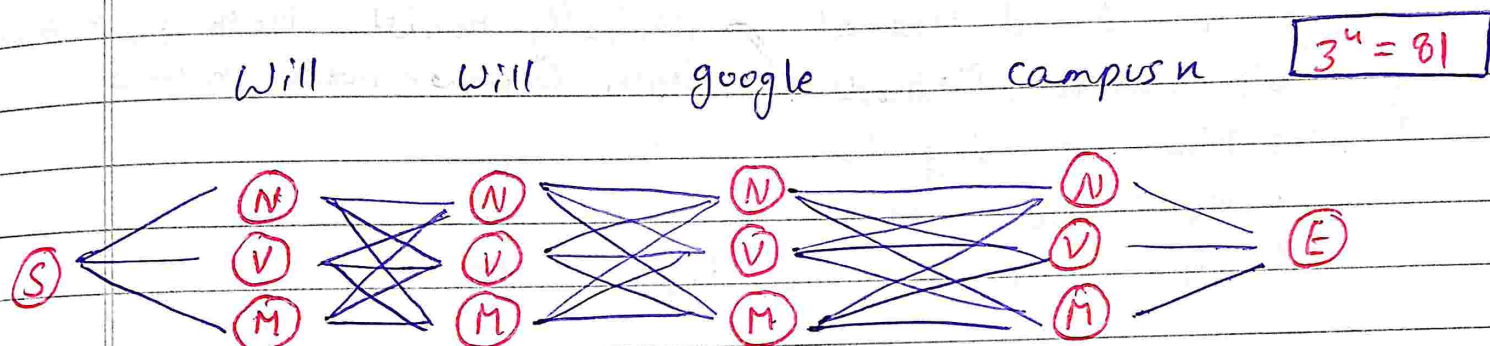
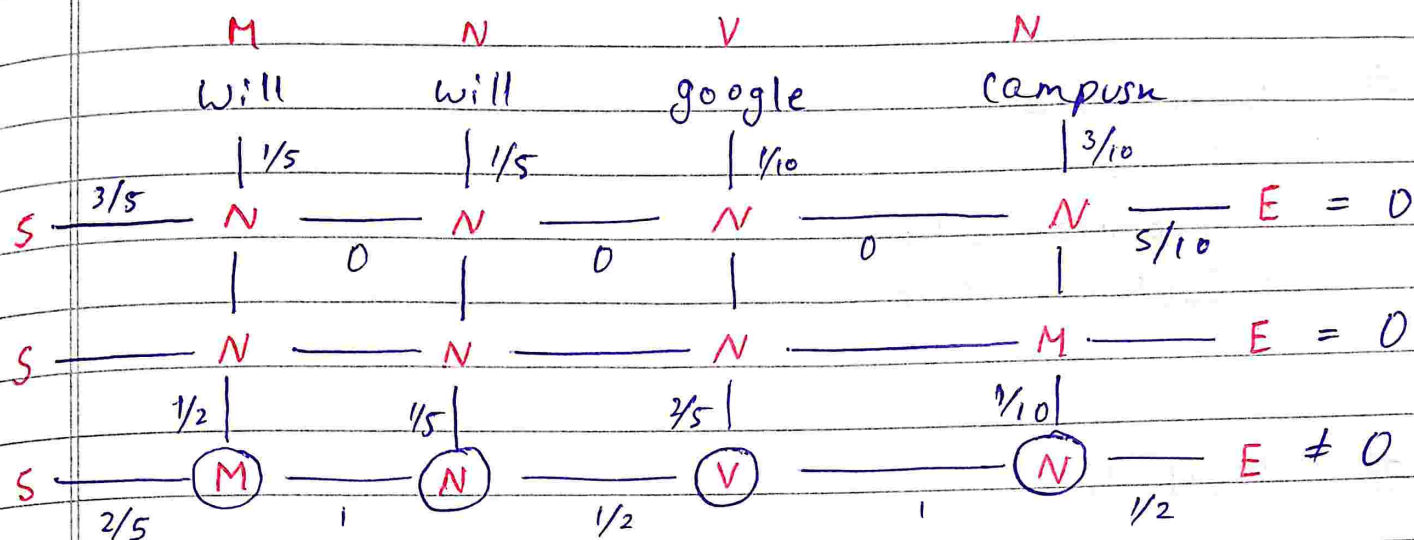
S- will loves google - E

Emission Prob			
	N	M	V
Nitish	2/10	0	0
loves	0	0	3/5
campus	3/10	0	0
Google	1/10	0	2/5
will	2/10	1/2	0
Ankita	2/10	0	0
can	0	1/2	0

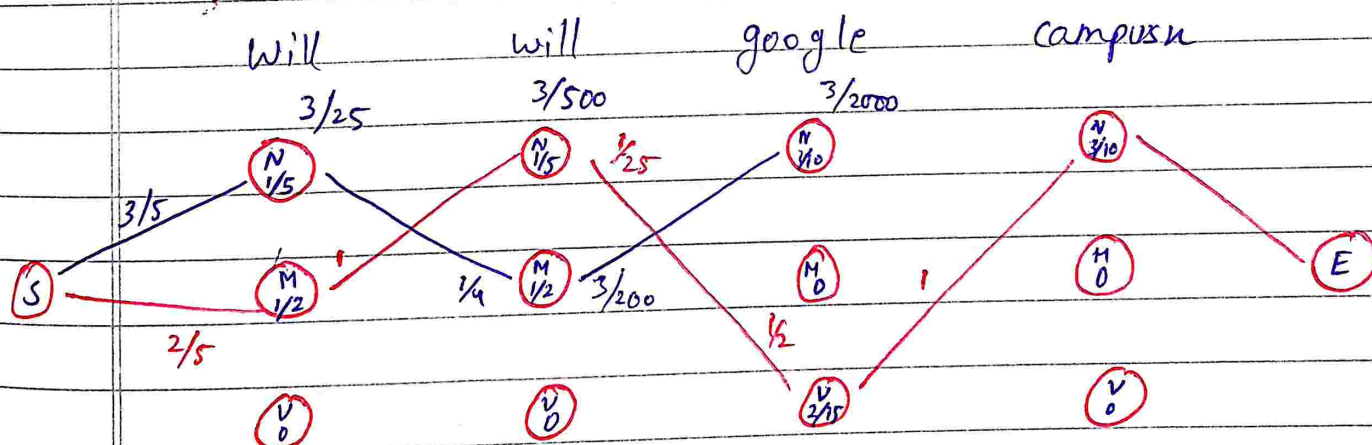
Transition Prob				
	N	M	V	E
S	3/5	2/5	0	0
N	0	0	5/10	5/10
M	2/2	0	0	0
V	5/5	0	0	0

Hidden Markov Model





Viterbi Algorithm



→ All probabilities calculated by multiplication at the end.

→ Next Topics: Topic Modelling, Named Entity Recognition

NER

Named Entity Recognition

IIMT

9/4/24

#

Applications

- Search
- Recommendations
- Customer Care

#

Methods

1)

Lexicon Based Methods

→ Basically Heuristic Methods, X-training

2)

Rule Based Methods

Regen, Grammar based patterns

3)

Machine Learning based Methods

◦ Multiclass classification

◦ Conditional Random Field (CRF)

4)

Dictionary/Lookup - Based Methods

→ Match words against a predefined dictionary or gazetteer

→ Good Precision, Low Recall (Misses unseen entities)