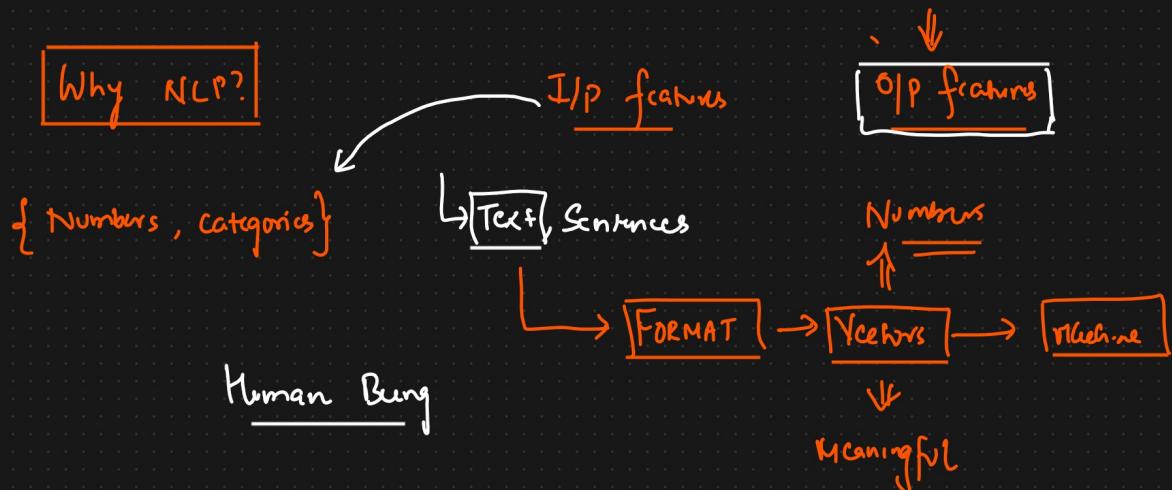


Natural Language Processing Machine Learning

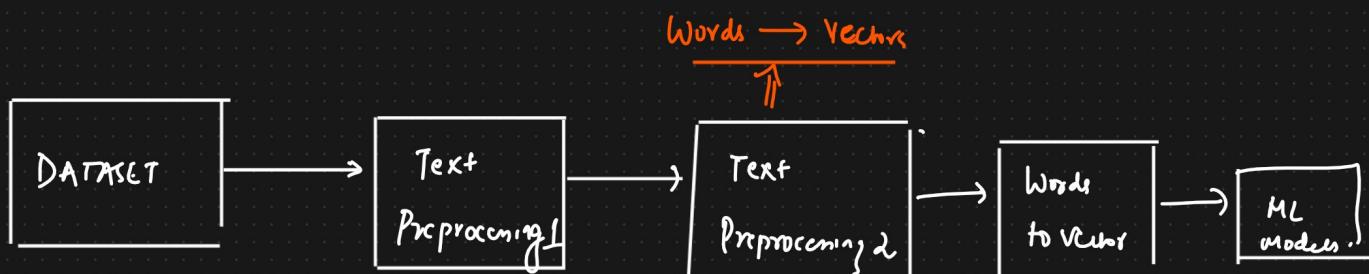


- {
- ① BOW
- ② TF-IDF
- ③ Word2Vec {Deep Learning Technique}

Words → Vectors.

Dataset

Text	O/P
D1	1
D2	0
D3	1
D4	1
D5	0



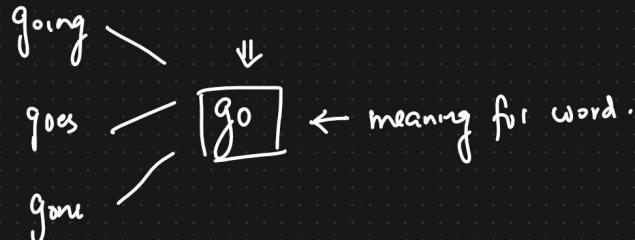
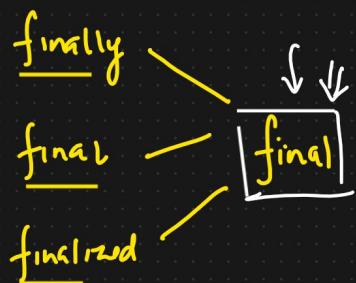
- Cleaning our text data:
- | | |
|--|--------------|
| ① Stopwords
② Stemming ✓
③ Lemmatization ✓
④ Tokenization { | ① BOW |
| | ② TF IDF |
| | ③ N Grams |
| | ④ Word2Vec } |

Stemming And Lemmatization

PROCESS OF REDUCING THE WORD
TO THEIR ROOT FORM



Stemming

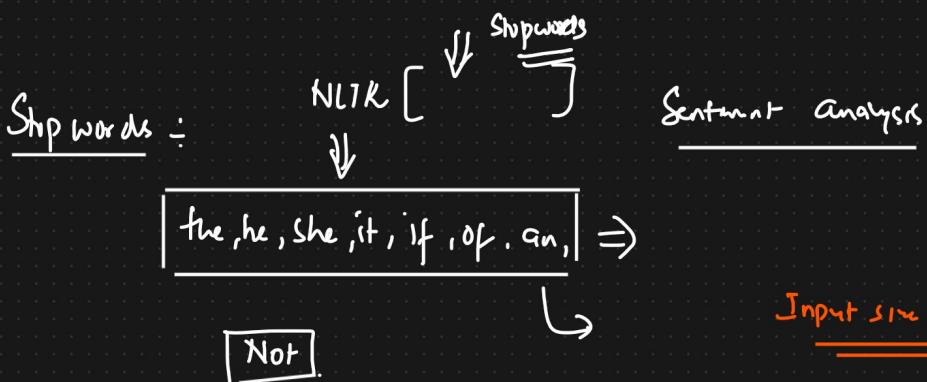


Meaning of the word may change

happy → happi

Lemmatization

Root Form of the Word → Meaningful Word.



① One hot Encoding

- D1 A man eat food]
 D2 [cat eat food]
 D3 [People like DataScience]
 (t+)

ML & DL
 → I/P size fixed

⇒ Train

⇒ Corpus

⇒ Model

$$D3 = \left[\begin{bmatrix} \quad \end{bmatrix} \right]$$

$$\boxed{4 \times 7} \checkmark \left[\begin{bmatrix} \quad \end{bmatrix} \right]$$

$$\Rightarrow \text{stop words} \left[\begin{bmatrix} \quad \end{bmatrix} \right]$$

$$\left[\begin{bmatrix} \quad \end{bmatrix} \right]$$

$\rightarrow D_4 [\text{KRISH} \text{ YT CHANNEL }] \leftarrow \text{Test Data}$

Vocabulary

[man, eat, cat, food, people, like, data science]

$$D_1 \left[\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix} \right] \xrightarrow[3 \times 7]{} \begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \Rightarrow \text{One Hot Encoding} \quad \checkmark$$

Advantages

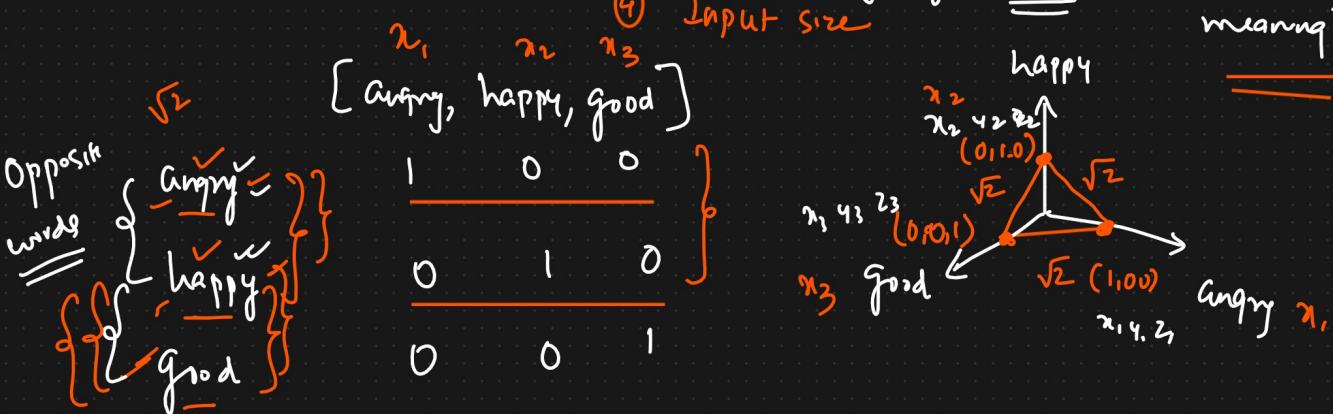
① Simple to Implement

④ Sparse Matrix {Overfitting} ✓

② OOV {out of vocabulary} ✓

{ ③ No capturing of semantic? {semantic meaning} }

④ Input size



$(0,1,0) \xrightarrow{\sqrt{2}} (1,0,0)$ $\{ \boxed{\text{No}} \}$ $(0,0,1) \xrightarrow{\sqrt{2}} (1,0,0)$

$$\sqrt{(x_1 - z_1)^2 + (y_1 - z_2)^2 + (z_1 - z_3)^2}$$

$$\sqrt{(1-0)^2 + (0-0)^2 + (0-1)^2}$$

$$= \sqrt{1+1} = \sqrt{2}$$

If size \rightarrow first

\Rightarrow

Test data \rightarrow [good : ~~KRISH~~ \boxed{YX}] \rightarrow New words

② Bow of Bag of Words } [Sentiment Analysis, Text Classification].

don't → Stopwards

D1 → He is a good boy ① Lower the words D1 → good boy boy
D2 → She is a good girl ⇒ ② Stopwords D2 → good girl
D3 → Boy and girl are good ③ Stemming D3 → boy girl good

$\{ \text{Binary} = \text{True} \}$

Vocabulary Frequency

Bow Vocabulary

Vocabulary

०/प

good	3
boy	2
girl	2

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

I/f \Rightarrow fixed

Advantages

Disadvantages

- ① Simple & Intuition
- ② Input fixed issue is Resolved

- ① Sparsity {Exists} It is low when Compa
↓ Reduced to one
- ② OOV → Exist → We should ignore the new words.
- ③ Semantic Relationship.

{ I like pizza }
{ I don't like pizza } ←

Captured Not Captured

$$\left\{ \begin{bmatrix} 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \right\}$$

Vocab

→ ORDERING

N-Grams

pizza 2
don't 1

2

Steve (1, 2)
 intelligent boy like pizza don't like pizza
 " " " " " "
 smart D1 I I O I O O O O O
 " " D2 I I I I O O O I O

Kenny

- ① Unigram \Rightarrow Bag of words
- ② Bigram \Rightarrow Bag of words + Bigram.
- ③ Tri Gram
- ④ Quad Gram
- ⋮
- n gram

(1,2)

ORDERING
↑

KRISH EATS FOOD

[FOOD EATS]

KRISH EATS FOOD

[KRISH EATS]

[EATS FOOD]

[KRISH FOOD]

[FOOD KRISH]

③ Term Frequency \rightarrow Inverse Document Frequency (TF-IDF)

Sent1: good boy

Sent2: good girl

Sent3: boy girl good

Term Frequency = No. of rep. of words in sentence

No. of words in sentence

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

Sent = Document

Term Frequency

*

Inverse Document Frequency
IDF

	Sent1	Sent2	Sent3	words	
good	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	good	$\log_e \left(\frac{3}{3} \right) = 0$
boy	$\frac{1}{2}$	0	$\frac{1}{3}$	boy	$\log_e \left(\frac{3}{2} \right)$
girl	0	$\frac{1}{2}$	$\frac{1}{3}$	girl	$\log_e \left(\frac{3}{2} \right)$

TF-IDF

	f_1	f_2	f_3	
	good	boy	girl	
S_{ent1}	0	$\frac{1}{2} \times \log_e(3/2)$	0	$\begin{bmatrix} good & boy \end{bmatrix}$
S_{ent2}	0	0	$\frac{1}{2} \log_e(3/2)$	$\begin{bmatrix} good & girl \end{bmatrix}$
S_{ent3}	0	$\frac{1}{3} \log_e(3/2)$	$\frac{1}{3} \log_e(3/2)$	$\begin{bmatrix} good & boy & girl \end{bmatrix}$

① Advantages

- ① Semantic meaning is well captured to some extent
- ② Intuitive.

① Disadvantage

- ① Sparsity
- ② Out of vocabulary
- ③ Ordering (Ngrams).
- ④ Computation is high.
↳ Complexity.