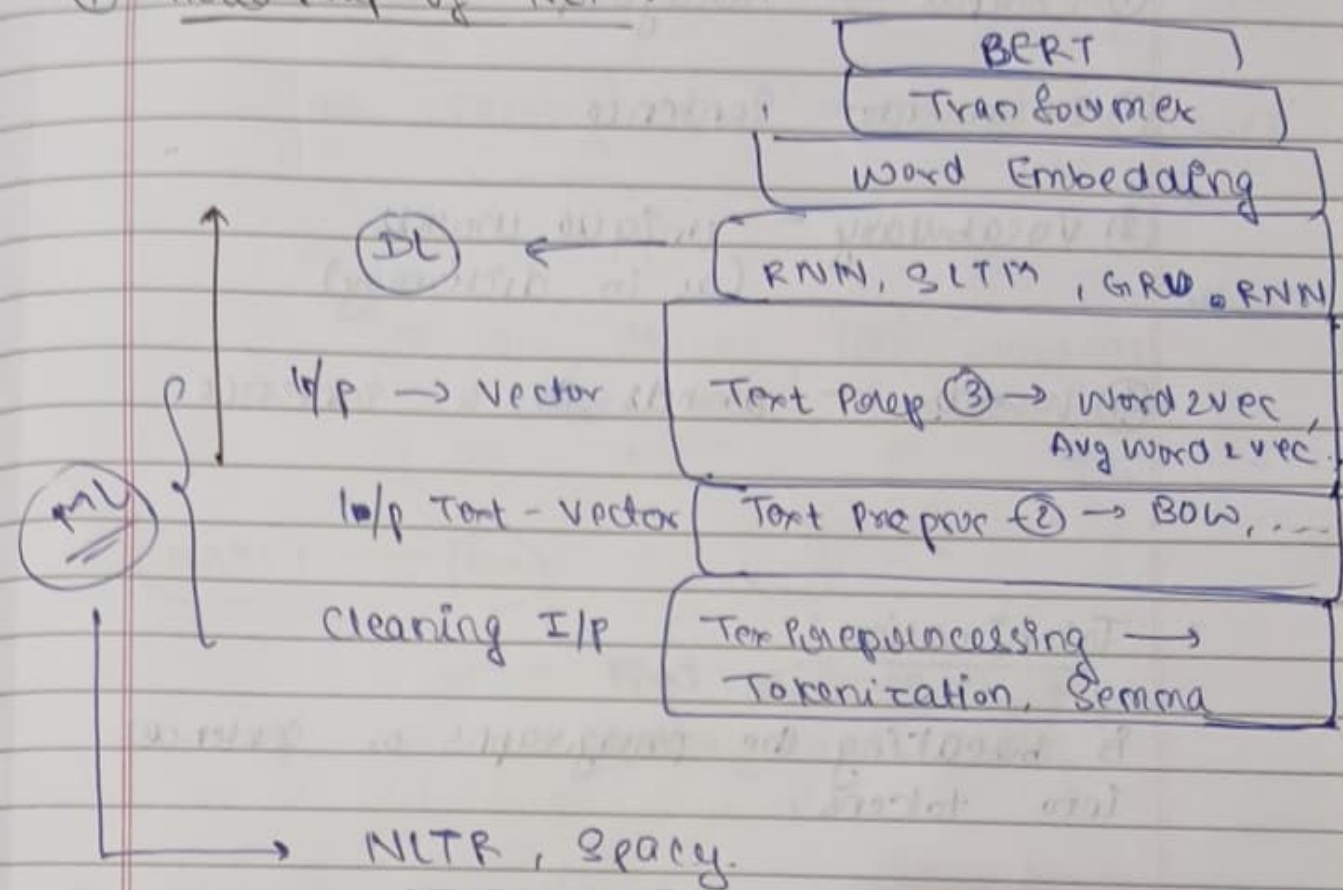


## NLP (Soc)

Complete NLP machine learning in one shot.

— krish naik (YouTube)

### ⊕ Road map of NLP :-



# # TOKENIZATION IN NLP :-

## ④ Topics

① Corpus — Paragraph

② Documents — Sentence

③ Vocabulary — Unique words  
(as in dictionary)

④ Words — words in a sentence.

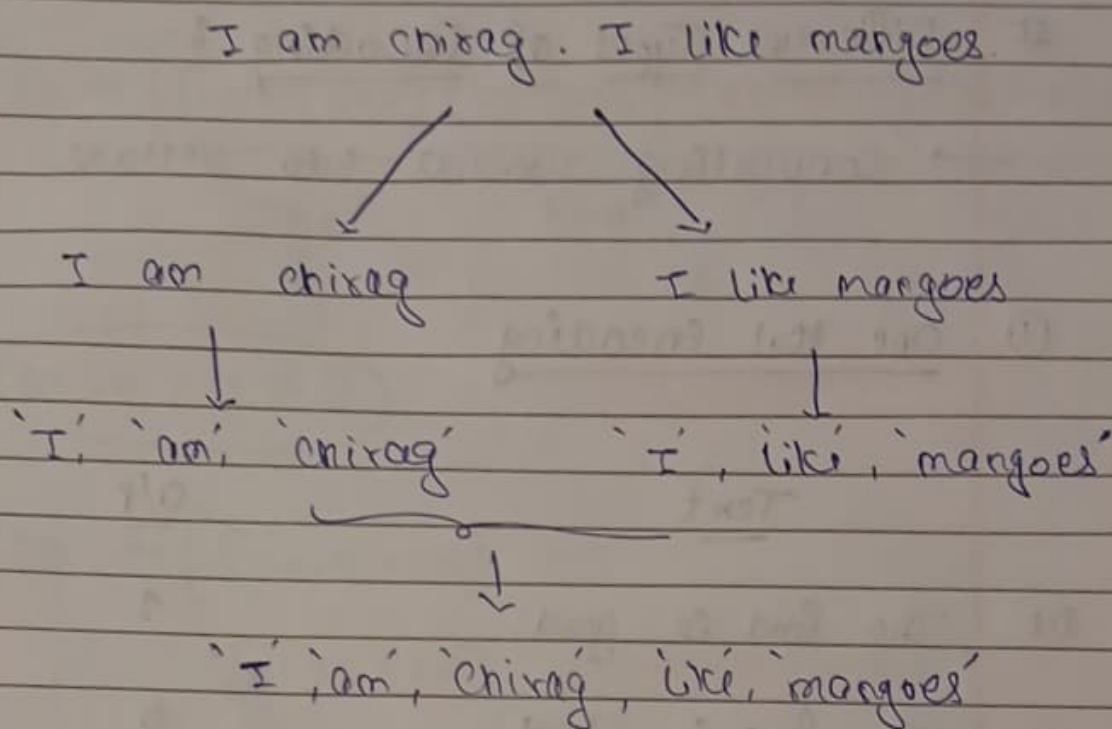
## • Tokenization :

Is breaking the paragraphs or sentences into tokens.

i.e.

{ paragraphs } → { sentences } → { words }

↓  
{ Unique words }



NOTE-

HW

→ Find differences b/w  
Spacey and NLTK

## # Stemming

→ Process of reducing a word to its word stem that affixes to suffixes or prefixes to the root of words known as lemma.

# # Different Types of Encoding :-

→ Converting words into vectors.

## ① One Hot Encoding

	<u>Text</u>	<u>O/P</u>
D1	The food is good	1
D2	The food is bad	0
D3	Pizza is amazing	1

Vocabulary {unique words}

↳ The food is good bad Pizza amazing

⇓

1	0	0	0	0	0	0
0	1	0	0	0	0	0
0	0	1	0	0	0	0
⋮	⋮	⋮	⋮	⋮	⋮	⋮

→ Number of unique vocabulary

$$V = 7$$

∴ we will have 7 vectors to represent



$$D1 \Rightarrow \left[ \underbrace{[1\ 0\ 0\ 0\ 0\ 0\ 0]}_{\text{The}}, \underbrace{[0\ 1\ 0\ 0\ 0\ 0\ 0]}_{\text{Food}}, \underbrace{[0\ 0\ 1\ 0\ 0\ 0\ 0]}_{\text{is}}, \underbrace{[0\ 0\ 0\ 1\ 0\ 0\ 0]}_{\text{good.}} \right]$$

$$D1 = 4 \times 7$$

$$D2 \Rightarrow \begin{bmatrix} [1\ 0\ 0\ 0\ 0\ 0\ 0] \\ [0\ 1\ 0\ 0\ 0\ 0\ 0] \\ [0\ 0\ 1\ 0\ 0\ 0\ 0] \\ [0\ 0\ 0\ 1\ 0\ 0\ 0] \end{bmatrix}$$

### Advantages

- Easy to implement with python  
(Sklearn and pandas)

### Disadvantages

- Sparse matrix → Overfitting  
!!  
Good accuracy for training but not for new data

- ML Algorithm → we need fixed size I/p
- No semantic meaning is getting captured
- out of Vocabulary → ~~if training data~~  
if test data has a new word.

## (2) Bag of words :-

- Sentiment classification,  
Mail Spam or Ham?, etc.

### Dataset

	Text	O/p
S1:	He is a good boy	1
S2:	She is a good girl	1
S3:	Boy and girl are good	1



lowering all word  
+  
stopwords

S1 : good boy      S2 : good girl      S3 : Boy girl good

vocabulary

frequency

good

3

boy

2

girl

2



→ maximum frequency is put first  
(i.e. descending order)

→ we select the top 10 or 20 features  
based on the frequency.

	good	boy	girl
S1 :	[ 1	1	0 ]
S2 :	[ 1	0	1 ]
S3 :	[ 1	1	1 ]

Binary Bow

and

Bow

{ 1 and 0 }  
only

{ Count is  
updated with  
frequency  
in statement }

## Advantages

- Simple and intuitive
- Fixed sized I/P - ML algorithm.

## Disadvantages

- Sparse matrix or arrays → overfitting
- Ordering of the word is getting changed
- Out of vocabulary (OOV)
- Semantic meaning is still not getting captured.

③ N - grams (Not in video)



## (2) TF - IDF

[Term Frequency - Inverse Document Frequency]

S<sub>1</sub>: good boy

S<sub>2</sub>: good girl.

S<sub>3</sub>: boy girl good.

$$\text{Term freq. (TF)} = \frac{\text{No. of rep of word}}{\text{No. of words in sent.}}$$

$$\text{IDF} = \log_e \left[ \frac{\text{No. of sentences}}{\text{No. of sentences containing the word}} \right]$$

Term Frequency

IDF

	S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>	words	
good	1/2	1/2	1/3	good	$\log_e(3/3) = 0$
boy	1/2	0	1/3	boy	$\log_e(3/2) =$
girl	0	1/2	1/3	girl.	$\log_e(3/2)$

## Final TF - IDF (Product)

	good	boy	girl
S1 :	0	$\frac{\log(3/2)}{2}$	0
S2 :	0	0	$\frac{1}{2} \log_e(\frac{3}{2})$
S3 :	0	$\frac{1}{3} \log_e(\frac{3}{2})$	$\frac{1}{3} \log_e(\frac{3}{2})$

### Advantages

- Intuitive
- Fixed size → Vocab Size
- Word importance is getting captured.

### Disadvantages

- Sparsity still exists
- OOV.

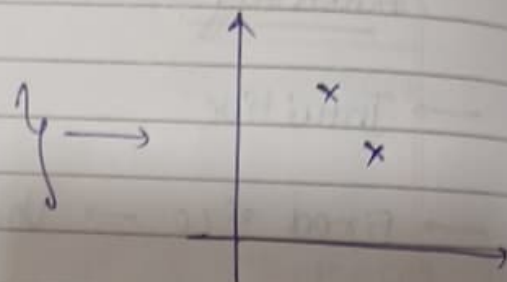
## # Word Embedding :-

→ 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 meaning.

↓  
[ wikipedia def<sup>n</sup> ]

Happy = [      ]

Excited = [      ]



Similar words.

### Word Embedding

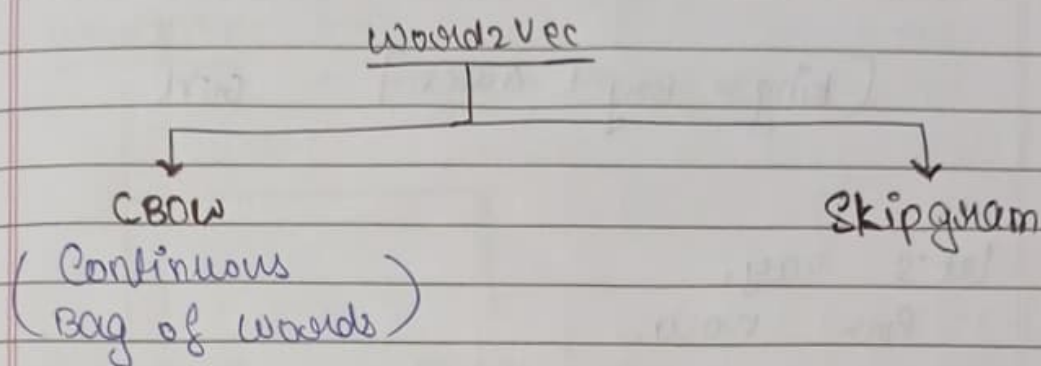
Count & Freq.

- ① OHE
- ② BOW
- ③ TF-IDF

Deep Learning Trained Model

⇕  
[ Word2Vec ]

→ Deep learning based models are very accurate and have very less the disadvantages.



# word2vec :-

→ by Google

→ Technique for NLP published in 2013.

→ Uses Neural Network model to learn word associations from a large corpus of text.

→ Each distinct word with a vector.

• Feature Representation

→ Vocabulary → unique words

Boy Girl King Queen Apple Mango

Gender  
Royal  
Age  
...

Features [300 dim]



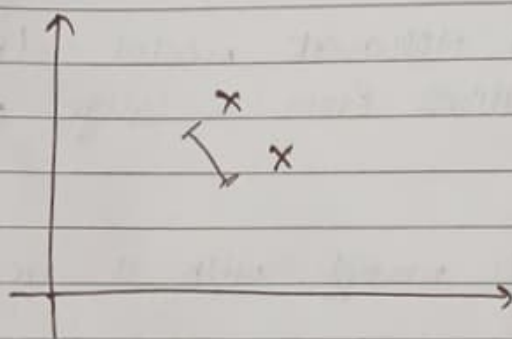
→ Each vector has a 300 dimension

$$[\text{king} - \text{Boy} + \text{Queen}] = \text{Girl}$$

let's say,  
for now,

we just have 2 dimensions,

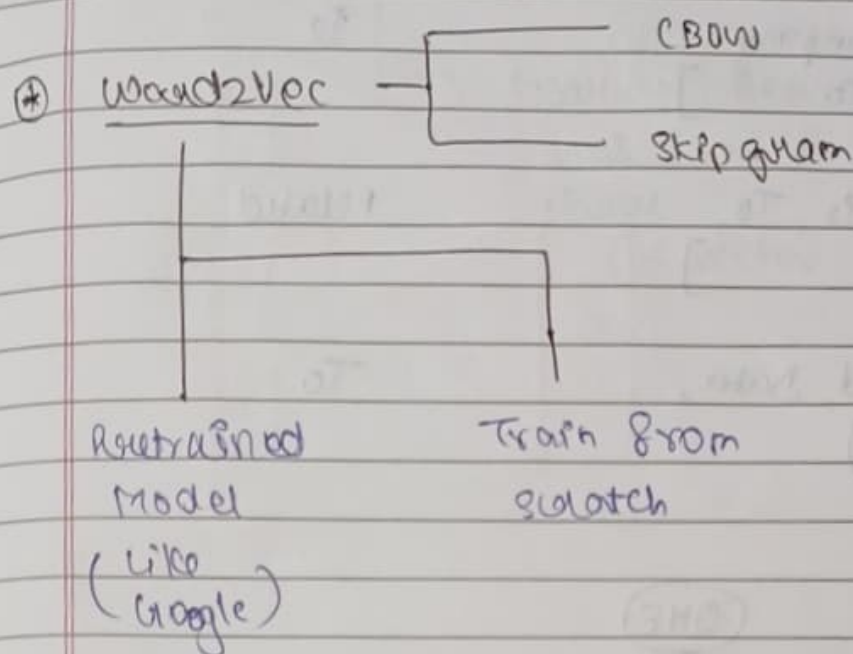
### • Cosine Similarity



$$\text{Distance} = 1 - (\text{Cosine Similarity})$$

~~Cosine~~

$$\boxed{\text{Cosine Similarity} = \cos 45^\circ}$$



# ① CBOW [Continuous Bag of Models]

CORPUS ÷ Dataset

[ "Ineuron Company is related To Data Science" ]

\* Choose a window size

Say [ window size = 5 ] (preferably odd)

then PTD →

I/p

O/p

→ [ Newton, Company  
Related To ]

Is

→ [ Company, Is, To,  
Data ]

Related

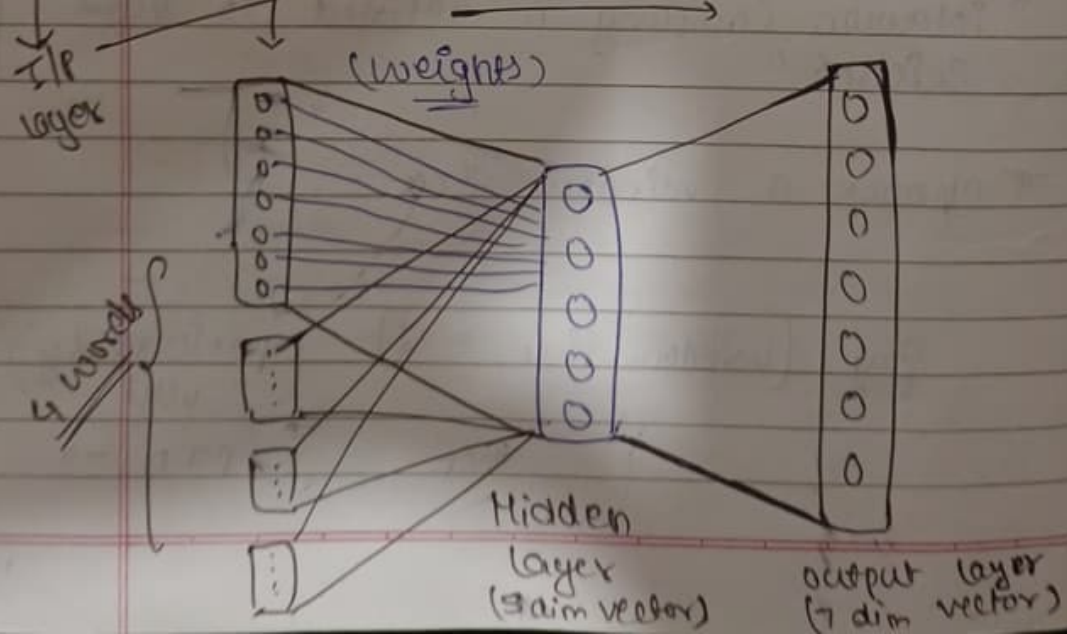
→ [ is, Related, Data,  
Science ]

To

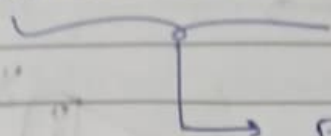
OHE

Newton	[ 1 0 0 0 0 0 0 ]
Company	[ 0 1 0 0 0 0 0 ]
Related	[ 0 0 0 1 0 0 0 ]
to	[ 0 0 0 0 1 0 0 ]

# CBOW { Fully Connected Neural Network }



→ Window Size = 5



Feature Representation  
has / or we want it to  
have  
(5 vector dim)

# Skpgram - word2vec

[window size = 5]

O/P

T/P

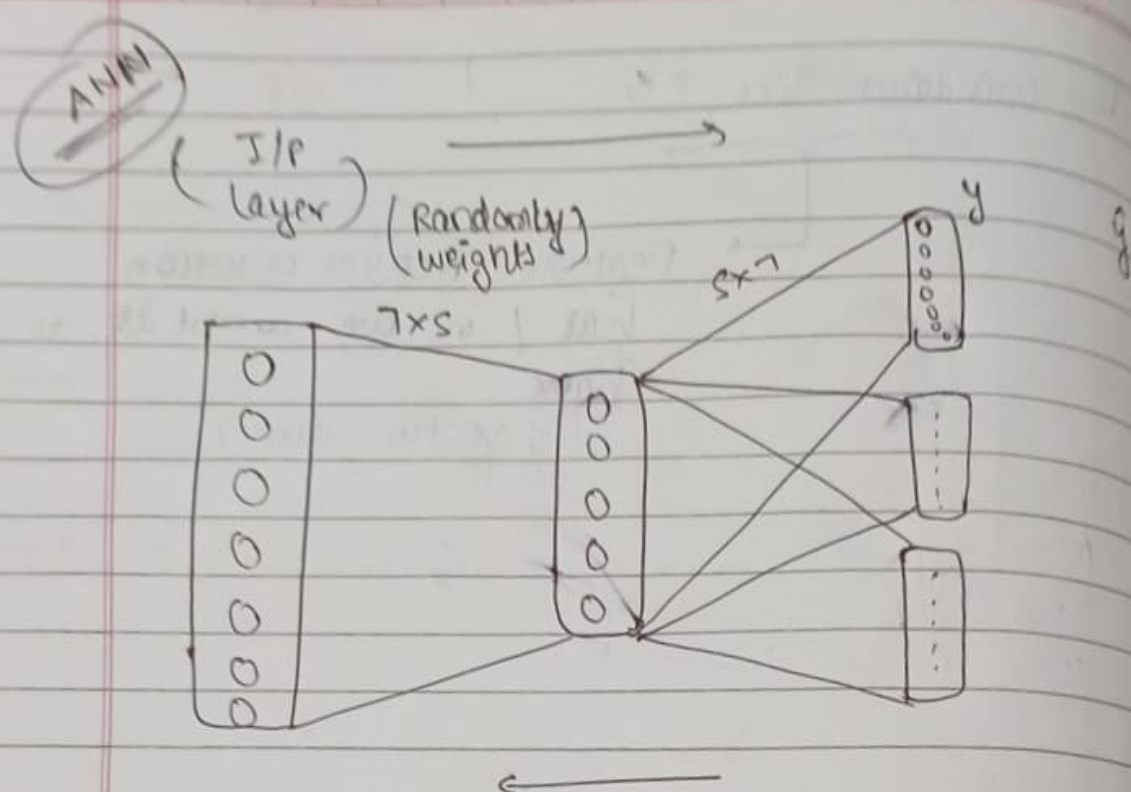
→ [iNewton, Company...  
.. Related, To]

IS

related

To





⊕ When to apply / use CBow or Skipgram

{ Small dataset — CBow  
 { Huge dataset — Skipgram

⊕ Improving

① Increase training data

② Increase the window size.

↳ Vector dimension is also increasing.

# Google Word2Vec.

3 - billion words — Google News

Feature representation of 300 dimension vector

cricket  $\rightarrow$  [ . . . . . ]

## ⊕ Avg Word2Vec.

[ The food is good ]  $\rightarrow$  D1

(o/p)

