

# NLP

Saturday, 28 December 2024 6:10 PM

Topics:

- ① Corpus — Paragraph
- 2) Documents — Sentences
- 3) Vocabulary — unique words
- 4) words.

## Tokenization:

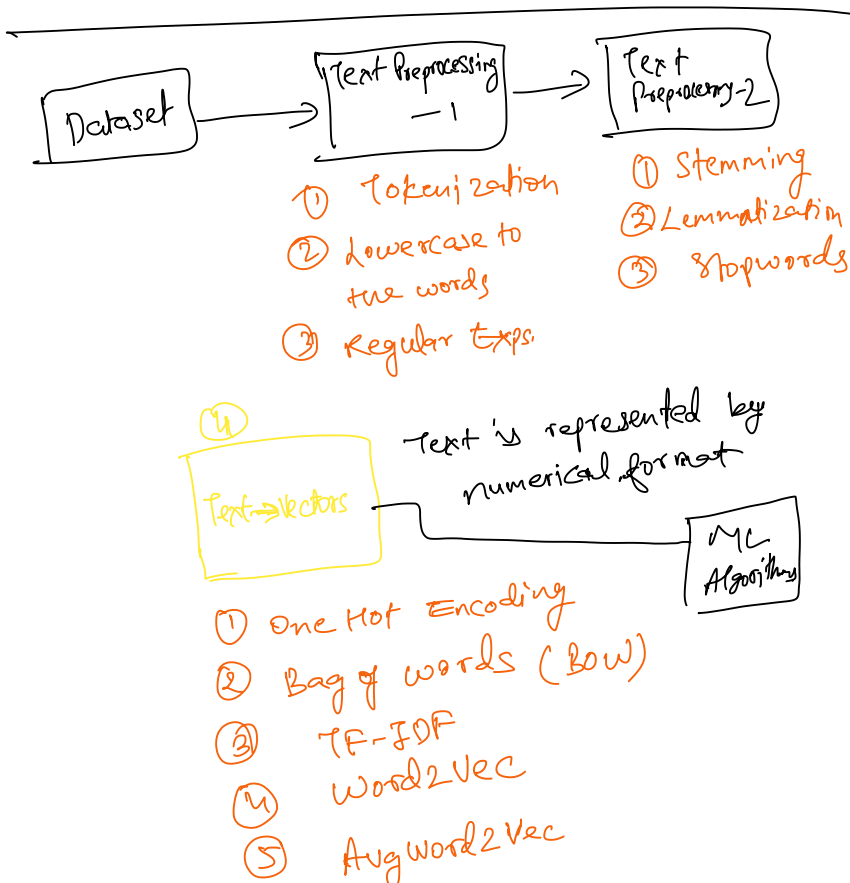
↳ convert either a paragraph or sentence into tokens.

## Vocabulary:

- ifobans
- ① I like to drink apple juice.
  - ② My friend likes mango juice.

total 11 words  
unique words: 10 { juice is repeating }

## # NLTK :-



# # One Hot Encoding:-

Text O/P

D1 The food is good 1

D2 The food is bad 0

D3 Pizza is Amazing 1

Unique vocabularies:-

→ The food is good bad pizza Amazing

the food → 1 0 0 0 0 0 0

D1 → 0 1 0 0 0 0 0

4x7 [ [1 0 0 0 0 0 0] → This  
[0 1 0 0 0 0 0] → food  
[0 0 1 0 0 0 0] → is  
[0 0 0 1 0 0 0] → good ]

D3 [ [0 0 0 0 0 1 0],  
[0 0 1 0 0 0 0],  
[0 0 0 0 0 0 1] ]

3x7

D2 [ [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] ]

4x7

→ one hot for each and every word.

Advantages	Disadvantages
① Easy to implement in python. (sklearn, OneHotEncoder...)	① Creates a sparse matrix → Overfitting → we need fix size for ML algorithm
	② No semantic meaning is captured. (cosine similarity...)
	③ Out of vocabulary. (unique words.)
	④ if we have 50k unique vocabulary size there will be huge sparse matrix.

## # Bag of Words! → Dataset:

Text	O/P
He is a good boy	1
She is a good girl	1
Boy and girl are good	1

lower all the words.  
stopwords

$s_1 \rightarrow$  good boy  
 $s_2 \rightarrow$  good girl  
 $s_3 \rightarrow$  Boy girl good

Vocabulary

good  
boy  
girl

frequency

3  
2  
2

⇒

	good	boy	girl	O/P
$s_1$	[1	1	0]	1
$s_2$	[1	0	1]	1
$s_3$	[1	1	1]	1

Binary BOW and BOW

{ 1 and 0 }

{ count will get updated based on frequency }

### Advantages:

- ① Simple & Intuitive
- ② fixed size input

### Disadvantages:

- ① sparse matrix  
↳ overfitting
- ② ordering of the words is getting changed.
- ③ Out of Vocabulary. (OOV).
- ④ Semantic meaning still not getting captured.

## # N Grams! → E.g. bigrams, trigrams

$s_1 \rightarrow$  The food is good

$s_2 \rightarrow$  The food is not good.

... is both the vector have a

The issue -  
 Little difference although both sentence  
 are opposite to each other.

Bigram - combination of two words.

	food	not	good	food good	food not	not good
S1 →	1	0	1	1	0	0
S2 →	1	<u>1</u>	1	<u>0</u>	<u>1</u>	<u>1</u>

→ Now we have a diff.

sklearn → n-grams (1,1) → unigrams  
 (1,2) → unigram bigram  
 (1,3) → unigram, bigram, trigram  
 (2,3) → Bigram, trigram.

# TF-IDF [Term Frequency - Inverse Document Frequency]

$$\text{Term freq} = \frac{\text{No. of rep. of words in sentence}}{\text{No. of words in sentence}}$$

$$\text{IDF} = \log_{10} \left( \frac{\text{No. of sentences}}{\text{No. of sentences containing the word}} \right)$$

Term freq.

	S1	S2	S3
good	1/2	1/2	1/2
boy	1/2	0	1/3
girl	0	1/2	1/3

IDF

$$\begin{aligned} \text{good} & \log_e(3/3) \\ \text{boy} & \log_e(3/2) \\ & \log_e(3/2) \end{aligned}$$

girl  $\log_e(1/2)$

Final TF-IDF

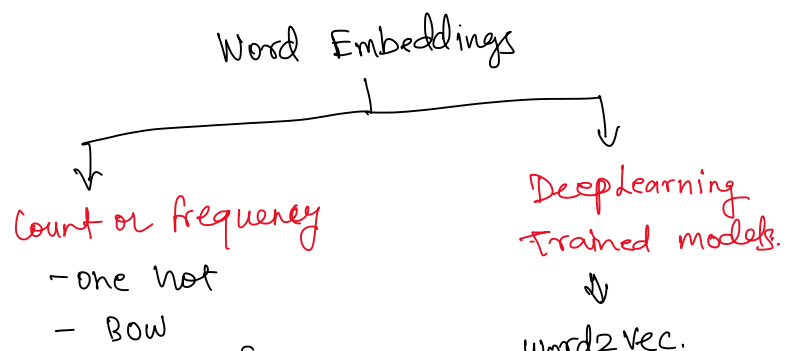
	$S_1$	$S_2$	$S_3$
good	$\frac{1}{2} \times 0$	$\frac{1}{2} \times 0$	$0$
boy	$\frac{1}{2} \times \log_e(\frac{2}{1})$	$0$	$\frac{1}{3} \times \log_e(\frac{2}{1})$
girl	$0$	$\frac{1}{2} \times \log_e(\frac{2}{1})$	$\frac{1}{3} \times \log_e(\frac{2}{1})$

$S_1 \begin{bmatrix} 0 & \frac{1}{2} \log_e \frac{2}{1} & 0 \end{bmatrix}$   
 $S_2 \begin{bmatrix} 0 & 0 & \frac{1}{2} \log_e \frac{2}{1} \end{bmatrix}$   
 $S_3 \begin{bmatrix} 0 & \frac{1}{3} \log_e \frac{2}{1} & \frac{1}{3} \log_e \frac{2}{1} \end{bmatrix}$

Advantages	Disadvantages.
① Intuitive ② Fixed Size ③ <u>word imp is getting captured.</u> (if word is in all sentences it is considered less important)	① Sparsity exists ② OOV.

# Word Embeddings! -

→ representation of words.  
 → in the form of real valued vectors.  
 that encodes the meaning of the words such that words that are closer in vector space are similar in meaning / semantics.



- TF-IDF



[Continuous Bag of Words]

# Word2Vec:- (Published by Google)

→ Word2Vec is a technique for natural language processing published in 2013 by Google.

→ The word2vec uses a neural network model to learn word associations from a large corpus of a text.

→ Once trained such a model can detect synonymous words or suggest additional words for a partial sentence. As the name implies, word2vec represents each distinct word with a particular list of numbers called a vector.

→ Google trained with 3 billion words.

Vocabulary → unique words in corpus.

Feature Representation

Features	Boy	Girl	King	Queen	Apple	Mango
Gender	-1	1	-0.92	0.93	0.01	0.023
Royal	0.04	0.02	0.95	0.96	0.5	0.68
Age	0.03	0.02	0.75	0.68	0.86	0.97
Food	:	:	:	:	0.98	0.92
...	:	:	:	:	:	:
nth	:	:	:	:	:	:

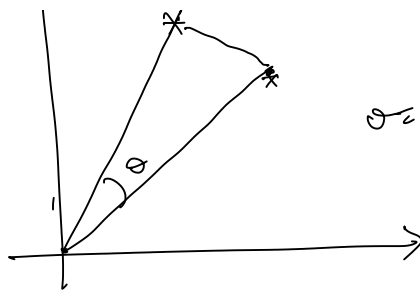
→ Each & every word is represented based on feature representation.

$$[KING - BOY + QUEEN] = GIRL$$

Cosine Similarity:-



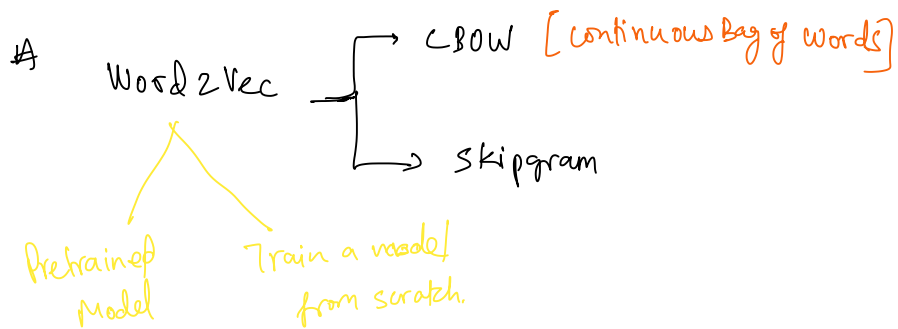
$$Distance = 1 - \text{cosine similarity}$$



$$\begin{aligned} \theta &= 45^\circ \\ 1 - \cos \theta &= 1 - \cos 45^\circ \\ &= 1 - \frac{1}{\sqrt{2}} \\ &= 1 - 0.707 \\ &\approx 0.29 \end{aligned}$$

$$\theta = 90^\circ \quad 1 - \cos 90^\circ = 1 - 0 = 1 \text{ (different)}$$

$$\theta = 0^\circ \quad 1 - \cos 0^\circ = 1 - 1 = 0 \text{ (similar)}$$



## ① CBOW

CORPUS:- Dataset

[iNeuron Company] is related to Data Science

↳ central word based on window.

window size = 5

Z/P      O/P      (Backward context & forward context is considered)

→ iNeuron Company, related to IS

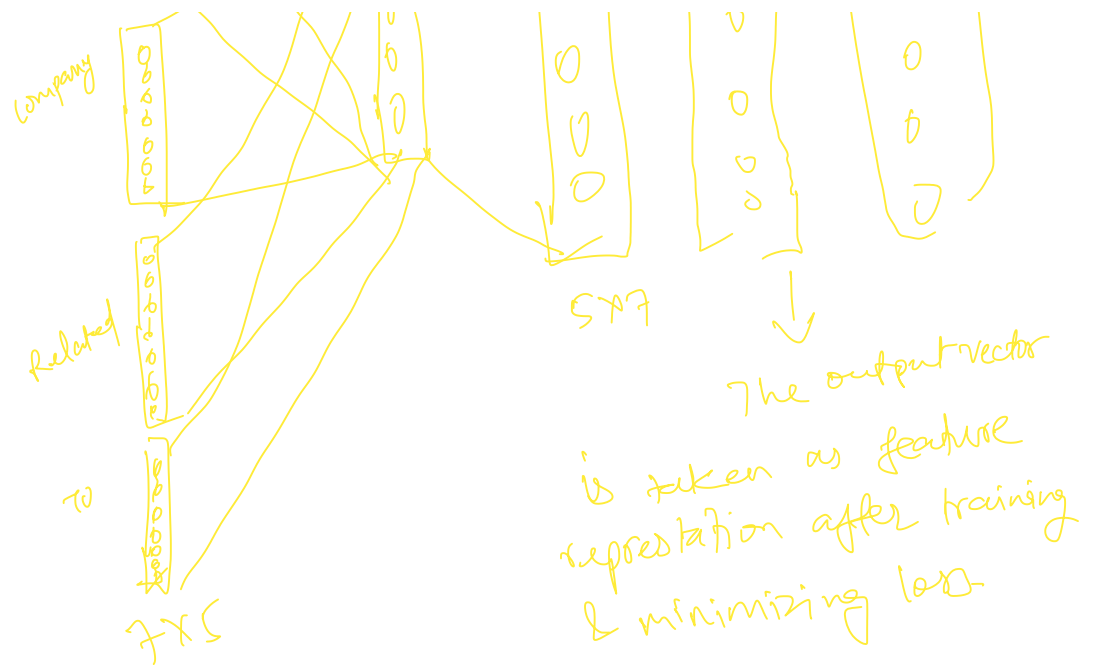
→ Company IS, To Data related

→ IS related, Data Science To

CBOW [Fully Connected Neural Net]

2 layers





## ② Skipgram!

$$\begin{matrix} \text{O/P} & & \text{I/P} \\ \text{[Inuron, Company, Related, To]} & \text{Is} & \\ \text{[Company, Is, Data]} & \text{Related} & \end{matrix}$$

Small Dataset  $\rightarrow$  CBOW

Huge Dataset  $\rightarrow$  Skipgram.

Increase acc. of CBOW or Skipgram.

- ① Increasing training data.
- ② Increase window size  $\rightarrow$  vector dimension increases.

Google word2vec

3 billion words  $\rightarrow$  Google News

Feature representation of 300 dim.

~  $\text{word} = [0. \dots 300^{\text{th}}]$



eg. circum

## # Advantages of Word2Vec.

- ① Dense matrix - earlier we were getting sparse.
- ② Semantic info is getting captured
- ③ Vocab. size  $\rightarrow$  we have fixed set of dimension in word2vec.
- ④ OOV  $\Rightarrow$  every word has some feature representation.