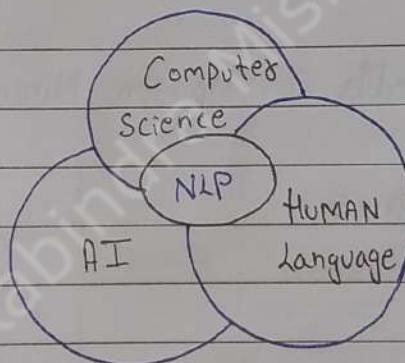


Introduction : (Wikipedia)

Natural language processing is a subfield of linguistics, computer science and AI 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.



NLP Diagram.

This field new and it has been just 50 years while 10 years :: gt has started booming
(current Year 2024).

- ① Natural Language : Language that has naturally evolved in humans.
It is different than constructed or Program language.

Job Role : NLP Engineers

PAGE NO.	11
DATE	

Main Objective : Making communication with machines using natural language for easier use.

(REAL WORLD Applications.)

1) Contextual Advertisements

{ displaying ad's based on your contextual info such as post, profile }

2) Spam filtering of Emails based on Contents

3) Social Media : Opinion Mining, Removing adult Content

4) Search Engines.

5) ChatBots.

COMMON NLP Tasks.

① Text / Document Classification

② Sentiment Analysis.

② Information Retrieval.

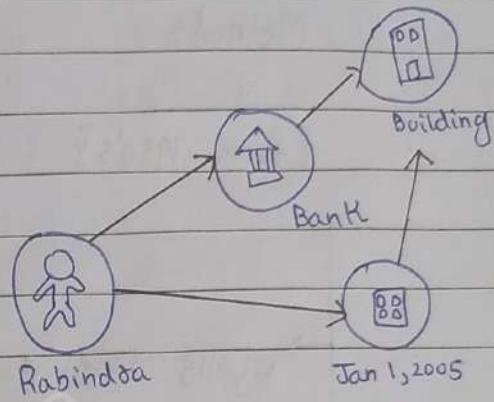
③ Parts of Speech Tagging

⑤ Language Detection
and Machine Translation

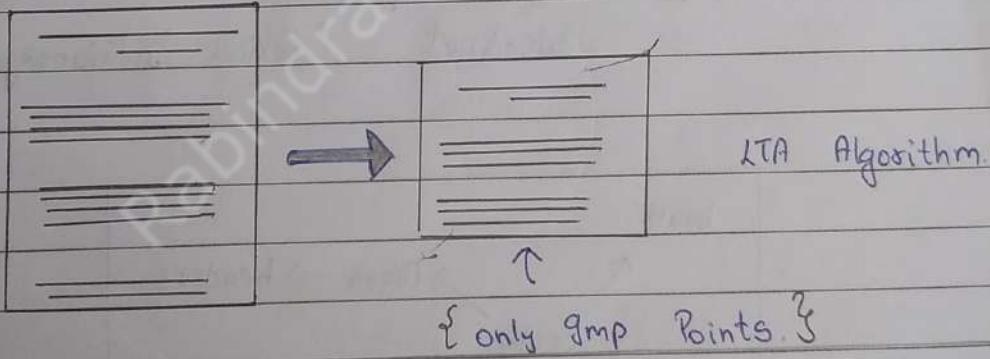
⑥ Conversational
Agents (Alexa).

⑦ Knowledge Graph and
QA Systems.

(Advance level)



⑧ Text Summarization



⑨ Text Generation { on Keyboards }

⑩ Spell checking and Grammar Correction

⑪ Text Parsing :

Sentence → Verb → "Tom"
→ Noun → Pronoun "ate"

⑫ Speech to Text :

APPROACHES To NLP

{ Started in 1950's }

① Heuristic Methods



{ from 1950's }

② Machine Learning Based Models.

{ from 2000's }

③ Deep Learning Based Models



{ from 2010 }

→ Using mental (Mental shortcut)

→ Count (+ve, -ve) words for sentiment Analysis.

→ Example : Regular Expressions

Wordnet (Lexical Dictionary).

book

Book → Pages → headers

Book

record → index → contents



Paragraphs.

Open Mind Common Sense.

Advantages :

① Quick

② Error is not so large

③ Used in current Scenario.

① Machine Learning Approach

(Big Advantage :) No need to make rules based on open approach and datasets
 (Most widely used).

Convert text to vector (ie numbers)

Algorithms Used

- ① Naive Bayes
- ② Logistic Regression
- ③ SVM
- ④ LDA
- ⑤ Hidden Markov Models.

② Deep Learning Approach

① Converting numbers to text or vice-versa may loose sequential information of texts (ie) sequence of words if we use ML model.

② Converts or generate Feature Selection.

i} RNN

ii} LSTM (Long Short Term Memory)

iii} GRU (Text Generation)

iv} CNN (Image)

v} Transformers (Google).

vi} Autoencoders (Set of Two neural networks)

(Encoder, Decoder).

NLP is Difficult

PAGE NO.	/ /
DATE	/ /

Challenges.

① Ambiguity

eg: I never tasted a cake quite like that one!

② Contextual Words.

eg: I ran to the store because we ran out of milk

③ Colloquialisms and Slang

eg: Cake Walk, pulling your leg.

④ Synonyms

eg: → What!! , What?

⑤ Spelling Errors.

⑥ Diversity of Languages.

NLP Pipeline

A set of steps followed to build an end to end NLP software.

Steps:

① Data Acquisition

② Text Preparation

- Text Cleanup
- Basic Preprocessing
- Advance Preprocessing

③ Feature Engineering (Bag of Word, TF-IDF)

④ Modelling (Model Building).

⑤ Deployment.

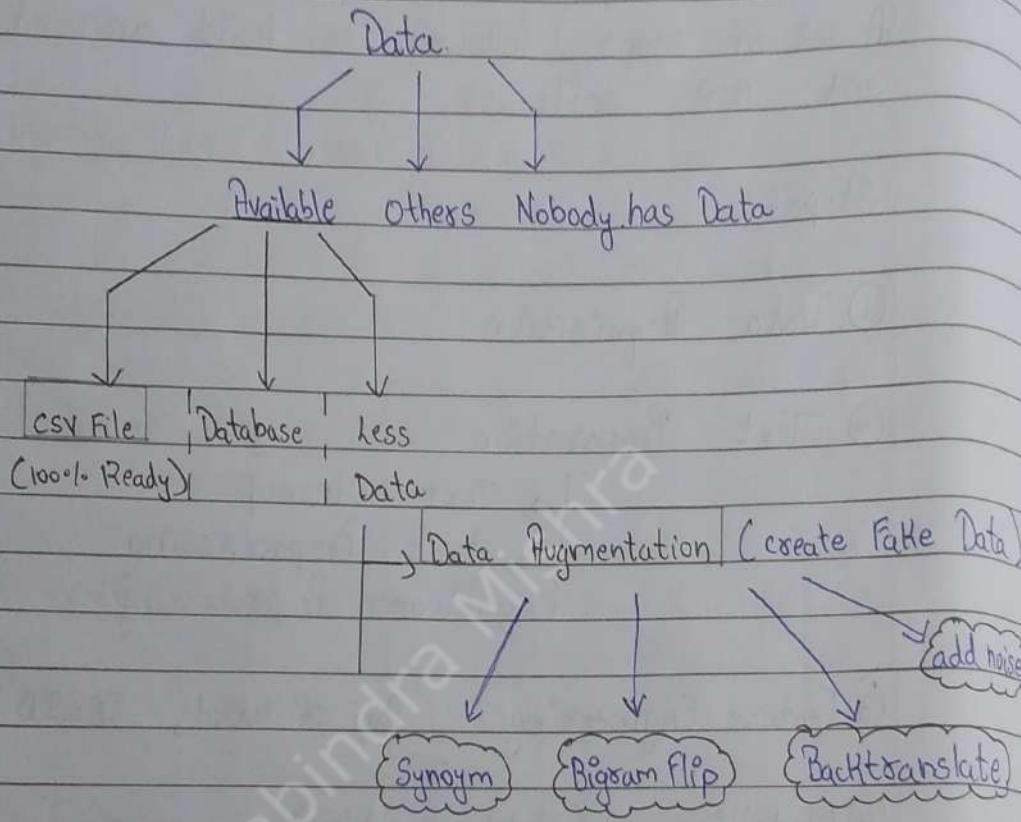
- Monitoring
- Model Update

Important :

① Sometimes Pipeline can also follow the Non-linear trends.

② Deep learning Pipelines are slightly Different.

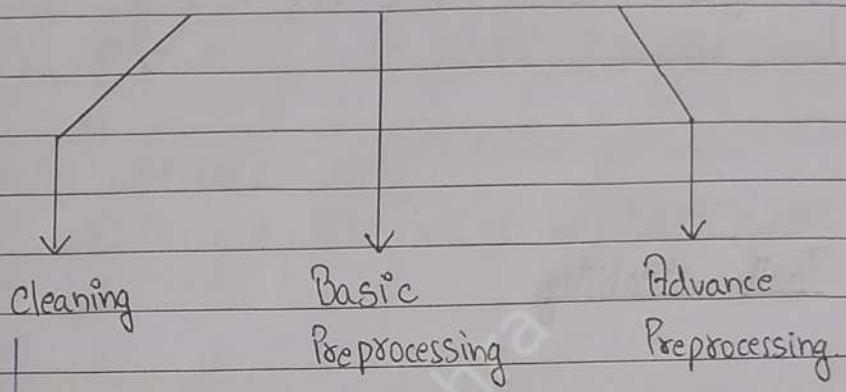
Data Acquisition.



Others :-

- Public Dataset
- Web scraping → (beautiful Soup)
- API
- Pdf
- Image
- Audio (Speech to Text).

Text Preparation :



→ html tag (in case Web scraping).

→ Unicode Normalization

→ Spelling Check.

Code:

```

import re
def clean(data):
    p = re.compile(r'<.*?>')
    return p.sub(' ', data)
    # It will remove all this
    
```

While web scraping we also extract emoji's, but need to convert them into machine understandable language.

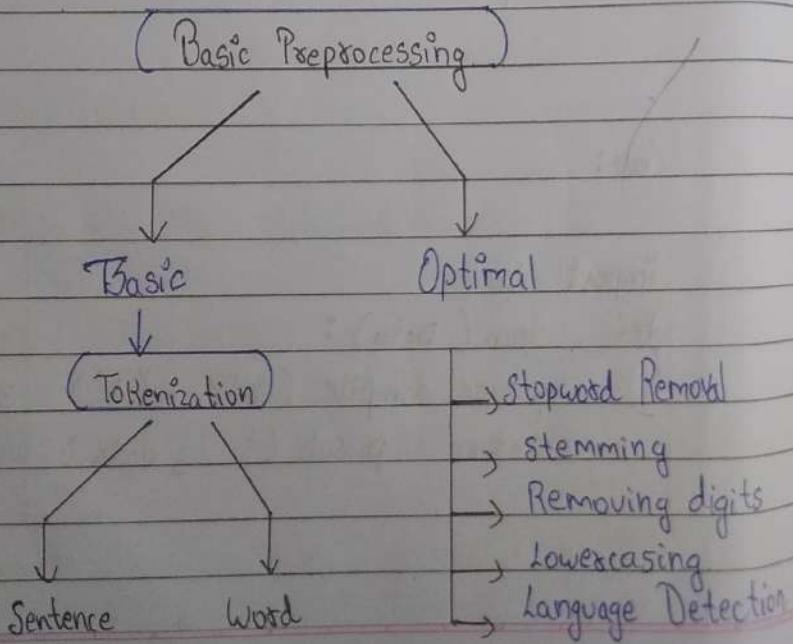
```
emoji_text.encode('utf-8')
```

III) Spell checking

incorrect_text = "manex ggeneration durringg"

```
from textblob import TextBlob
```

```
textB1b = TextBlob(incorrect_text)  
textB1b.correct()
```



■ Understand Tokenization using Code.

dummy = "Lorem Ipsum....."

```
from nltk.tokenize import sent_tokenize, word_tokenize
```

```
sents = sent_tokenize(dummy).
```

↳ It will break paragraphs into sentences.

Word Tokenization

```
for sent in sents:
```

```
    print(word_tokenize(sent))
```

Stemming :

Different words having same meaning

e.g. Dance, Dancing only difference is tense.

Thus in stemming we convert them into root word.

→ Advance Preprocessing

→ Pos tagging

→ Parsing

→ Coreference Resolution.

In Parsing we create a parse Tree for syntactic combination.

FEATURE ENGINEERING.

Basically converting texts to numbers because machines can process number in a more better way.

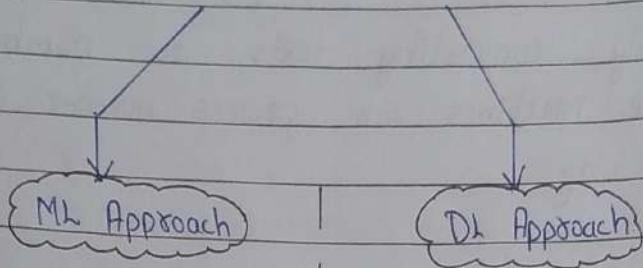
Example of feature engineering.

Review_Text	Sentiments.
	-1 (negative)
	1 (positive)

No. of +ve words	No. of -ve words	Neutral
5	3	2

Final outcome: Positive Review. (But these are basic)

Feature Engineering



Advantage: Interpretability

Disadvantage: Create Feature

on own.

Feature created by own need
not be helpful for model

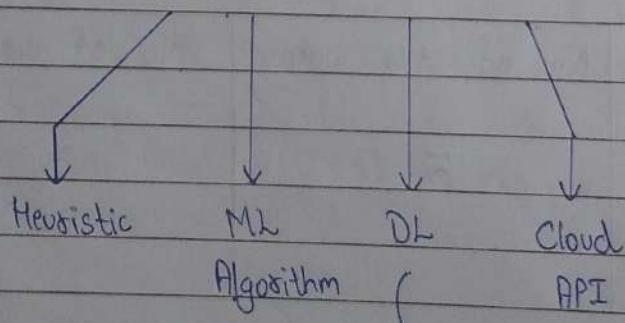
Advantage: No manual feature creation

Disadvantage: No interpretability ie

you never know why model is
performing good or bad.

MODELLING

Training the Model by applying algorithms.



Approach depends on

- ① Amount of Data
- ② Nature of Problem.

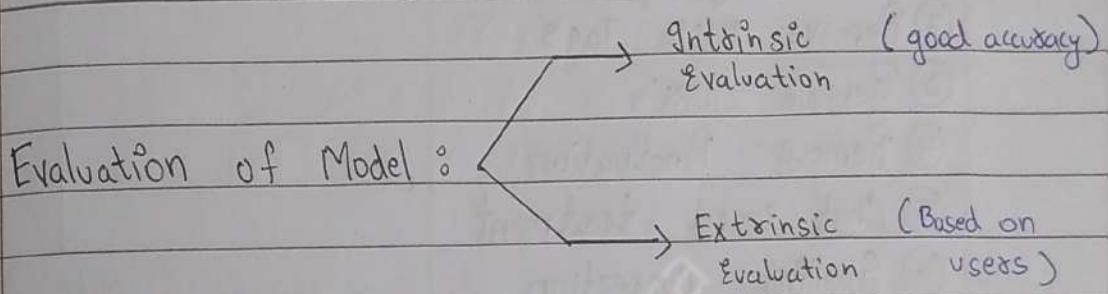
Transfer learning

BERT: A Transformer which is already trained on 40GB Data on internet.

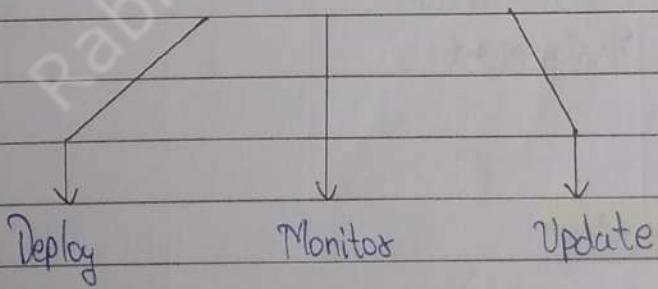
PAGE NO.	/ / /
DATE	

If enough data is not present use heuristic approach.

Deep Learning Requires very large amount of Data.



Deployment



API (ON cloud service)	Monitoring as Dashboard
chatBot	Using Graph
Android App.	KPI Cards.

update as per location / environment

TEXT PREPROCESSING

- ① Lowercasing
- ② Remove HTML Tags
- ③ Remove URL's
- ④ Remove Punctuation
- ⑤ Chat word treatment
- ⑥ Spelling Correction
- ⑦ Removing Stopwords
- ⑧ Handling Emoji's
- ⑨ Tokenization
- ⑩ Stemming
- ⑪ Lemmatization
- ⑫ Assignment.

→ Lowercasing:

∴ Python is a case sensitive language we convert our text either in lower or uppercase.
Because (Team, tram) in such cases Python will identify this as two different words.

① df['review'].str.lower()

Remove Punctuation

It is necessary because whenever we do tokenization our tokens must be appropriate.

'Hello! how are you?'

Hello(1), how, are, you, (2)

Hello!, how, are, you?

[Here this is being considered as separate word and without any meaning our tokens got increased.]

[But this can also create problem. Because Hello!, Hello will be considered two different words and 99% of use case we remove punctuation]

import string, time
string.punctuation

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

[All this include total punctuations.]

Code :

```
import string, time
string.punctuation
```

```
exclude = string.punctuation
def remove_punc(text):
    for char in exclude:
        text = text.replace(char, '')
    return text.
```

We have used Time module to calculate execution time especially for Jupyter NoteBook.

```
start = time.time()
remove_punc(text).
Time1 = time.time() - start
print(Time1).
```

But the above method is a slower approach.

Method 2 :

```
def remove_punc1(text):
    return text.translate(str.maketrans(' ', ' ', exclude))
S
function
```

df.sample(5)

o/p: # gt will give any 5 random sample from Datasets.

Chat word Treatment

eq: →

- ① ASAP
- ② GN
- ③ lmao
- ④ imho.

We need to replace this with Good Night.

#upper(): Returns all characters in uppercase.

Code:

Chat_words = { 'THX': 'Thank You', 'ASAP': 'As soon as Possible', 'WB': 'Welcome Back' } .

```

def Chat_conversion(text):
    new_text = []
    for w in text.split():
        if w.upper() in chat_words:
            new_text.append(chat_words[w.upper()])
        else:
            new_text.append(w)
    return " ".join(new_text).

```

Spelling Correction

Eg: "Please read the notebook and also like ntebook."

We can use `(textblob)` for this. | There are different other libraries.

Removing Stop Words

a the of are my } All these words are used for sentence formation

But do not contribute to the meaning much. So for different tasks such as "sentiment analysis" we remove these words.

We use `(NLTK.)`

↳ It contains list of stopwords.

Code:

```
from nltk.corpus import stopwords  
stopwords.words('english').
```

We need to specify language.

If we want Pos tagging we don't remove stopwords.

Code:)

```
def remove_stopwords (text):
```

```
    new_text = []
```

```
    for word in text.split():
```

```
        if word in stopwords.words('english'):  
            next_text.append ('')
```

```
    else:
```

```
        new_text.append (word)
```

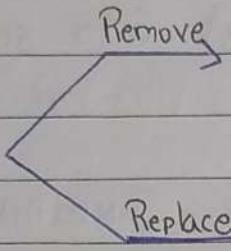
```
x = new_text [:]
```

```
new_text.clear()
```

```
return " ".join(x).
```

remove_stopwords ('probably my all time favourite movie, a story of selflessness and sacrifice')

HANDLING EMOJIS



Replace → Happy {we can replace it with words?}

U"\U0001F600-\U0001F64F"

↳ This is unicode where U denotes unicode string
range between U0001F600 - U0001F64F

{#} Printing all emojis in Python 3.x

```
emojis= [chr(code) for code in range(0x1F600, 0x1F610)]
print(" ".join(emojis))
```

Method 2) :- Replacing Emoji.

```
import emoji
print(emoji.demojize('Python is 🐍'))
```

O/P: Python is: Fire

Tokenization

Process of breaking the text into smaller parts called as tokens.

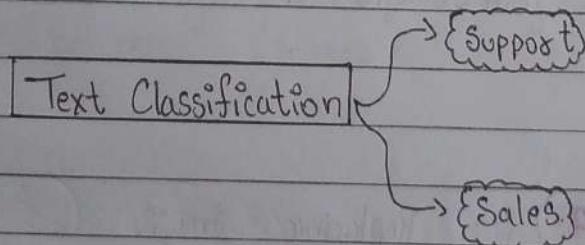
I am an INDIAN

[I , am , an , INDIAN .]

This is called as word tokenization.

Suppose we have Text and we need to do text classification.

and we need to classify the class based on text.



Problem: I am new in new delhi?

In this case it will create 2 tokens for "new" but the second new is for new Delhi.

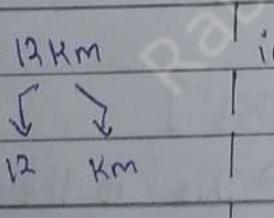
Problems during Tokenization:

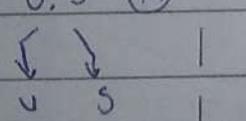
1) Prefix : character(s) at the beginning \$ ("i

2) Suffix : character(s) at the end KM), - !"

3) Exception : Special-case rule to
split a string into several
tokens. U.S , let's

4) Infix : character(s) in between. --- / ---

eg: ① 12 Km ie 12 & Km must be tokenized differently


② U.S (X) Here U.S must me tokenized as a


Using Split() Method:

①

Sent1 = 'I am living in Mumbai'

Sent1.split()

O/p: ['I', 'am', 'living', 'in', 'Mumbai']

Problem:

②

Sent3 = "I am going Delhi!"

Sent3.split()

O/p: ['I', 'am', 'going', 'Delhi!']

↳ This is the problem
because while analyzing we want [Delhi, :]

Method 2) Regular Expression

import re

Sent3 = 'I am going to Delhi!'

tokens = re.findall (" [w]+ ", sent3)

tokens.

Other Examples :

sent5 = " I have a Ph.D in A.I "

sent6 = " help! We're mail us nK@gmail"

sent7 = " A 5Km ride cost \$10.50"

Problem with nltk

- ① sent6 It breaks email into different tokens
- ② sent7 5Km does get tokenizes differently.

(Difference between math_operations.py)

def add(a,b): return a+b	# Collections of functions, class, variables.
-----------------------------	---

def subtract(a,b): return a-b	# Collections of functions, class, variables.
----------------------------------	---

Class calculator :

def multiply (Self,a,b): return a*b
--

Object= math_operations.calculator()

Method 3

SPACY

```

import spacy → nlp library
nlp = spacy.load('en_core_web_sm') → english dictionary
doc1 = nlp(sent5) } → convert sentence to documents
doc2 = nlp(sent6)
doc3 = nlp(sent7)
doc4 = nlp(sent1)
    
```

```

for token in doc4:
    print(token)
    
```

- for doc3 spacy works better. # it does not tokenizes 5RM differently
- # for doc2 also spacy works better.

For some cases spacy will work better while in other cases nlp.tokenize will work better.

No Technique is Fullproof.

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.

Stemming

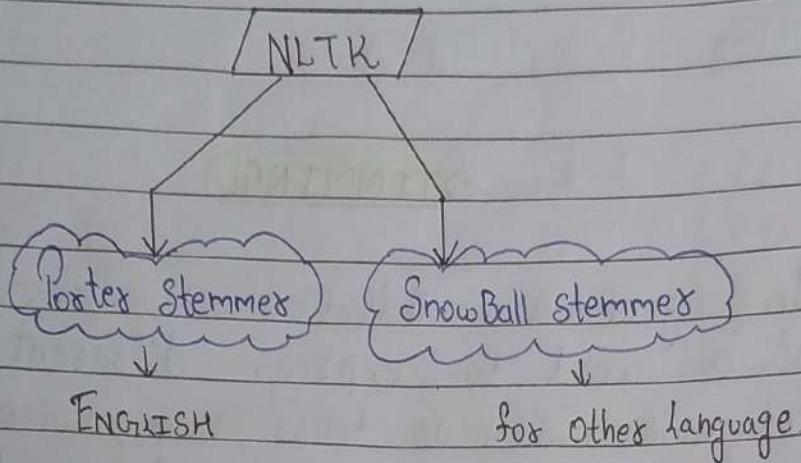
"Process of reducing inflection in words to their root forms such as mapping a group of words to the same stem even if the stem itself is not a valid word in the language"

These are mainly used for Information Retrieval System. eg: Google.

Suppose you are searching: fish, fishing, fishnet.
But whenever you are searching any one of them results must be displayed for all ie include all for wide explorations



We Use Stemmer (Algorithm) for stemming.



Output after stemming would sometimes not be even a valid english word. This is the biggest problem with stemming.

Example :-

```

from nltk.stem.porter import PorterStemmer
ps = PorterStemmer()
def stem_words(text):
    return " ".join(ps.stem(word) for word in
text.split())
text = "Dance like Dances"
  
```

stem_words(text)

O/P :- Dance like Dances

Thus, we need to apply lemmatization.

LEMMATIZATION

Difference in lemmatization is that the output of lemmatization will always be an English word.

But lemmatization is slow as compared to stemming.

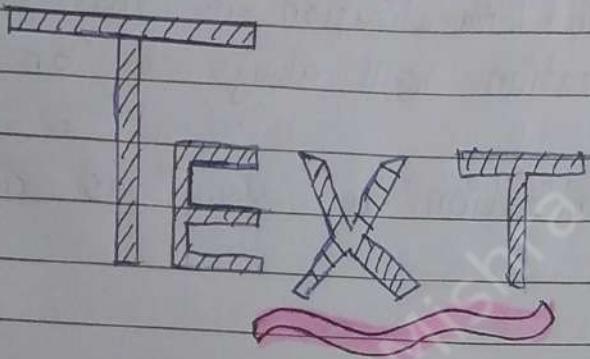
If speed is your concern use stemming.

In lemmatization root word is called as ~~st~~ lemma and lemma always belongs to the language.

→ Here we use Wordnet lemmatizer
(lexical Dictionary)

Icon: Lemmatization ⚡ Searching शब्दी शब्दों को सॉल्व करना

CHP : 2



R
EPRESENTATION.

- Converting text to Number

FAMOUS Quote in

ML "Garbage In Garbage Out"

PAGE NO.	/ /
DATE	

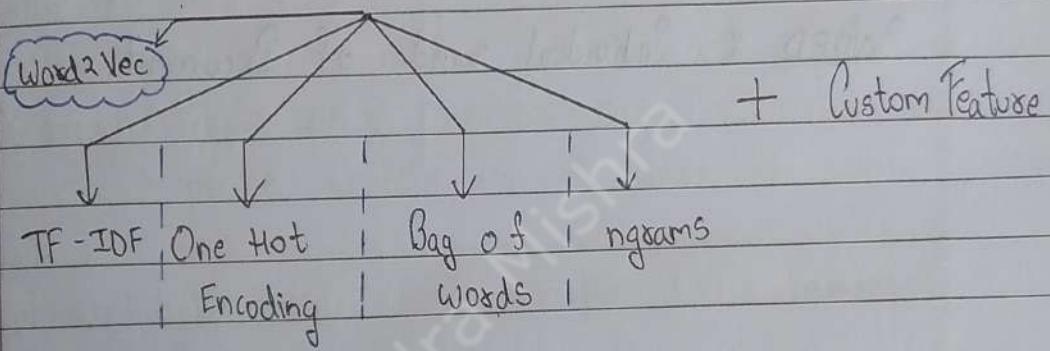
What is Feature Extraction from text?

Text vectorization / Text Representation

Ans

eg: ① Input :- Audio Clip Output:- Emotions.

② Hello how are you? (sentiment Analysis)



COMMON TERMS

- 1) Corpus
- 2) Vocabulary
- 3) Document
- 4) Word.

Combining all words of Dataset is called as "**Corpus**"

① Hello

② Hi

③ Come College

⇒ Corpus: "Hello Hi Come College"

Vocabulary: Unique Words of Corpus

Document: Individual Review

Here,	eg:	Hello	1 st Document
		Hi	2 nd Document

WORD: Individual word of Document.

ONE HOT ENCODING

D1	People watch campusx
D2	Campusx watch Campusx
D3	people write comment
D4	Campusx write comment

Corpus : People watch campusx campusx watch campusx
 people write comment campusx write comment

Vocabulary : People watch campusx write comment.



$$V = 5$$

T _i :	People	watch	campusx	write	comment

People watch campusx
 ↗ ↗ ↗
 $D_1 = [[1, 0, 0, 0, 0], [0, 1, 0, 0, 0], [0, 0, 1, 0, 0]]$

$$\text{Shape of } D_1 = (3, 5)$$

We are numbering "1" as per above T_i representation.

$$D_2 = [[0, 0, 1, 0, 0], [0, 1, 0, 0, 0], [0, 0, 1, 0, 0]]$$

Pros

- Intuitive
- Easy to Implement

(Cons)

- Sparse Array
(causes Overfitting)

→ Variable Document Size
(ML won't work.)

→ OOV (Out of Vocabulary)

→ No capturing of Semantic

Understanding Sparse Arrays.

Array in which only one element is one and rest all are zero's.

eg:	[[1,0,0,0] [0,0,1,0] [0,1,0,0]]
-----	--

It becomes difficult for ML model to handle Sparse Array because it contains too many zero's causing overfitting. Thus in real world scenarios you won't find using One Hot Encoding. eg: $V = 50,000$.

Understanding Variable Document Size

PAGE NO.	
DATE	///

e.g. D1 : Hello Rabindra

D2 : Rabindra is Reading.

Vocabulary → Hello Rabindra is Reading.

For Document	D1 :	Hello	Rabindra	is	Reading
	↓	1	0	0	0
"Hello Rabindra"		0	1	0	0

D1 : [[1,0,0,0], [0,1,0,0]] (3,4)

For Document	D2 :	Hello	Rabindra	is	Reading
		0	1	0	0
"Rabindra is Reading"		0	0	1	0
		0	0	0	1

D2 = [[0,1,0,0], [0,0,1,0], [0,0,0,1]] (3,5).

Variable length of D1 and D2.

⑩ Out of Vocabulary Problem.

		Prediction
D1	Hello Rabindra	0
D2	Rabindra is Reading	1

But suppose you have given input as "Hi Rabindra" for prediction it would not be able to predict it properly.
Because

Hi is not present in Vocabulary.

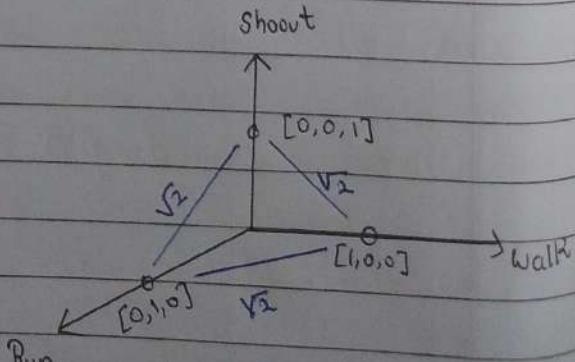
Hello	Rabindra	is	Reading
-------	----------	----	---------

⑪ No capturing semantic.

Walk | Run | Shoot

Walk [1, 0, 0] } (vectors)
 Run [0, 1, 0]
 Shout [0, 0, 1]

Plotting Vectors



Now, although all vectors are equidistant but they differ in meaning thus unable to capture semantic meaning

Because Run and walk are different and has distance of $\sqrt{2}$ but shout also has same distance although shout is completely dissimilar than (Run and walk).

Meaning of words is unable to be represented by numbers.

Bag of Words

Mostly used Technique in NLP

Used for Text Classification

D1	People watch	1
D2	Campusx watch	0
D3	People write	•1
D4	Campusx write	0

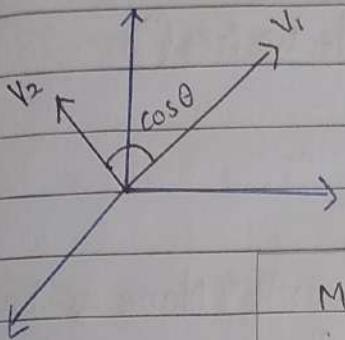
$$V = 4$$

	People	Watch	Campusx	Write	
D1	1	1	0	0	
D2	0	1	1	0	
D3	1	0	0	1	
D4	0	0	1	1	

Intuition: Same word being repeated same frequency
of times (order of words does not matter)
(context does not matter)

{ Vdim Plotting Here, $V=4$. }

Name is given Bag of words due to it is understood that collecting all the words inside a Bag and then assigning it as per Vtable (vocabulary table).



If angle between v_1 and v_2 is small then they will come under one category.

More the angle More is dis similarity

In this we will use `CountVectorizer()`

from sklearn.feature_extraction.text import CountVectorizer
`cv = CountVectorizer()`

`bow = cv.fit_transform(df['text'])`

fitting data and transforming it.

`print(cv.vocabulary_)`
 ↓ attribute of cv.

`print(bow[0].toarray())`

`print(bow[1].toarray())`.

Handles OOV problem.

`cv.transform(["campus is Rabindra"]).toarray()`

► Hyperparameter of CountVectorizer()

CountVectorizer(*, stop_words = None, lowercase = True)
we can mention language to remove
stopwords.

Preprocessor = None, binary = False, max_features,
n_gram = (1,1)
↳ #unigram.

① Understand binary = True.

If hyperparameter binary is set to True it means all non zero counts are set to 1

ex:

Vocabulary Table

Rabindra	is	Coding
1	2	0

binary = False

Vocabulary Table

Rabindra	is	Coding
1	1	0

binary = True

(either 1 or 0)

- ★ If we want to find frequency of words binary = False
- ★ If we only want to check if word exists or not binary = True.

for sentiment analysis we make binary = True as per Research.

⑩ Understanding max_feature

Suppose if, max_feature = 1

In this case it means we only want one feature as an output whose frequency is the highest

(eg :-)

Rabindra would be output

max_feature = 2

Rabindra	is	Coding
2	1	0

In this case here it will only include "Rabindra" and "is" while eliminate "coding"

This is used for eliminating rare occurrence words.

Advantages.

→ Simple and Intuitive

→ Fixed Document size

→ OOV

→ Some semantic meaning is preserved

Disadvantages.

→ Sparsity

→ OOV ('Hi' is new word)

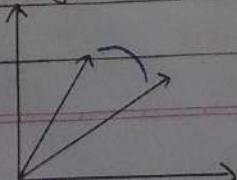
In this we ignore new word
However, Hi is bringing new information.

→ Ordering is not considered

Understanding Major flaws of Bag of Words

e.g. → This is a good movie } In both the sentence
 This is not a good movie. } there is only a
 minor difference in
 vocabulary.

In this case Bag of Words will ignore the "not" word. and both the vectors would be very much close showing similarity.



"But meaning is very much different although vocabulary is very much similar"

This problem is been handled by

Bag of N-grams

N-grams :

D ₁	Rabindra watch cartoon	1
D ₂	Cartoon watch Cartoon	1
D ₃	Rabindra write comment	0
D ₄	cartoon write comment	0

Bigram:

Take two words at a time.

Rabindra watch cartoon.
 ↑ ↑

Tri-gram:

Rabindra watch cartoon
 ↑ ↑

Similarly, n-gram

Bag of bi-gram:

Here, One vocabulary is combination of Two words.

e.g.:

D1	Rabindra watch Cartoon	1
D2	Cartoon watch Cartoon	1
D3	Rabindra write Comment	0
D4	Cartoon write Comment	0

Vocabulary:

Rabindra	watch	watch	Cartoon
Cartoon	watch	watch	Cartoon
Rabindra	write	write	Comment
Cartoon	write	write	Comment

Vocabulary Table:

People	Watch	Watch	Campus	x	Cartoon	watch	watch	Cartoon	Rabindra,
									write

D1	1	1	0	0	0	0	0	0	0
D2	0	1	1	0	0	0	0	0	0
D3	0	0	0	0	1	1	0	0	0
D4	0	0	0	0	0	1	0	0	1

If that Bigram is present in Document increment the count else keep zero.

Trigrams:

people watch campusx | campusx watch campusx | D3 | D4

D1	1	0	0	0
D2	0	1	0	0
D3	0	0	1	0
D4	0	0	0	1

① CountVectorizer (n_gram_range = (2,2))



This means only Bigram

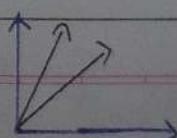
ngram_range = (1,2) # This means both unigram
and Bigrams

* Problem Overcome! *

Bag of Words

	This	movie	is	very	good	not
D1	1	1	1	1	1	0
D2	1	1	1	0	1	1

Now, vectors are very close.



Using Bigram Approach :-

This movie is very good
 This movie is not good

Using Bigram Vocabulary :-

This movie movie is is very very good is not not good						
D1	1	1	1	1	0	0
D2	1	1	0	0	1	1
		↑	↑	↑	↑	↑

Rabin Karp

Different Values

Now, vectors are more far away.

Advantages

1) Able to capture semantic of the sentence

2) Easy implement.

Disadvantages

2} For large Datasets

unigram

Bigram
Trigram.

May increase dimension
(slows down Algorithm)

4} No solution for Out of Vocabulary.

TF-IDF

D1	People watch campusx	1
D2	campusx watch campusx	1
D3	people write comment	0
D4	campusx write comment	0

People	watch	campusx	write	Comment
1	1	1	0	0

Unlike other Techniques Tf-IDf assigns different "Weights"

Idea: / Intuition

जिसके लिए वह word ज्योति ही आ रहा है। इन अलग से words के लिए वर्तमान में यह बहुत अच्छा है कि यह एक particular Document में उपलब्ध होना चाहिए लेकिन यह एक whole corpus में उपलब्ध न हो।

In such scenario that word becomes important for the Document and is given more importance.

TF (Term Frequency)

IDF (Inverse Document Frequency)

Q.) How to Calculate TF, IDF?

Ans: $TF(t,d) = \frac{\text{Number of occurrences of term } t \text{ in document } d}{\text{Total number of terms in document } d}$

$IDF(t) = \log_e \left(\frac{\text{Total number of Documents in the Corpus}}{\text{Number of documents with term } t \text{ in them}} + 1 \right)$

e.g. $TF(\text{people}, D1) = 1/3$

{ Ek aur People aya hai
Total 3 terms hai }

$TF(\text{campus}, D2) = 2/3$

Note:

$0 < TF < 1$

can consider it as probability.

$IDF(\text{campus}) = \log \left(\frac{4}{3} \right)$

	IDF
People	$\log(4/2)$
Watch	$\log(4/2)$
CampusX	$\log(4/3)$
Write	$\log(4/2)$
Comment	$\log(4/2)$

Analyzing Trend behind TF-IDF :-

if the term is more frequent in Corpus \Rightarrow IDF ↓

$$\text{eg: } \log\left(\frac{4}{4}\right) = \log(1) = 0 + 1$$

if the term is less frequent in corpus \Rightarrow IDF ↑

TF is basically probability of frequency of element in a Document.

(Vocabulary Table) FORMULA: $TF \times IDF$

	People	Watch	CampusX	Write	Comment
D1	$1/3 \times 0.3$	$1/3 \times 0.3$	$1/3 \times 0.125$	0	0
D2	0	$1/3 \times 0.3$	$2/3 \times 0.125$	0	0
D3	$1/3 \times 0.3$	0	0	$1/3 \times 0.3$	$1/3 \times 0.3$
D4	0	0	$1/3 \times 0.125$	$1/3 \times 0.3$	$1/3 \times 0.3$

Code Implementation :

```
from sklearn.feature_extraction.text import TfidfVectorizer
tfidf = TfidfVectorizer()
```

```
rabindra = tfidf.fit_transform(df['text']).toarray()
rabindra.
```

```
print(tfidf.idf_) # Prints IDF Values.
```

```
print(tfidf.get_feature_names_out())
```

"+1" is added because terms occurring with all documents may get ignored without "+1"

e.g. $\log\left(\frac{N}{n}\right) = 0 \quad TF \times IDF = 0 \quad (\text{ignored})$



BRAINSTORMING

Q. Why do we use log for calculating IDF?

Ans

Suppose, there is one word which is very rare out of 10,000

$$IDF = \left(\frac{N}{n} \right) = \frac{10,000}{1} \quad | \quad \log\left(\frac{10000}{1}\right) = 4$$

Using log will prevent unnecessary big value for IDF

In such case $Tf \times \text{IDF}$ while $0 < Tf < 1$ and
 \downarrow
 1000

in this case IDF will dominate over Tf and then
 Tf will not have any significance.

Advantages

→ Information
 Retrieval

Disadvantages

→ Sparsity
 → OOV
 → Dimension
 → Semantic

Custom Features

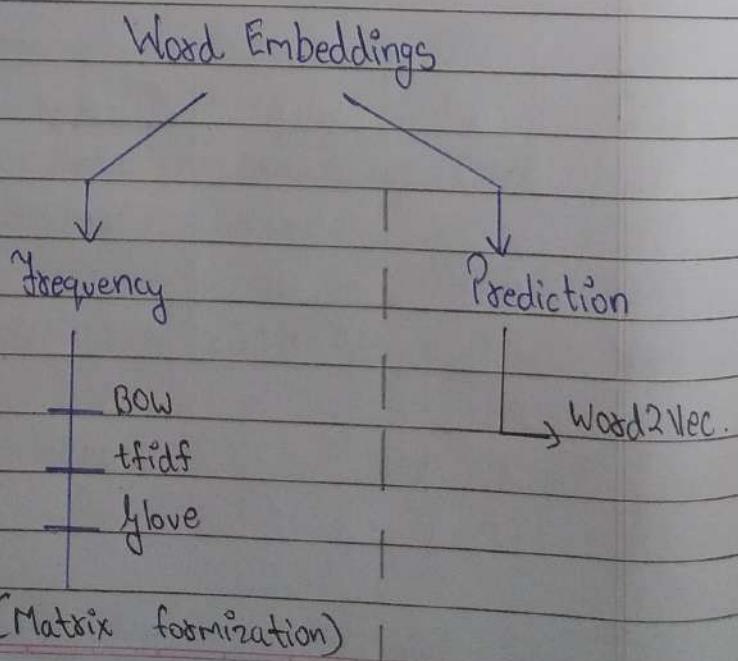
Such as: ① Ratio of +ve/-ve ② Word Count
 ③ No. of +ve words, -ve words

WORD2Vec (Deep Learning)

Understanding Word Embeddings:

In NLP processing, word embedding is a term used for the representation of words for text analysis, typically in the form of a real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in the meaning.

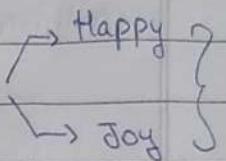
eg: We have a word "King" and we want to convert it to vector ie [24 19]



Q.] What is Word2Vec? (# Complex)

a) Developed by Google engineers in 2013

b) Word2Vec can capture Relationships like



ie it can capture semantic meaning.

c) Creates low dimensional vectors usually in range of [100 - 300] as compared to Bow/TF-IDF

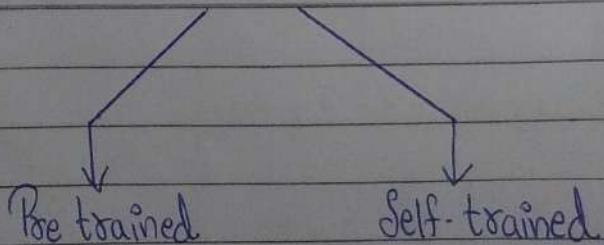
d.) In case of Bow/TF-IDF you'll see a very sparse Matrix [0 0 0 2 0 0.]

But in Word2Vec. this matrix is more dense ie less no. of zero and thus reduces overfitting problem.

e.) It is a deep learning Technique.

DEMO

There are two ways of using Word2Vec.



We will use the pre-trained weights of word2vec that was trained on Google News corpus containing 3 billion words. This model consists of 300-dimensional vectors for 3 million words and phrases.

Here, it is trying to say that we will have vector form of 30 lakh words as a dimension of 300
ie. [0.1... 300]

Thus it makes this word2vec accurate.

```
import genism
from genism.models import Word2Vec, KeyedVectors
```

Important Commands

① pip install

② pip list

③ pip show package-name. Search by package name

Shows all packages installed

pip: preferred installer program.

PyPi: Python Package Index.

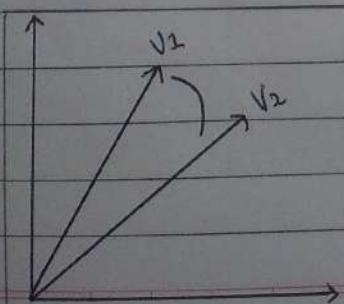
→ It is a central repository for python programming language.

Pip installs packages from this default repository from PyPi.

Word2Vec

For time being consider 300 Dimensional vector as 2D so we can visualize.

Now, we can use these values to find similarity of two words



We can also perform
 ① Vector Addition
 ② Vector Sub
 and other vector operations

most_similar(): Function

It returns list of all similar vectors to the input word

eg: `model.most_similar('Man')`

O/P: `[('woman', 0.76640129...),
 ('men', 0.8367...),
 ('teen', 0.7347...),
 ('girl', 0.7446...),
]`

→ vector values

doesn't_match():

gives oddman out.

eg: `model.doesnt_match(['cricket', 'NASA', 'ISRO'])`

O/P: Cricket

⑩ Vector Operation.

$\text{vec} = \text{model}[\text{'Wicket'}] + \text{model}[\text{'RUN'}] - \text{model}[\text{'Ball'}]$
 $\text{model. most_similar}([\text{vec}])$

$\text{vec2} = \text{model}[\text{'INR'}] - \text{model}[\text{'USD'}]$
 $\text{model. most_similar}([\text{vec2}])$

⑪ INTUITION of Word2Vec

- It captures semantic meaning
- We need to convert words in such a way that it should capture the meaning.

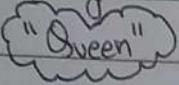
In this we convert words based on different features.

Features.	King	Queen	Man	Woman	Monkey.
Gender	1	0	1	0	1
Wealth	1	1	0.3	0.3	0
Power	1	0.7	0.2	0.2	0
Weight	0.8	0.4	0.6	0.5	0.3
Speak	1	1	1	1	0

King [1 1 1 0.8 1]

Monkey [1 0 0 0.3 0]

Understanding Vector Operations

[King] - [Monkey] will give us something close to
 "Queen"

Problem

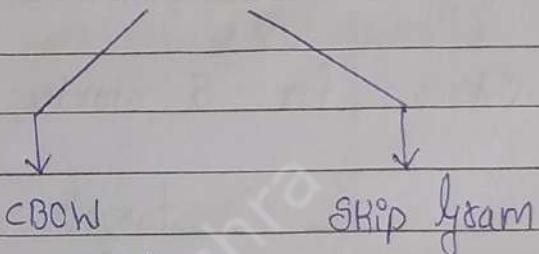
- 1) In real-world instead of 5 we will have lakhs of vocabulary so handcrafting it will be impossible
- 2) Thus we need neural network to create features automatically. But disadvantage is you will not know what exactly are the features.
- 3) It has generally 300 features (300 Dimension)

The underlying assumption of word2vec is that two words sharing similar contexts also share a similar meaning and consequently a similar vector representation from the model.

(eg) Football Player took a shot
 Hockey Player took a shot.

Here, Football \leftrightarrow Hockey is been used for same context.

(Two Architectures Word2Vec)



Both are shallow neural networks.

CBOW Architecture

Here, we try to solve a fake problem and solve the problem as a by product we get our vectors.

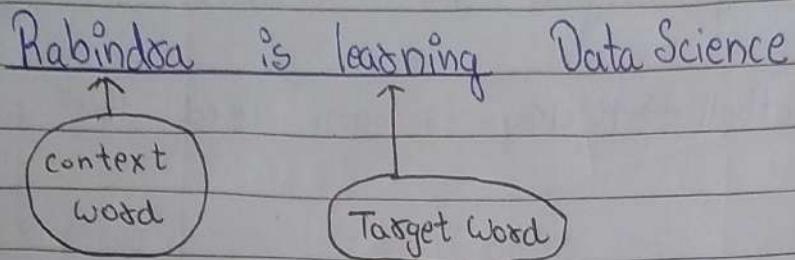
Fake Problem

\hookrightarrow Solve \rightarrow vector \rightarrow By product

Rabindra is learning Data Science.

To convert all the words into vector we will use a dummy problem.

Assuming a window of 3 words.



target words always comes with context word

Windows for 5 words

target

a) Now we will take this window for training our Data

X	Y
watch	
Rabindra, is	learning.

b) Now, we will take another window

Rabindra is learning DataScience

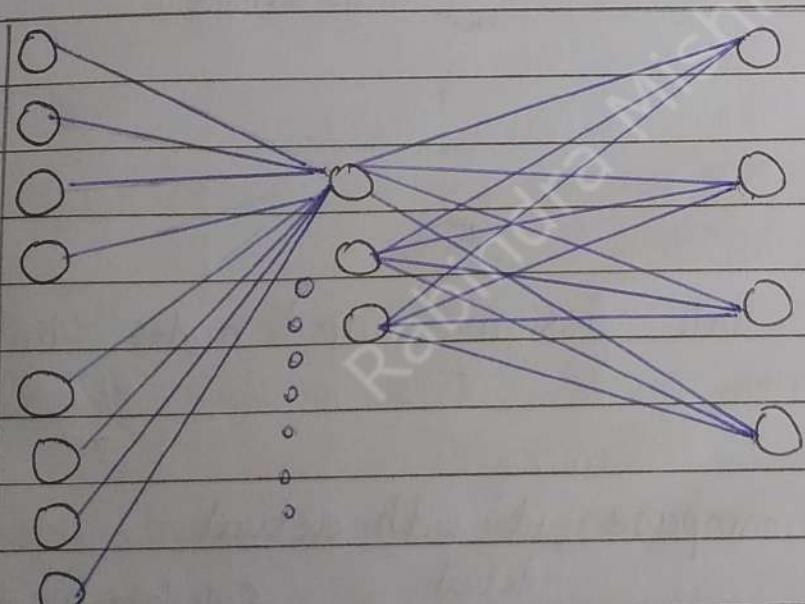
X (Training Data)	Y
Rabindra, is is, learning	Learning DataScience

Then we will convert our Input to "One hot Encoded vectors."

1000	}
0100	
0010	
0001	

One Hot Encoded Vectors.

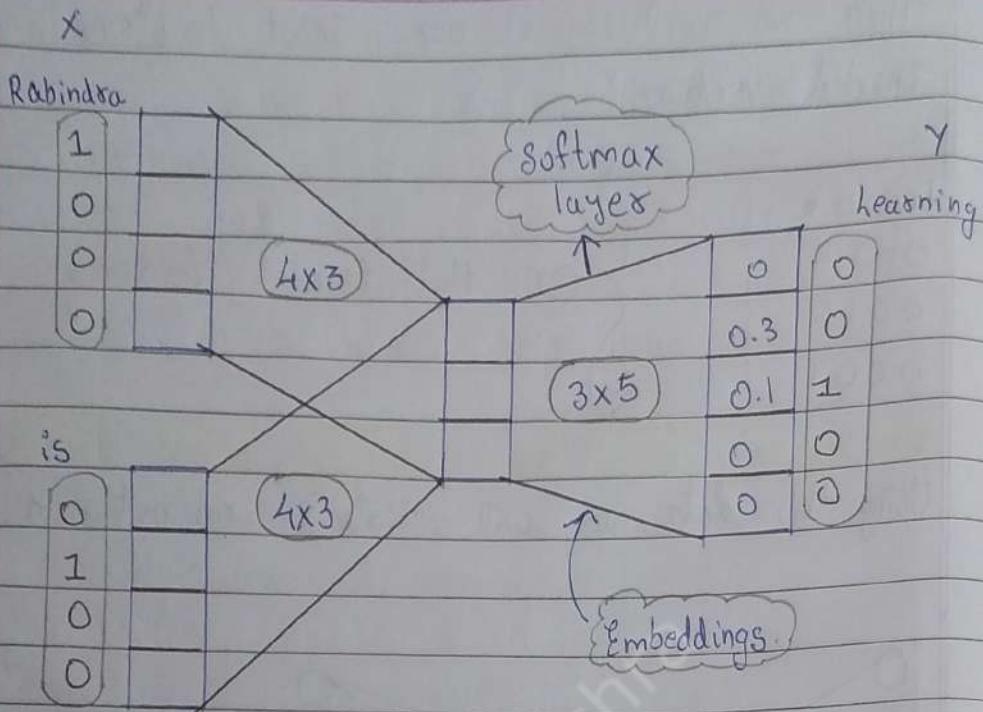
Using this data we will create neural network.



(Hidden layer)

Here 3 nodes are there in hidden layer representing Window size.

In output layer we have 4 words represented by individual node.



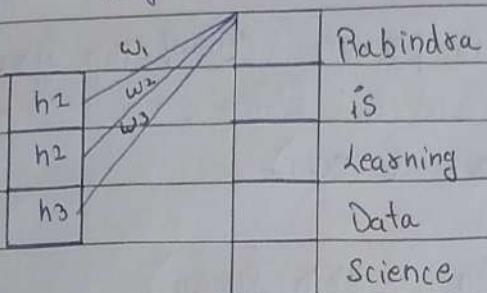
Softmax layer will assign output nodes in terms of probability.

Now, it will compare it with actual

	Actual	Calculate Loss
ie	0	0
	0.3	1
	0.1	1
	0	0
	0	0

After calculating loss we will backpropagate for multiple epoch. to adjust the weights.

Fully Connected Neural Network



Vector Representation of Rabindra: $[h_1 * w_1 + h_2 * w_2 + h_3 * w_3]$

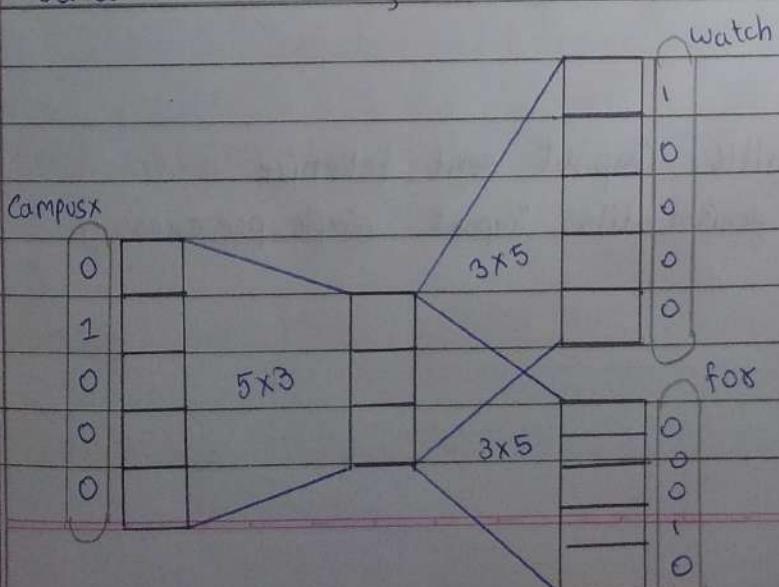
(SKIP-GRAM ARCHITECTURE)

Opposite of CBOW

In this also we will use dummy Problem

From target word you need to predict Context word

Campusx	\rightarrow	watch, for
for	\rightarrow	campusx, data
Data	\rightarrow	for, science



From Research it is proven CBOW/Skip gram

→ Small Data (CBOW)

→ Large Data (Skip gram)

Word2Vec Embeddings (Improve)

- Increase the training Data
- Increase Dimension of Vector
- Increase Window Size

"Training your own Model"

import genism

import numpy as np

import pandas as pd

import os

from nltk import sent_tokenize

import genism.utils import simple_preprocess

Story = []

Data = r"C:\Users\Game of Throne"

for filename in os.listdir(Data):

f = open(os.path.join(Data, filename))

corpus = f.read()

raw_sent = sent_tokenize(corpus)

for sent in raw_sent:

story.append(simple_preprocess(sent))

defining our Model structure.

model = gensim.models.Word2Vec(

window=10

min_count=2 # sentence having atleast 2 words

model.build_vocab(story) # Vocab build

model.train(story, total_examples=model.corpus_count,
epochs=Model.epochs) # training Model

total_examples: tells total sentences in corpus

Accessing our Model

→ model.wv.most_similar('King')

L> word2vec

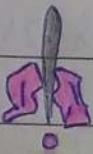
doesn't_match():

→ model.wv.doesnt_match([**'jon'**, '**tikon**', '**robb**', '**arya**'])
O/P → **'arya'**

Calculating Similarity

1) model.wv.similarity(**'arya'**, **'sansa'**)
0.8534194

2) model.wv.similarity(**'tywin'**, **'sansa'**)
0.26445806

PROBLEMALERT

Plotting 100 dimensional vectors will be complex
thus we will need to reduce it to 2 Dim / 3 Dim.

→ Apply PCA.

Getting Vector Representation of all words

⇒ model.wv.get_normed_vectors().shape
(17453, 100)

Here, we have 17453 words in our corpus.

V.VIMP
② model.wv.get_normed_vectors()

O/P :-

[[-0.119, 0.3176, 0.77, ..., 0.138],
[0.765, 0.117, 0.712, ..., 0.551],
:
:
:
:
:
[0.6179, 0.932, 0.003, ..., 0.05]

])

i.e 1743 words of 100 Dimension.

To see which words are presented for this particular vector.

y = model.wv.index_to_key

for docstring in jupyter Notebook
click Shift + tab

PAGE NO.	/ / /
DATE	

Plotting the graph

→ Applying PCA

from sklearn.decomposition import PCA

pca = PCA (n_components=3) #3D

X = pca.fit_transform (model.wv.get_normed_vectors)
X [:5]

import plotly.express as px

effect of PCA

X. Shape

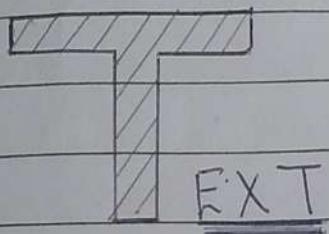
O/p: (17453, 3)

↳ Reduce Dimension

fig = px.scatter_3d (X [:100], x=0, y=1, z=2, color=y [:100])
fig.show ()

Important Visualization libraries

- Matplotlib
- Seaborn
- Plotly (for interactive plots, Dashboard)



C CLASSIFICATION

Classification on Different Data.

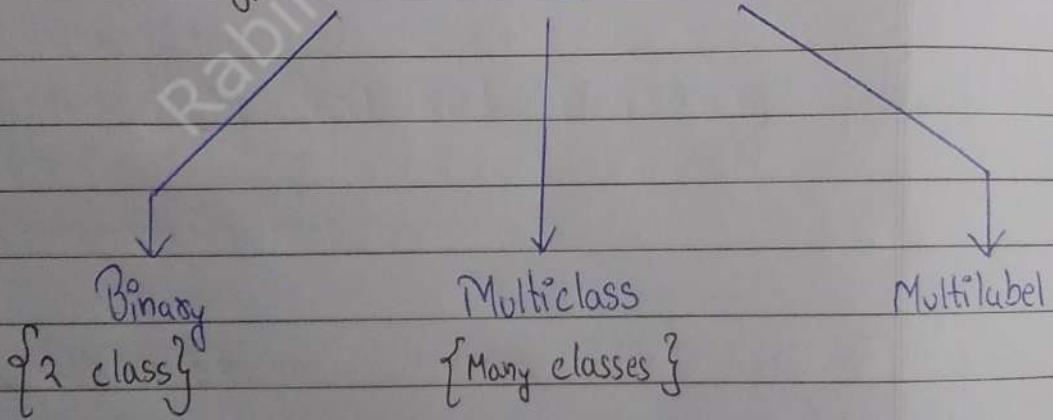
tabular

image

Text

It is a supervised ML Task

(Types of Text Classification)



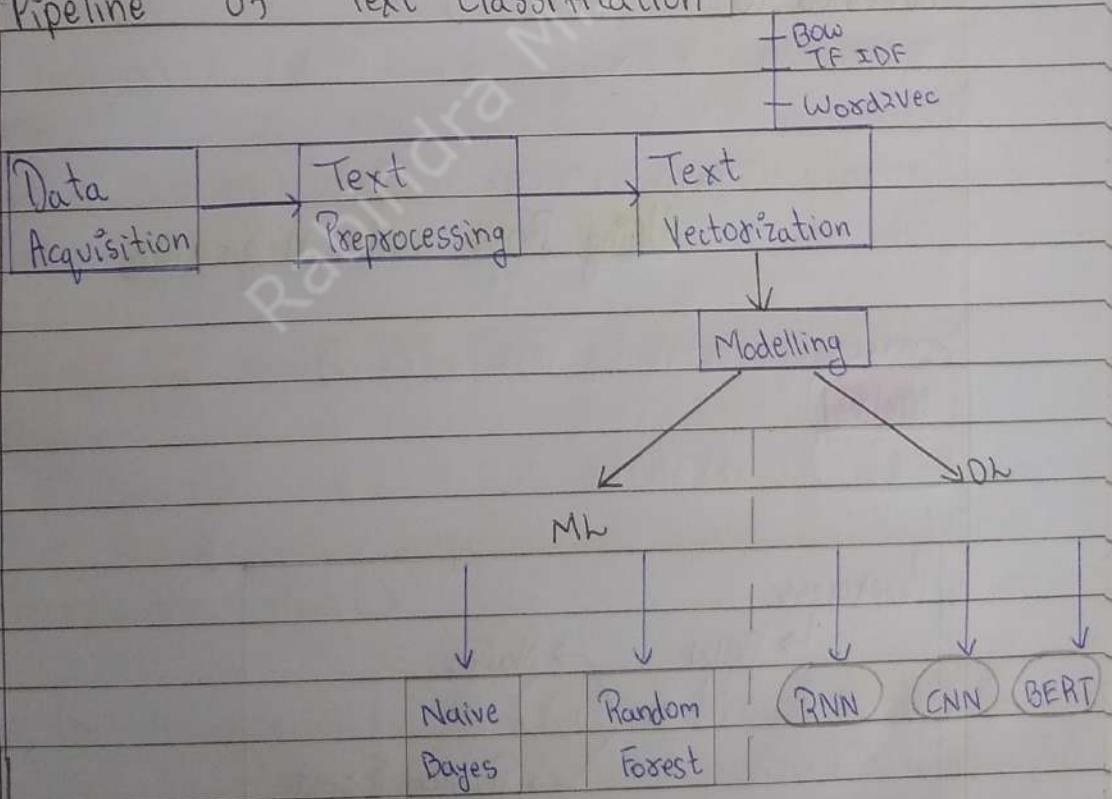
Imp: NLP is mostly used for text classification.

Applications

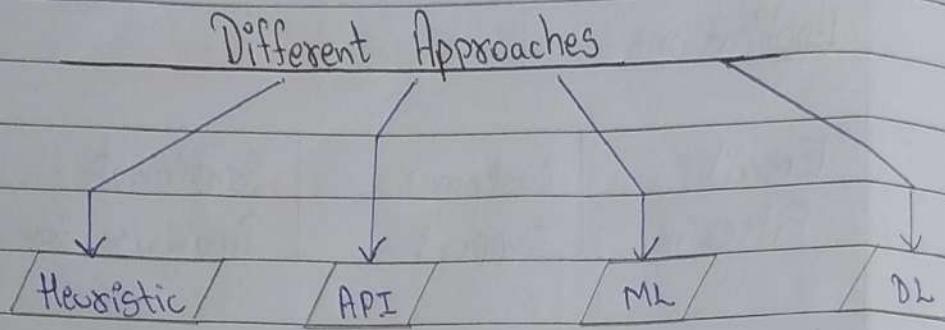
Email Filtering	Customer Support	Sentiment Analysis for Review
-----------------	------------------	-------------------------------

Language Detection	Fake news Detection
--------------------	---------------------

⑩ Pipeline of Text Classification



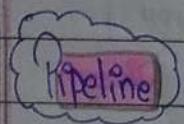
Evaluation → accuracy
→ confusion



When we don't have enough data we need to use heuristic approach

NLPcloud.io : API

Using Bow and n-grams.



Preprocess

↳ Bow

↳

→ Naive

→ Random Forest

```
import numpy as np
import pandas as pd
import os
```

① df['sentiment'].value_counts()

positive : 5028

negative : 4972

② df.isnull().sum() # checking if null

review 0

sentiment 0

③ df.duplicated().sum() # checking duplicates

17

XX Fitting our model as BOW

```
from sklearn.feature_extraction.text import CountVectorizer
cv = CountVectorizer()
```

```
x_train_bow = cv.fit_transform(x_train['review']).toarray()
```

```
x_test_bow = cv.transform(x_test['review']).toarray()
```

```
print(x_train_bow.shape, x_test_bow.shape)
```

Apply Naive Bayes

```
from sklearn.naive_bayes import GaussianNB  
gnb = GaussianNB()
```

```
# fitting our model  
gnb.fit(x_train_bow, y_train)
```

```
# Predicting our model  
y_pred = gnb.predict(x_test_bow)
```

```
# Calculate Accuracy of Model
```

```
from sklearn.metrics import accuracy_score,  
confusion_matrix  
  
accuracy_score(y_test, y_pred)
```

O/P: → 0.62844266

Apply Random Forest

```
from sklearn.ensemble import RandomForestClassifier  
rf = RandomForestClassifier()
```

```
# fitting Model  
rf.fit(x_train_bow, y_train)
```

Predicting the Model

$y_{pred} = \text{rf.predict}(x_{test_bow})$

Calculate Accuracy.

accuracy_score(y_test, y_pred)

Using Word2Vec

Pre trained
Word2Vec

Train on your
own Data →
Enough
Data

Vocab

Generally we convert all the words as a vector

example :-

Delhi → []

↳ 300 Dimension

But we haven't seen anything getting like eg Document as a vector

ie Document 1 : []

Example →

I	Love	NLP
↓	↓	↓
[]	[]	[]

Adding all these vectors will give us new vectors

[-----]

```
def document_vector(doc):
    # Remove out of vocabulary word
```

```
doc = [word for word in doc.split() if word in
       model.wv.index_to_key]
```

```
return np.mean(model.wv[doc], axis=0)
```

```
from tqdm import tqdm
```

tqdm is a popular python library used to display progress bars for loops, making it easy to monitor the progress of lengthy computations.

* Practical Advise *

- 1) Ensemble Technique
- 2) Heuristic features
- 3) Deep Learning.

{ POS TAGGING }

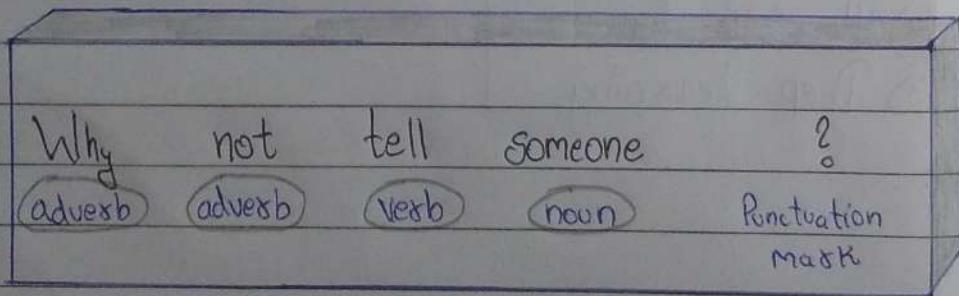
It is necessary for complex system like chatbots, information retrieval system.

Here, we will learn about how does POS tagging work internally.

Pos tagging means giving a part of speech to a sentence word by word.

Definition : In simple words, we can say that Pos tagging is a task of labelling each word in a sentence with its appropriate part of speech.

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



Note : It is a preprocessing step

★ Major Applications ★

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

Done using "(Spacy.)"

import spacy

nlp = Spacy.load('en_core_web_sm')

doc = nlp(u"g will google about tabindra")
 ↳ (This denotes unicode string)

doc[0]

O/P: → I

doc[0].pos_	# g is Pronoun
-------------	----------------

O/P: 'PRON'

doc[0].tag_	# Present Participle.
-------------	-----------------------

O/P: 'PRP'

Note: pos_ gives course great parts of speech

tag_ gives find great parts of speech

(This tells us about in more detail)

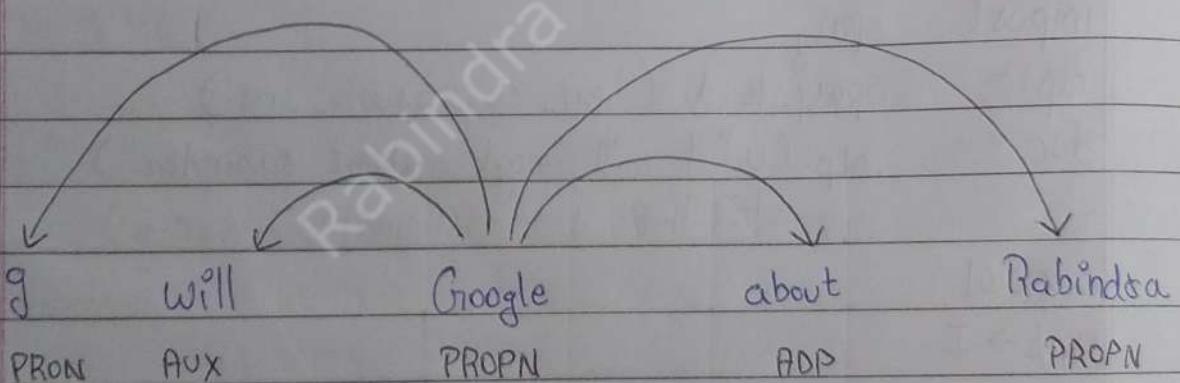
* To understand Spacy we can write explain() function *

- ① Spacy.explain('NNP')
- ② spac.explain('PRP')
- OR...

→ Pos Tagging Visualization

→ from spacy import displacy.

displacy.render(doc, style='dep', jupyter=True)



Working of POS TAGGING ?

Soln:

Famous Algorithm HMM



Hidden Markov Model

Consider a statement

i) Will will google chess

↳ (here Will is name of a person)

But for predicting the output for above statement
we need to train it on some Data.

eg:	Nitish loves chess
	Can nitish google CampusX
	Will Ankita google chess
	Ankita loves Will

Using this data we will train our HMM model to predict
above statement "Will will google chess."

ie we will create pos tagging of above statements and
then try to create pos of our statement
"Will will google chess"

Understanding our Statement wise POS

(N)	(V)	(N)
Nitish	loves	Chess
(M)	(N)	(V)
Can	Nitish	google
(M)	(N)	(V)
Will	Ankita	google chess
(N)	(V)	(N)
Will	Ankita	Loves Will
(N)	(V)	(N)

NS Noun { Rows will contain all unique words
 VS Verb { while column will contain all unique Pos
 MS Model

	N	M	V
Nitish	2	0	0
loves	0	0	2+1
chess	3	0	0
Google	1	0	2
Will	2	1	0
Ankita	2	0	0
can	0	1	0

Here, we need to calculate Emission Probability.
 Divide by total of the column.

Emission Probability

	N	M	V	
Nitish	2/10	0	0	
loves	0	0	3/5	
chess	3/10	0	0	
google	4/10	0	2/5	
will	2/10	1/2	0	
Aankita	2/10	0	0	
can	0	1/2	0	

उपर दिए गए word is Noun (N) probability of it being Nitish is [2/10]

TRANSITION Probability

(N) → (V) → (N)

i.e. probability of transition i.e. Nitish loves chess

Noun की अंगता word Verb होने की संभावना

- Verb की अंगता word Noun होने की संभावना
- Model की अंगता word Verb होने की Probability
- Noun की अंगता word Model होने की Probability

To calculate this we will assign 'Start' and 'End' to each and every sentence.

	(N)	(V)	(N)	
S	Nitish	loves	Chess	E
(M)	(N)	(V)	(N)	
S	Can	Nitish	Google	Chess E
(M)	(N)	(V)	(N)	
S	Will	Ankita	Google	Chess E
(N)	(N)	(N)		
S	Ankita	loves	Will	E
(N)	(V)	(N)		
S	Will	Loves	Google	E

TRANSITION TABLE

		N	M	V	E	END
Start	S	3	2	0	0	
	N	0	0	5	5	
	M	2	0	0	0	
	V	5	0	0	0	

ie Start Ke baad Noun 3 bar aya hai

Noun Ke baad Verb 5 bar aya hai

:

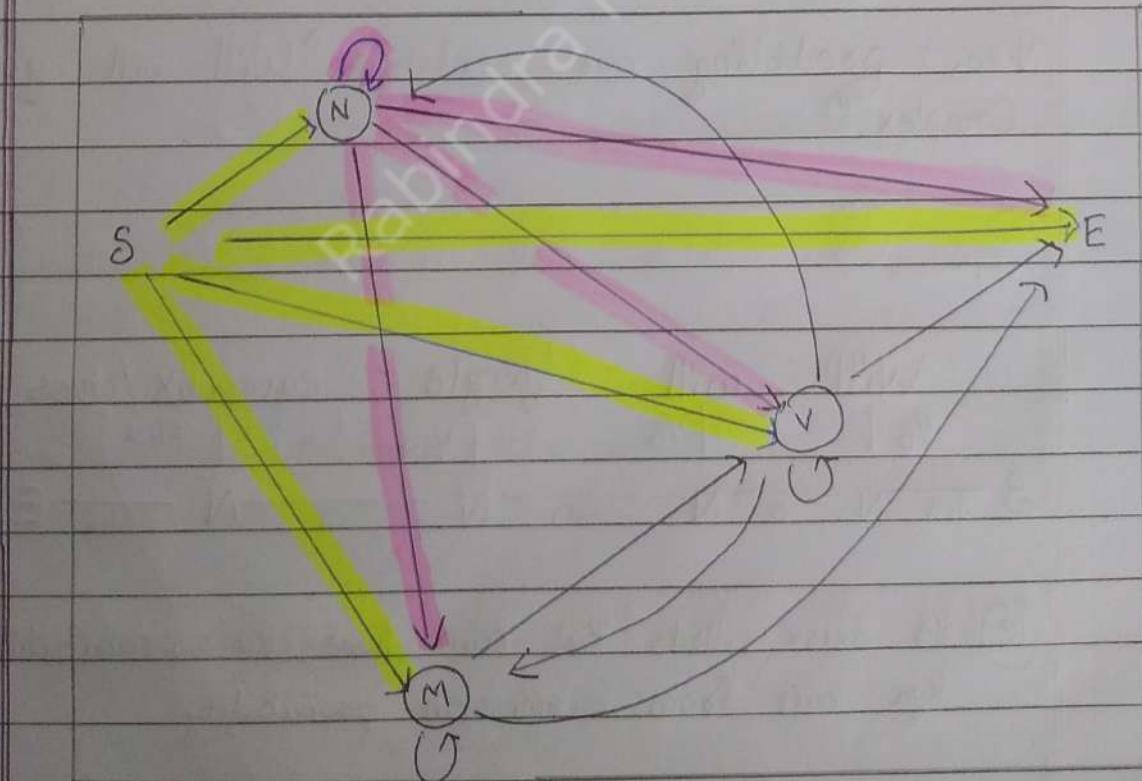
Calculating Probability

Divide it by total of "Row" and not Column.

Transition Probability

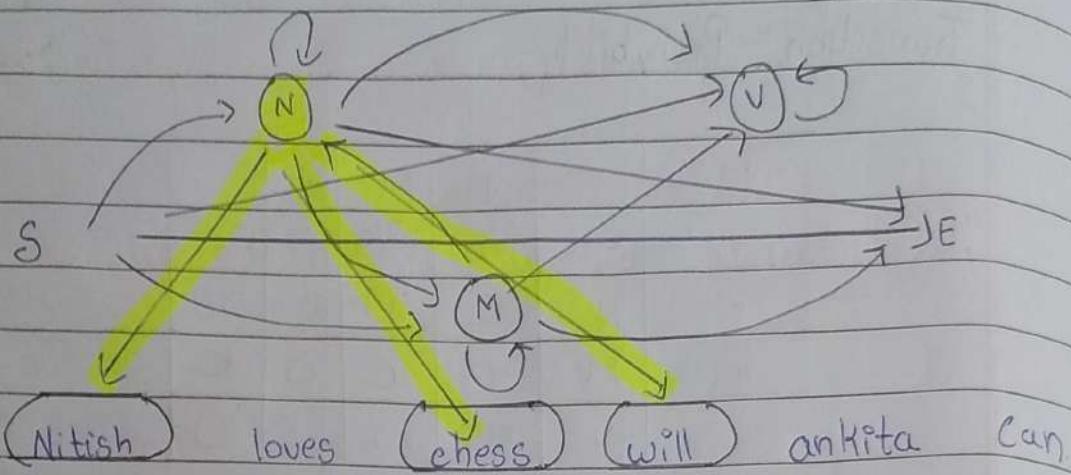
	N	M	V	E
S	3/5	2/5	0	0
N	0	0	5/0	5/0
M	2/2	0	0	0
V	5/5	0	0	0

Using "Emission" Table and "Transition" Table we can create HMM Model



Highlighting All vertices from start, and all vertices originating from N (NOUN).

This what is called Transition Diagram.



i.e Probability of **“Noun”** with respect to “**Nitish**”,
“chess” and **“will”**

Now, predicting our sentence “Will will google
 Complex.”

Will	will	google	campusx / chess
1/5	1/5	1/10	3/10

S — $\frac{3}{5}$ N — O N — O N — O N — $\frac{5}{10}$ E

(ie) ye saare words ka noun hone ka probability.
 Kya hai from emission probability.

will will google chess

(3500 200 900)
3500 100 3500

$$S - \underset{8/10}{N} - \underset{10}{N} - \underset{0}{N} - \underset{0}{N} - \underset{5/10}{N} - E = \emptyset$$

$$S - N - N - N - N - \cancel{N} - E = \emptyset$$

$$S - \underset{2/5}{M} - \underset{1}{N} - \underset{1/2}{V} - \underset{1}{N} - \underset{1/2}{N} - E \neq \emptyset$$

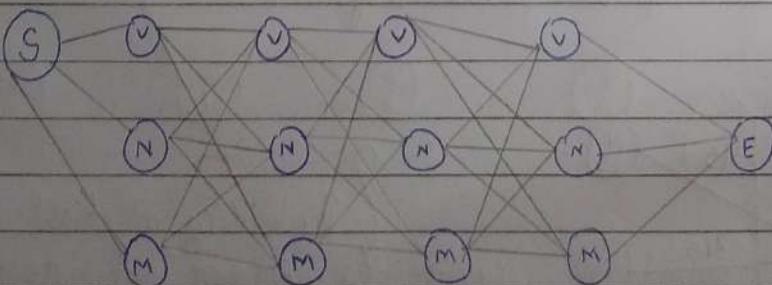
In this case, we will use the above

(S - M - N - V - N - E) combination for the statement
"Will will google chess"

Problem:

For Large Number of words No. of combinations would increase

e.g. Will Will google chess



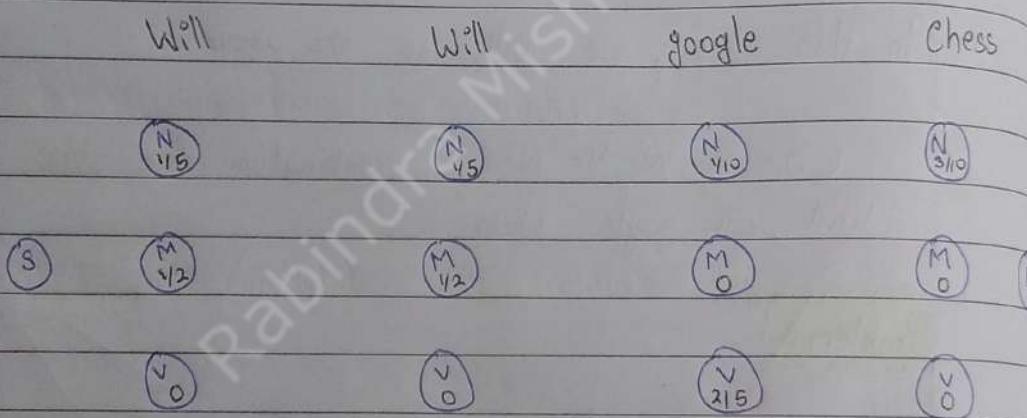
Total combinations is = V^W (V = No. of POS
 W = No. of words)

$$= 3^4 = 81$$

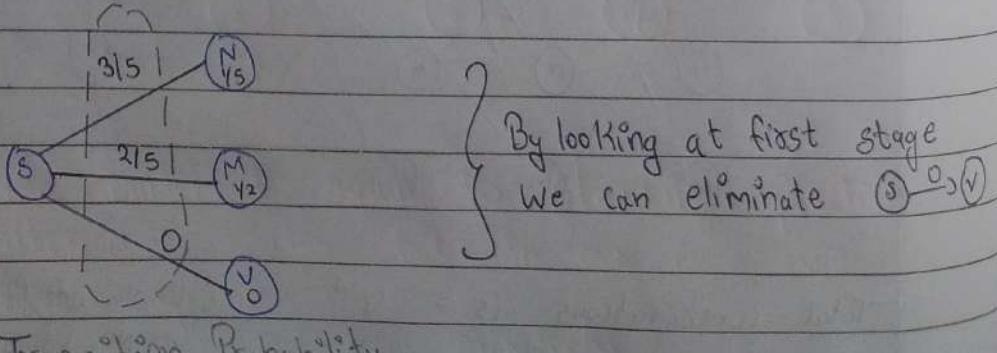
So as soon as number will increase our combinations will grow exponentially.

For optimization we can use Viterbi Algorithm.

Viterbi Algorithm



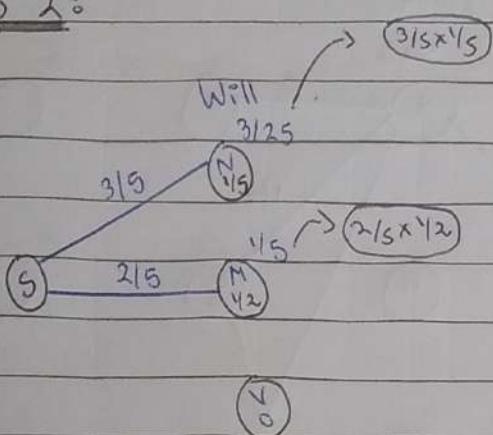
Step 1: We will first consider the part with highest probability.



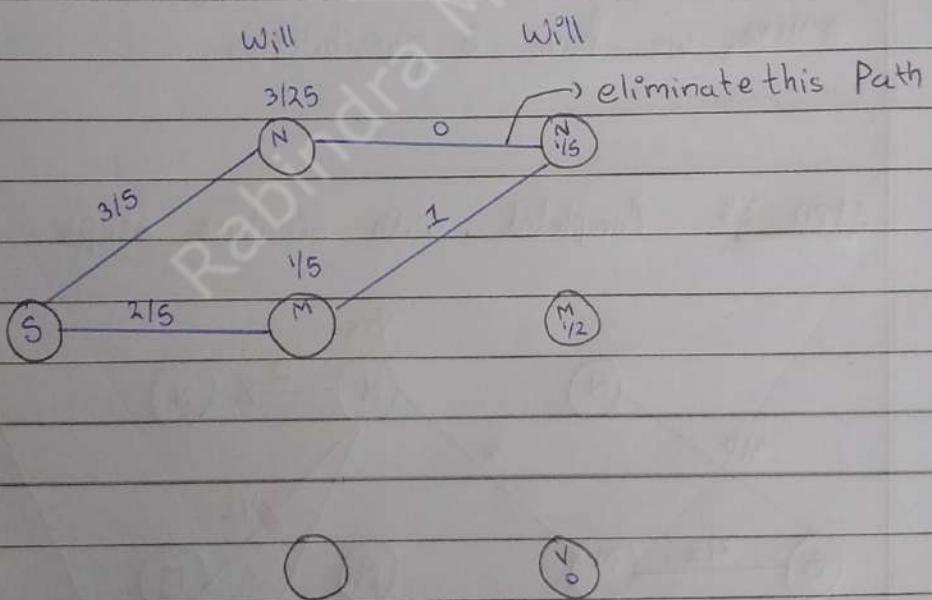
Transition Probability.

Step 2:

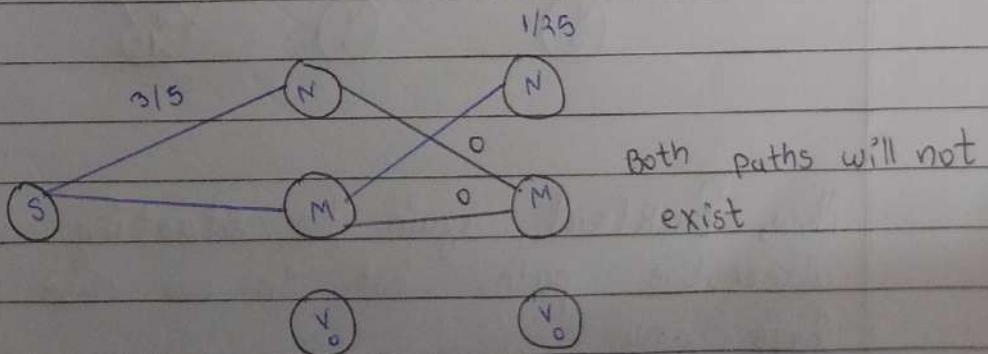
(i)



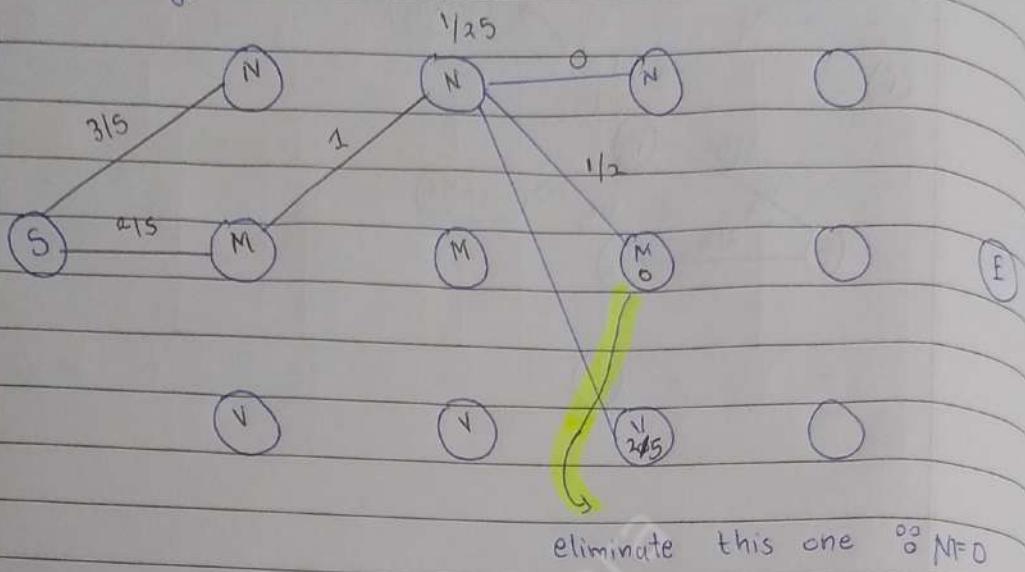
(ii)



(iii)

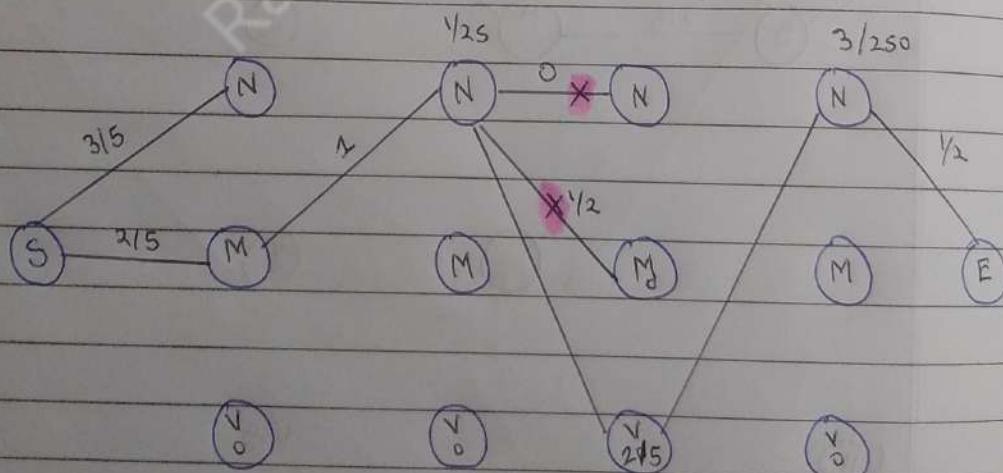


Similarly, for second layer.



Here, we have 3 possibilities

Step 3: Completed with similar steps.



Now, Backtrack again to starting position for alternative paths but in our case we don't have any.

So Final Output : M/N/V/N → POS Tagging.

So in this algorithm we are able to reduce certain possibilities.

(PROJECT) :- i) Duplicate Question Pairs

2) Spam Classified.

Dataset: Quora (A very good project).