

Fundamental NLP tasks extract linguistic information such as syntax and semantics from words, sentences and documents. While they can provide useful information for user end application such as m/c translation and automatic dialogue systems.

On the other word level, morphology is a branch of linguistics that investigates the structure and formation of words, the task of morphological analysis study automatic prediction of morphological features of input words.

Task	Input	Output
Morphological Analysis	walking	walk+ing
Tokenisation	Mr. Smith	Mr. Smith
Pos tagging	I can open this can	PRP MD VB DT NN

Phases in NLP There are phases in NLP in order to extract meaningful information from text.

Once these phases are completed you are ready with your refined text and then you can apply some machine learning algorithm model to predict something.

i) Lexical Analysis:- In this phase text is broken down into paragraphs, sentence and words. It includes techniques as follows:

(i) Stop word remover (removing 'and' 'of' 'the' etc)

(ii) Tokenisation (breaking the text into sentence or word)
→ word tokenizer

→ Sentence tokenizer

→ Tweet tokenizer

(iii) Stemming (removing ing, es from tail of word)

(iv) Lemmatization (converting word to their base form)

2. **Syntactic Analysis**:- This is used to check grammar, arrangement of words and relationships between words.

Example: Truck is eating oranges

The word does not make sense.

Here we need to analyse the intent of words in Sentence. Some of techniques used in this phase

→ Dependency Parsing

→ Parts of speech (POS) tagging

3. **Semantic Analysis**:- Once tagging and word dependencies are analysed, semantic analysis extract only meaningful information from the text and reject/ignore the sentence that do not make sense.

Eg:- Truck is eating oranges will be ignored from information summary.

4. **Discourse Integration**:- Its scope is not only limited to word or sentence, rather discourse integration helps in studying the whole text.

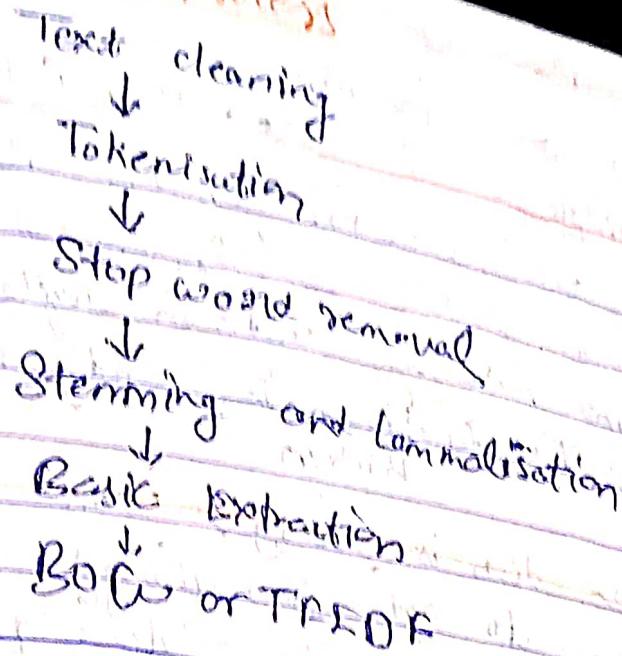
Eg:- John got ready at 9 AM. Later he took the train to California.

Here, the MC is able to understand that the word "he" in the second sentence is referring to "John".

5. **Pragmatic Analysis**:- It is a complex phase where MC should have knowledge not only about the provided text but also about the real world.

Eg:- Thank you for coming so late, we have wrapped up the meeting. (Contain sarcasm)

NLP Process



* NLP Based problem usually handle unstructured data and when the data is unstructured form, then data processing become difficult.

Unstructured :- which is not proper structured like - Video, text, image, Audio.

Application of NLP

- (i) Autocomplete feature in Anail
- (ii) Voice Recognition
- (iii) Text to speech
- (iv) Chatbot
- (v) Voice Assistance (Alexa, Siri)
- (vi) Sentiment Analysis
- (vii) Spam and Ham Email

Component of NLP

Two component
One for (i) NLU (Natural language Understanding)

(ii) Natural language generation

→ Syntax and Semantics are learned by nlp. It is the step where nlp learn actual meaning and context of sentence. But there are few problems mentioned below which occurs while understanding the text.

(a) Lexical Ambiguity

(b) Syntactic Ambiguity

(a) Lexical Ambiguity:- When a word has more than one meaning.

Example:- I saw Bats (the Mammal Bat or wooden cricket Bat)

(b) Syntactic Ambiguity:- Presence of two or more possible meaning with in single sentence

e.g. The ~~child~~ is ready to eat chicken

(Chicken dish is ready to eat or chicken himself is ready to eat something)

(c) Referential Ambiguity:- Here, when it is not clear about the object you are referring to.

Example:- John called Jay. Later, he laughed.
(Here who is referring to John or Jay)

Natural language generation work in 3 phase

(D) Text Planning:- Here useful content is selected from Text.

(I) Sentence Planning:- Selection of word forming meaningful phrase

(II) Text realisation :- Execution of sentence plan

Phases of NLP - We have already discussed in previous page

Text cleaning

When the text corpus is given to us, it may have following issues:-

(1) HTML Tags

(2) Upper/Lower case inconsistency

(3) Punctuations

(4) Stop words

(5) Words in their root form

(6) Before using the data for prediction, we need to clean it.

(1) HTML Tags removal :- While scrapping data from a website, you may get HTML tag.

e.g. ` am neha from Allahabad `

(2) Upper and lower case inconsistency :- Remove inconsistency and convert everything into lower case

`text_data = text_data.lower()`

(3) Remove Punctuations :- Punctuations in the text do not make much sense, hence we can remove e.g. %, %

Tokenisation :- Process of splitting Text, phrases, sentences into smaller unit is called Tokenisation.

- e.g. (1) Splitting of text into sentences (Sentence is considered as Token)
(2) Splitting of ~~sentence~~ into words. (Word is considered as Token)

We can import different type of Tokenizer from nltk library

(1) Sentence Tokenizer :- Text data will be splitted into sentence from `nltk.tokenize import sent_tokenize`

(2) Word Tokenizer :- Text data will be splitted into word

(3) whitespace Tokenizer :- based on white space word are splitted

e.g. I love singing, dancing and cooking

After Tokenization → "I", "Love", "Singing", "dancing", "and", "Cooking"

~~Stop word removal~~ - There are words in our sentence which do not provide relevant information from the text.

Example:- and, of, it, the etc.

There are multiple library which deals with Stop word removal → NLTK.
→ Spacy
→ Gensim

Stopwords like → I, me, myself, am, is, an, our, over, you, you're, you've, you'll, you'd etc.

~~Stemming & Lemmatization~~ - When we work on textual document, removal of punctuation, stop word are not just enough. The word in sentence can be used as present tense, past tense or future tense. e.g. Word 'go' is used as 'go'/'goes' in present form but 'went' in past tense.

The inconsistency of data can affect the model training and predictions, hence we need to make sure that the words exist in their root form.

To handle this, there are two methods :-

(1) Stemming :- Stemming is the process of converting the inflected words to their root form. In this method suffixes are removed from inflected word so that it becomes root. e.g. 'Going' → Inflectional suffix will get removed and inflected 'going' word will become 'go'.

Developed → Develop

Developing → Develop

Development → Develops

Python library nltk is used for stemming

from nltkstem import PorterStemmer

(Stemming to be continued)

However there are some word which do not get properly handled by stemming process.
e.g. - "went", "flew", "saw"
Ans (past tense ("went"))

Op^t, But

Stemming process is not smart enough it just know how to trim suffix part.

Fortunately we have Lemmatization

Pros Good part of stemming is not only useful for English language , it also useful for other language

Pros - Computationally fast, bcz it just trim the suffix

Cons - It is not useful enough if you are concern about the valid part. Stemmer can give you some word which do not have any meaning
e.g. "goes" → "goe"

* Lemmatization ! - It is where the words are converted to their root form by understanding the context of word in sentence.

Pros - The root word which we get after conversion holds some meaning and word belongs to Dictionary.

Cons - It is computationally expensive.

Code :-

```
from nltk.stem import WordNetLemmatizer
word_lemmatizer = WordNetLemmatizer()
```

print (wordnet.lemmatizer.lemmatize("going"))
print (wordnet.lemmatizer.lemmatize("goes"))
print (wordnet.lemmatizer.lemmatize("went"))

Output:- going
go
went

But you might be wondering that lemmatizer is unable to normalize the words "going" and "went" into their root forms.

This is because we have not passed the context to it. Part of speech is a parameter which we need to specify. By default it is noun.

Code:- print (wordnet.lemmatize("going", pos="Vb"))

print (wordnet.lemmatize("goes", pos="Vb"))

print (wordnet.lemmatize("went", pos="Vb"))

O/P:- go
go
go

So far we have looked into many token cleaning and normalization technique. Now it's time to study feature engineering.

A Basic Extraction:-

↳ How many words are present

↳ How many character " ",

↳ How many lowerCase " "

↳ How many Uppercase " "

↳ How many stop words " "

↳ What is avg length of each word

★ Bag of words : BOW is a technique to extract features from the provided text data. This technique counts the occurrence of words in sentence. If the word is found in sentence, then occurrence value increase by 1, else 0.

- (*) This technique is simple and easy to implement but it has own limitation too; we will discuss later.
- (**) Once all steps are done which we have discussed earlier like cleaning, stopword, lemmatization etc., then apply BOW.

Cols Text Cleaned Text

I am too lazy doing anything	dazi anything	
Today is Friday	today friday	
Friday is good	friday good	
Friday is near weekend	friday near weekend	
I am too lazy		

	today	is	Friday	near	weekend	good
Sent1	1	1	1	0	0	0
Sent2	0	0	1	0	0	1
Sent3	0	1	1	1	1	0

O/P: 'Today': 1, 'is': 3, 'Friday': 3, 'near': 1, 'weekend': 1, 'good': 1

(*) now make new column of named as word count

→ This result for uncleaned text

when we apply BOW in cleaned text

	today	Friday	good	near	weekend
Sent1	1	1	0	0	0
Sent2	0	1	1	0	0
Sent3	0	1	0	1	1

Limitation of Bow!

① count vectorizer does not understand the meaning of word

e.g. ① I go to sleep at 10 PM and ~~get~~ go to walk at 7 AM.

② I go to walk at 10 PM and go to sleep at 7 AM.

After feature engineering both will ~~get~~ result in same column values.

2) There exist so many zeros in matrix and is called sparse matrix

TF-IDF (Term Frequency - Inverse Document Frequency)

But only focuses on frequency of words in sentence
Consider the scenario where BOP is not good enough :-

(1) Suppose we don't want to remove stop words in this case frequency of "is", "the", "a" will also be very high but actually these words don't make sense in the sentence.

(2) Suppose we are processing Product reviews of Amazon / Flipkart then the terms like "Product", "Item" are domain dependent and are used too often in each review. Hence these word model will not help model in learning anything.

(3) The mobile keyword in the phone data set is not giving any value. Keyword like "5G", "Splash Proof", "Android" will make more sense.

Hence TFIDF where most often words are suppressed (given lower importance) and unique words (less frequent one) provided a higher weightage in the sentence.

I do not like Vanilla cake

I do not like Vanilla Ice cream

In above both sentences unique words are "Vanilla" and "Ice cream".

Example :-

Suppose no. of words in document A : 6

No. 1 1 1 1 4 8 3 6

Term	Document A - Term Frequency	Document B: Term Boolean	TDF
I	1/6	1/6	$\log(2) = 0$
do	1/6	1/6	$\log(2) = 0$
not	1/6	1/6	$\log(2) = 0$
like	1/6	1/6	$\log(2) = 0$
Vanilla	1/6	1/6	$\log(2) = 0$
Cake	1/6	0	$\log(2/1) = 0.69$
Icecream	0	1/6	$\log(2/1) = 0.69$

TF-IDF = Term * DF	I	do	not	like	Vanilla	Icecream	GRe
DOC A	$(1/6) \times 0$	$(1/6) \times 0$	$1/6 \times 0$	$1/6 \times 0$	$1/6 \times 0$	$1/6 \times 0.69$	0×0.69
DOC B	$1/6 \times 0$	$1/6 \times 0$	$1/6 \times 0$	$1/6 \times 0$	$1/6 \times 0$	$1/6 \times 0.69$	$1/6 \times 0.69$

TFIDF	I	do	not	like	Vanilla	Icecream	Cake
DOC A	0	0	0	0	0	0.118	
DOC B	0	0	0	0	0	0	0.115

Above example you can see, less frequent word have more weight

We can also observe this by sklearn

↳ do not likes vanilla cake

↳ do not likes vanilla Icecream

After using sklearn:

OP: [6, 'like': 1, 2, 6, 'Vanilla': 2, 'cake': 0, 'Icecream': 1]

Cleaning

When we apply TFIDF

	Cake	Icecream	like	Vanilla
Cake	1.405	0	0.50	0.50
Icecream	0	1.405	0.50	0.50
like	0.50	0.50	1.405	0
Vanilla	0.50	0.50	0	1.405

N-gram: - N-gram is continuous sequence of words or symbol or token in document.

① N-gram is continuous sequence of n items from given sample or text.

Remember the days we learned how to input an array by first inputting its size ('n')

② n Same for n-gram, n-gram is classified as

n Term

1

Unigram

2

bigram

3

Trigram

n

n-gram

Why we need different types of n-gram?

Ans. Because different type of n-gram suitable for different types of Application. You should try different type n-grams on your data in order to confidently conclude which one works the best among all for your text analysis.

Sl. No.

Type of n-gram

Generate n-gram

Unigram

["I", "reside", "in", "Allahab

Bigram

["I reside", "reside in", "in Allahab

Trigram

["I reside in", "reside in Allahab

Here I haven't consider removal of stop words. Now it is clear that unigram means taking one word at a time, Bigram means two words at a time and so on.

- Steps
- Explore the dataset
 - Feature extraction → shape (→ colour, rows)
 - Train - Test Split
 - Basic Pre-processing → remove punctuation, stop words
 - Code to generate n-gram
 - Creating Unigram
 - Creating Bigram
 - Creating Trigram

Doubt: What is the difference b/w Tokenized word and Unigram

Numerical on n-gram

Q: Check the probability of $I \text{ am } not$ using bigram.

$\Delta \langle S \rangle I \text{ am a human} \langle I \rangle S$

$\langle S \rangle I \text{ am not a stone} \langle I \rangle S$

$\langle S \rangle I \text{ I live in Mumbai} \langle I \rangle S$

SOP Probability of I when S is given $\rightarrow P(I|S)$

$P_{II} \rightarrow I \text{ when } I \rightarrow P(I|I)$

$I \text{ am } \rightarrow I \text{ am } \rightarrow P(\text{am}/I)$

$I \text{ not } \rightarrow I \text{ not } \rightarrow P(\text{not}/I)$

$I \text{ not } \rightarrow P(I|not)$

So we calculate all of the above for this we count the all individually (→ $C(S)$, $C(I)$, $C(I|am)$, $C(I|not)$)

$$\Rightarrow \frac{C(\langle S \rangle | I)}{C(S)} \cdot \frac{C(I)}{C(S)} \cdot \frac{C(I|am)}{C(am)} \cdot \frac{C(am|not)}{C(not)}$$

$$C(\text{not} \langle I \rangle S) \\ C(\text{not})$$

→ Now EH Sentence 1, 2, 3 & check $\frac{P(I|S)}{P(\text{not}|S)}$ $\frac{P(I|S)}{P(\text{not}|S)}$

$\langle S \rangle$ and I ने एक sentence में ~~सारे~~ 3 वाक्य
एक divide करने total वाक्य $\langle S \rangle$ के 4
एक $\frac{C(\langle S \rangle/I)}{C(\langle S \rangle)}$ 3 वाक्य Same

$$\frac{C(I/I)}{C(I)} \text{ First} \rightarrow \frac{C(I/cm)}{C(I)}$$

$$\Rightarrow \frac{3}{3} \times \frac{1}{4} \rightarrow \text{I am 1st sentence of } \\ \frac{1}{4} \times \frac{2}{4} \rightarrow \text{I am 2nd sentence of } \\ \frac{1}{2} \times \frac{0}{1} \rightarrow \text{I am not 3rd sentence of } \\ \downarrow \text{Total 3rd I.g.} \\ \rightarrow 0$$

Q1 Consider the following data

$\langle S \rangle$ I am Jack $\langle /S \rangle$

$\langle S \rangle$ Jack I am $\langle /S \rangle$

$\langle S \rangle$ Jack I like $\langle /S \rangle$

$\langle S \rangle$ Jack I do like $\langle /S \rangle$

$\langle S \rangle$ do I like Jack $\langle /S \rangle$

Assume that we use bigram language model

based on above data what is most probable
next word predicted by model??

Q1 - $\langle S \rangle$ Jack —

$$(3/4 \text{ or } \log(3/4) + \text{const})$$

(2) = $\langle S \rangle$ Jack I do —

(1/1 + 1/2 + 1/4)
(1/1 + 1/2 + 1/4)

Sol

$$P(I/S) = C(\langle S \rangle/I)/C(\langle S \rangle) = 4/5$$

$$P(Jack/S) = C(\langle S \rangle/Jack)/C(\langle S \rangle) = 3/5$$

$$P(do/S) = C(\langle S \rangle/do)/C(\langle S \rangle) = 1/5$$

$$P(am/I) = C(I/am)/C(I) = 2/5$$

$$P(like/I) = C(I/like)/C(I) = 2/5$$

$$P(do/I) = C(I/do)/C(I) = 1/5$$

$$P(\langle s \rangle | Jack) = C(Jack / \langle s \rangle) / C(Jack) = 2/5$$

$$P(\langle s \rangle | like) = C(like / \langle s \rangle) / C(am) = 1/2$$

$$P(I | Jack) = C(Jack / I) / C(Jack) = 3/5$$

$$P(\text{like} | do) = C(do / like) / C(do) = 1/2$$

$$P(Jack | like) = C(Like / Jack) / C(Like) = 1/3$$

$$P(Jack | am) = C(am / Jack) / C(am) = 1/2$$

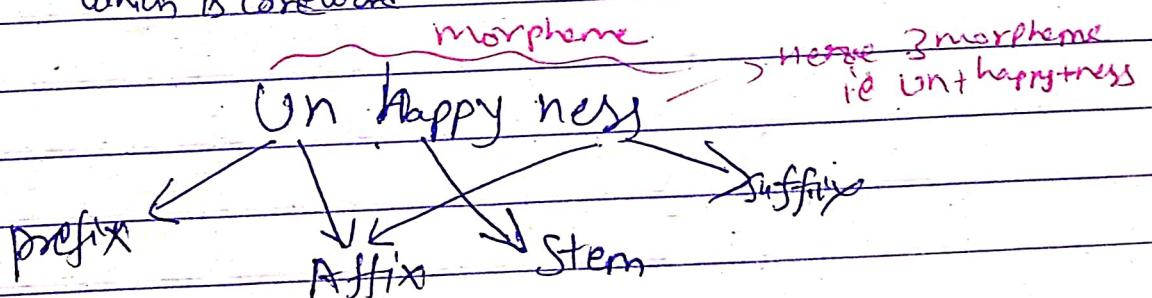
2) $\langle s \rangle$ Jack \xrightarrow{I}
 2:- $\langle s \rangle$ Jack $\xrightarrow{I, like}$

Spell-check

Finite-state Xducer :- A Xducer is piece of software that maps one stream of symbols on to another stream of symbols.

What is morphology \Rightarrow morphology is study of ways words are built from morphemes

* Morphemes :- morpheme made up of two things
 ① Stem ② Affix \rightarrow which either change meaning of word
 which is coreward \rightarrow or grammatical function



Process of determining morphemes :-

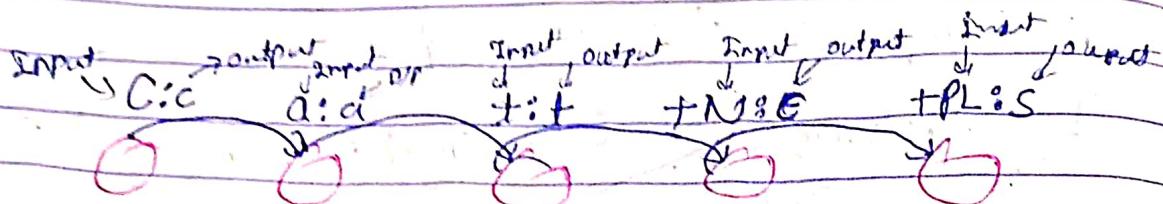
Cats := cat + N + Pl
 ↓ ↓ ↓
 word noun plural

Cat :- cat + N + S []

FST (Two Tapes)

Lexical set $\{C|alt|+N|+PL\}$

Surface $\{C|alt|S\}$



(Note: G shows empty)

Notes - ① FSA represent set of strings

Ex: {walk, walks, walked, long, long, long}

② An FST represents a set of pairs of strings
 $\{(walk, walk + V + PL)\}, (walk, walk + N + S), (walk,$
 $\text{walk} + V + PAST)\}$

③ FST have Input & output pair on tapes
 where a) PFA automata only has Output pair

* WordNet - A Lexical Taxonomy of English words

- When it comes to NLP, understanding meaning of words, as well as pre-processing textual data, can be a challenging task. To support this, we often use lexicons.
- We often map the text in our data to lexicon which in turn help us understand the relationship between those words.
- A really useful lexical resource is WordNet. It's a unique semantic net help us to find word relation, synonymy, grammar etc.

WordNet is large lexical database of English words.

Nouns, Verbs, adjectives and adverbs are grouped into set of cognitive synonymy called synsets.

- The NLTK module includes the english WordNet with 155,987 words and 117,659 synonym sets.
from nltk.corpus import wordnet

word:

Semantics:- Study of meaning

Semantic Analysis:- Text की meaning findout करने का process

होता है | Text की meaning find out करने का process

जो को meaning नहीं होता तो को समझाते हैं कि को meaning होता है

computer की समझता है।

Ambiguity:- word \in multiple meaning जैसे

that is known as ambiguity

3rd को Text किया जाता है जो की semantic analysis

3rd Text की को Ambiguity जैसे होता है। 3rd

जो 3rd होता है wordNet की में है, so basically wordNet is a dictionary which consist of several words.

And analysis with the help of wordnet semantic analysis can identify the actual meaning word in sentence.

Wordnet is offshoot of Semantic Analysis.

What is hypernym and hyponym

(i) Class etc. is called Hypernym and class is

(ii) element etc. is called Hyponym etc.

e.g., Animal → It is Hypernym

Dog → It is Hyponym

Ambiguity:-

e.g., She killed him with a baseball bat

wooden mammal object

bat के दो meaning and word sense disambiguation

प्रश्न के लिए यह समस्या है कि वर्ड के multiple meaning के लिए sentence के According के meaning को कौन कौन से meaning को दें।

- Word Sense Disambiguation is done via Approach
- Knowledge based Approach
 - Hybrid
 - Supervised and Unsupervised

What is Synonym & Antonym?

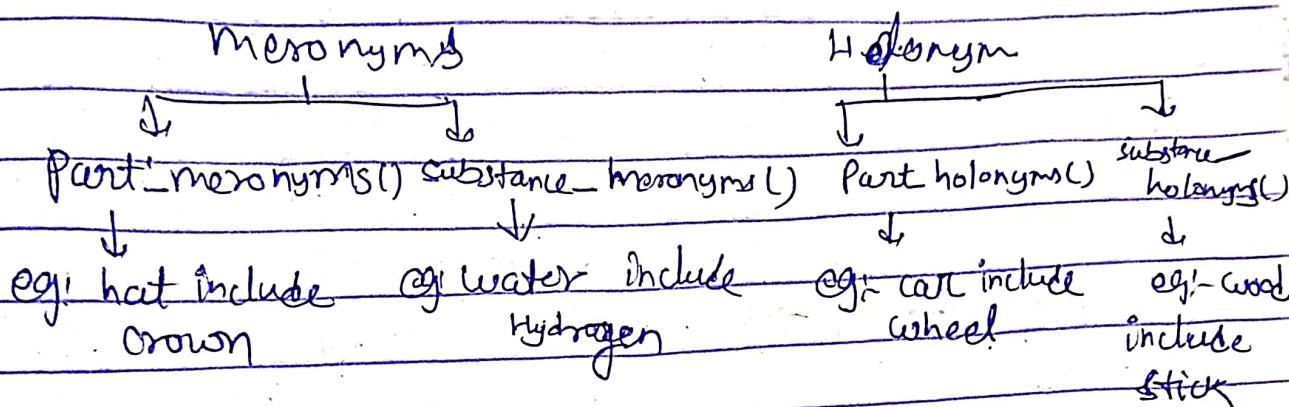
Word net covers this quite well. Synonyms mean which have similar meaning in some or all sense. For e.g. Car → motorcar, automobile etc.

Antonyms mean opposite meaning of word; e.g. Fat is opposite of thin

Meronyms and Holonyms

Meronymy expresses the 'component-of' relationship. That is, a relation between two concepts where Concept A makes up a part of Concept B.

Holonym expresses the membership of relationship relation b/w two concept, where Concept B is member of concept A.



Entailment - It is semantic relationship b/w two verbs. The relation is unidirectional.
e.g. Snoring entails sleeping but sleeping does not entail snoring

Multiilingual WordNet - It is even utilized by google translate as a part of translation process. These include not only best known ones such as French, German, etc. but also less known Kannada, Albanian etc.

WordNet can be used to perform language translation.

POS Tagging

(1) Why not tell someone?
 Adverb Adverb Verb noun ^{punctuation mark}

(2) We use short representation referred to as "Tags" to represent the categories.

Lexical Term	Tag	Example
Noun	NN	Paris, France, someone, Ruby
Verb	VB	work, train, learn, run, skip
Determiner	DT	the, a

प्रारंभिक नाम का टैग होता है NNS

(3) It is a process of converting a sentence to a form. (What is form?)

Any → list of words or, list of tuples (word, pos tag)

उदाहरण से suppose कि words हैं तो POS tag की लिखेंगे

(1) यह एक संख्या है जो उन शब्दों का एक ट्रॉप है जैसे ("Neha", "NN"), ("ate", "VB")]

(2) इसका अर्थ है कि शब्द का एक POS tag है जो उसकी मुख्य क्षेत्रीय वर्गीकरण में से एक है।

(1) Open Class :- Nouns, Verbs, Adjectives, Adverbs

(2) Closed Class :- Prepositions, determiners, Conjunction, pronouns

e.g. - 'the' determiner है जो एक बहुत ही

* POS word को नियम देता है।
वे Tagger decide करता है। उसके लिए अन्य वाक्यों
की जिसी 2 pos tag हो सकती है। तो - यह words को देता है।
eg! I cry everyday
Verb

My company's name is CRY
Noun

पर ये कैसे decide करता है कि कौन सा Verb है या Noun?
Answer is neighbour word को बोलता है CRY में
(CRY noun के जरूरी ही यह हो सकता है) and I cry में
(cry verb के जरूरी है)

* Application:- (1) Text to speech : भारतीय शापें
→ Type करता है computer द्वारा voice के
द्वारा पर करता है spelling same होता है
जोड़ता है part में कौन सी वाक्यांशील pronunciation होगा
and present ही कौन सी वाक्यांशील से

Eg:- road → Part
→ Present

(1) Search Bar

(2) Parser

* Problem:- Major problem is ambiguity

* Tags set for English:-

→ Open class :- new-new words include things
like

closed class :- fixed

Rule Based POS Tagging → In next

Input → Tag Dictionary or + Handwritten rule → output
Lexicon (word, tag)

संबंधी पद्धति इनपुट एवं वाक्य के बीच POS
string or sentence एवं वाक्य के बीच POS

3rd Step में Tag Dictionary की सहायता से
वाक्य के बीच इसकी वाक्य के बीच Sentence
में विभिन्न सभी वाक्यों को POS tag
देते हैं। But कुछ वाक्यों को दो या तीन
विभिन्न POS tags देते हैं (जिसका कहा जाता है
Ambiguity)
Ex:- I play cricket everyday.
I want ^{Verb} do perform a play.

4th वा Ambiguity की Problem 3rd के Step 3rd
में होती है Hand written rules. Hand written
rules में विभिन्न rules होते हैं तो वे Particular
Sentence के POS की सभी वाक्य Appropriate होते हैं तो
Printout उत्तर होता है

I play cricket everyday.

^{Verb} _{Noun}

बहुत बहुत ये sentence को को (I/p) \rightarrow तो ये ये tag dictionary
or Lexicon के पास आया विभिन्न वाक्य, cricket, everyday
एवं को POS tagging मिलते हैं लेकिन play को वाक्य
verb, noun ये ये ये Hand written rule के पास
मिलते हैं तो Printout हमें ये sentence के According
most appropriate वाक्य ही देते हैं ये verb को play Verb के
लिए use के अनुदेश

Application of POS tagging

- ① Named Entity Recognition
- ② Question Answering System
- ③ Chatbot
- ④ Word sense Disambiguation

POS Tagging

In the
of Part-of-speech tagging steps

Implementation - Hidden Markov model is very important model. This generate model based on list of step sequence.

For this we will study Markov property

- (1) memory less (future state depends on previous state)
- (2) .

Watch Capes X YouTube video for POS tagging + HMM

Suppose we have:-

Sentence 1	Can	Nitish	goggle	Camp	TDTW
Senten(2)	Will	Ankita	google	Run	TDTW
Sentence(3)	Ankita	loves	will		
Sentence(4)	Will	loves	Google		
Senten(5)	Nitish	loves		TDTW	

First we know the question is first train above after drawing of above 4 sentence predict the sentence \Rightarrow Will will Google TDTW

Soln	First we manually train above 4 sentence and					
Sentence(1)	an	Nitish	①	Google	①	TDTW
Sent(2)	will	Ankita	②	google	②	TDTW
Sent(3)	Ankita	loves	③	will	③	TDTW
Sent(4)	will	loves	④	Google	④	TDTW
Sent(5)	Nitish		⑤	loves	⑤	TDTW

Now we will plot Emission table through above analysis.

Total 20R
which contains
and Transitions
so 81

	Non	Adj	Verb
Noun	8/10	0	0
Adjective	0	0	2/5
Adverb	2/10	0	0
Verb	1/2	0	0
Adverb	2/10	0	2/5
Adjective	0	2/2	0
Verb	3/10	0	0

will model
will go to
use for
Total model 2
Total 81

Now we will plot Transition Table:-

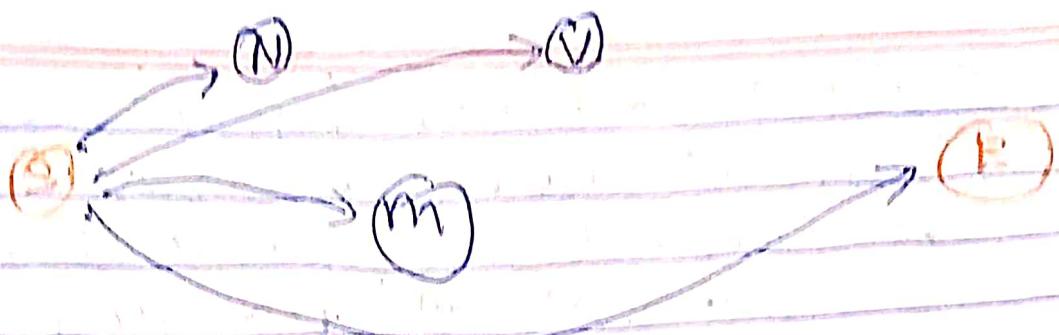
To draw this table we will assume
that every sentence starts with 'S' and
ends with 'E'

S	N	M	V	E
3/5	2/5	0	0	0
0	0	5	5	5
2	0	0	0	0
5	0	0	0	0

Now ET Row \rightarrow sum & divide 5 to get
Actual Transition Table

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

Now we will make Hidden markov model

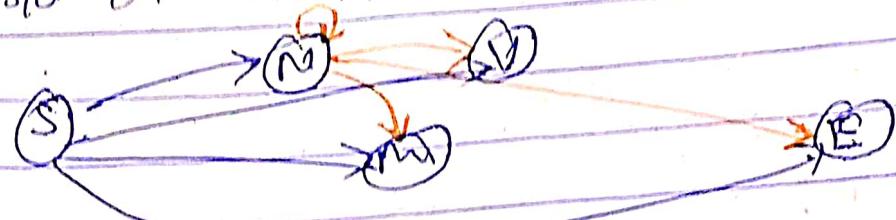


EMISSION part according to Emission Table in above \Rightarrow
Nouns, Verbs, Adverbs will generate can generate

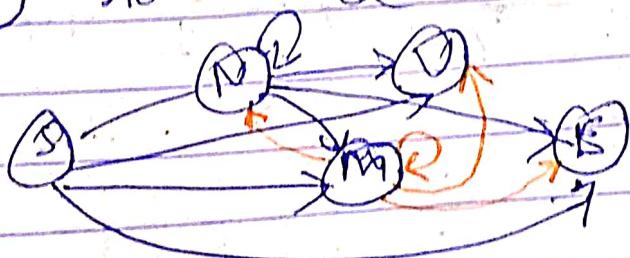
Step (1)

1) First draw our HMM start with initial state
2) Then (from translation to English)

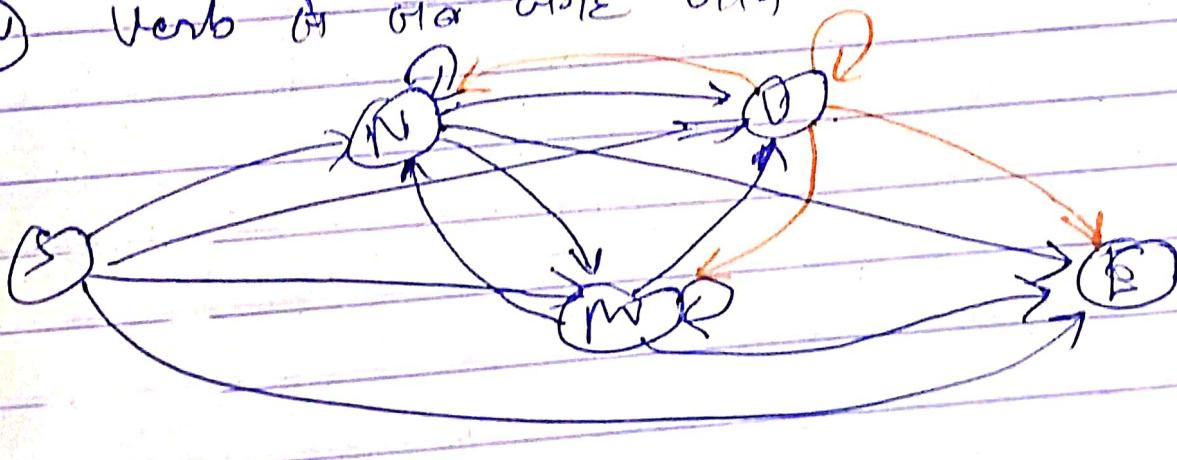
Step (2) 3rd EMT Noun to English words (verb)



Step (3) 3rd model \downarrow Hindi - English translation



Step (4) Verb to other language (Hindi)



Now after it Will will say google

Prob.

- "All item & combination try at it".
- 1st we will assume all words are noun now we will calculate first 3 nouns & off probability

