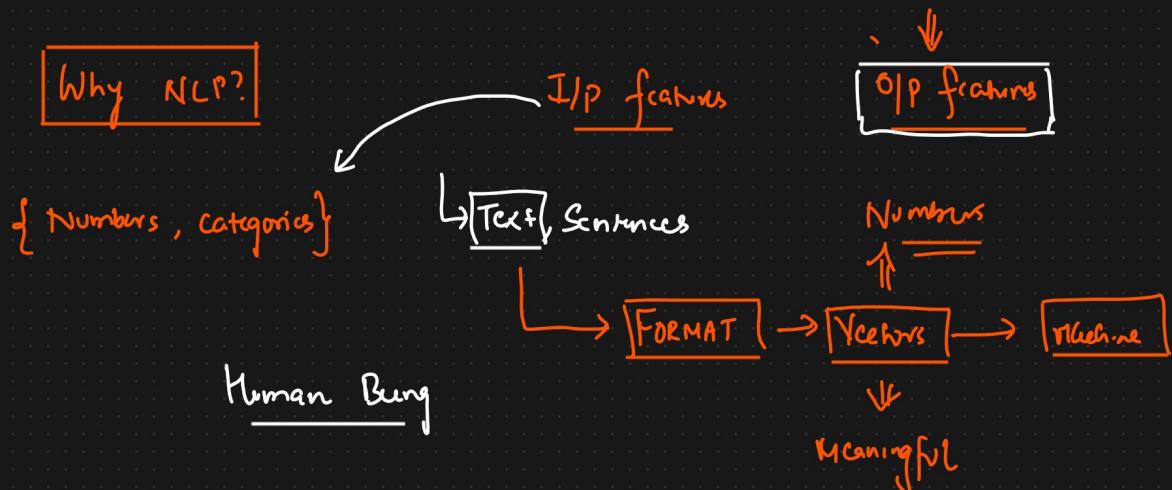


# Natural Language Processing Machine Learning

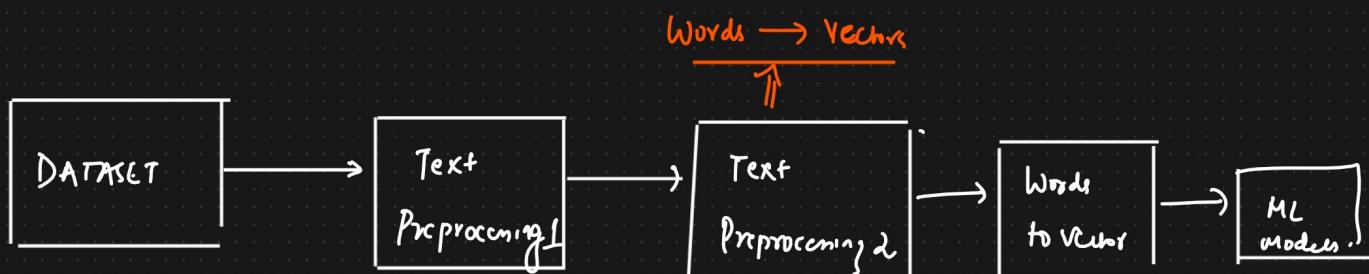


- {
- ① BOW
- ② TF-IDF
- ③ Word2Vec {Deep Learning Technique}

Words → Vectors.

Dataset

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



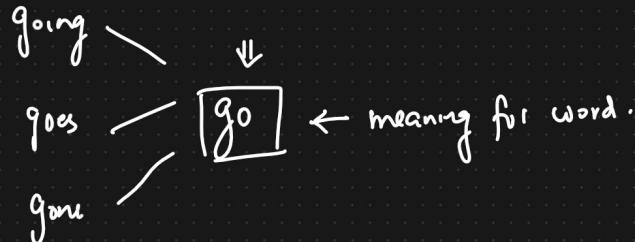
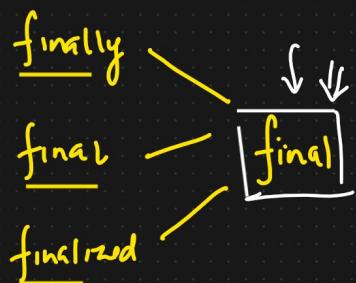
- Cleaning our text data:
- |                                                                      |                                                |
|----------------------------------------------------------------------|------------------------------------------------|
| ① Stopwords<br>② Stemming ✓<br>③ Lemmatization ✓<br>④ Tokenization { | ① Bow<br>② TF IDF<br>③ N Grams<br>④ Word2Vec { |
|----------------------------------------------------------------------|------------------------------------------------|

## Stemming And Lemmatization

PROCESS OF REDUCING THE WORD  
TO THEIR ROOT FORM



Stemming

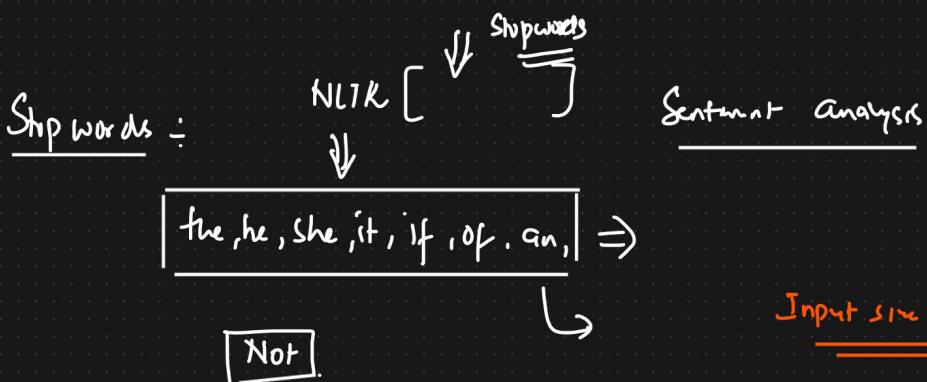


Meaning of the word may change

happy → happi

## Lemmatization

Root Form of the Word → Meaningful Word.



## ① One hot Encoding

D1 [ A | man | eat food ] ⇒ Train

D2 [ eat eat food ] ⇒ Corpus

D3 [ People like Datasource ] (tt) ⇒ Model

ML & DL  
 → I/P size fixed

$D_3 = \begin{bmatrix} [ \quad ] \\ [ \quad ] \\ [ \quad ] \\ [ \quad ] \end{bmatrix}$

$4 \times 7$  ✓ [        ]

[        ]

[        ]

[        ]

→ Du [KRISN YT CHANNEL] ← test Data

man, eat, cat, food, people, like, data Science

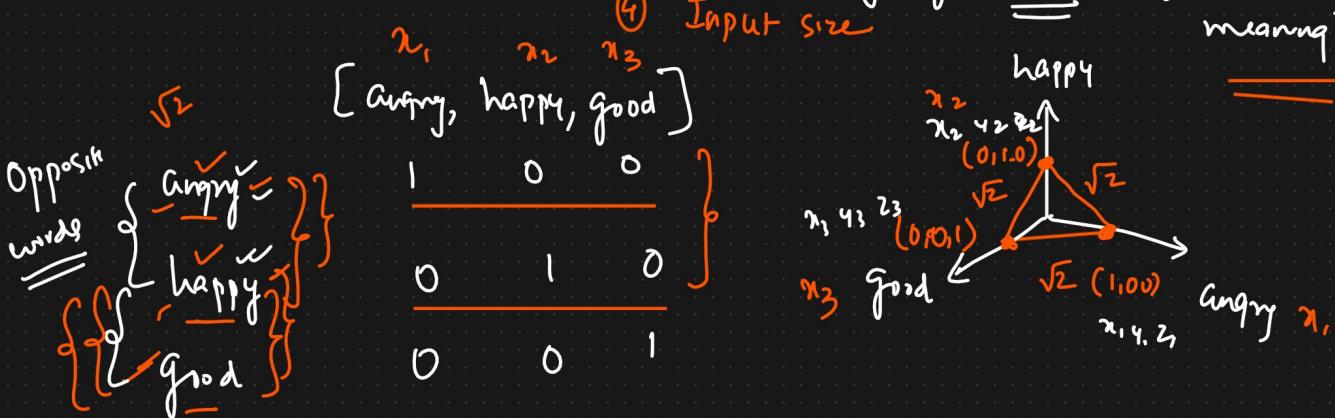
$$D1 \begin{bmatrix} [1 & 0 & 0 & 0 & 0 & 0 & 0] \\ [0 & 1 & 0 & 0 & 0 & 0 & 0] \\ [0 & 0 & 0 & 1 & 0 & 0 & 0] \end{bmatrix} \begin{bmatrix} [0 & 0 & 0 & 0 & 0 & 0 & 0] \\ [0 & 0 & 0 & 0 & 0 & 0 & 0] \end{bmatrix}_{\boxed{3 \times 7}} \Rightarrow \underline{\text{One Hot Encoding}}.$$

## Advantages

## Disadvantages

- ① Simple to Implement      ④ Spark Marin {Overfitting} ✓

- ② OOV {out of vocabulary} ✓
  - ③ No capturing of Semantic? {semantic meaning?}
  - ④ Input size =



$$\sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2 + (z_1 - z_2)^2}$$

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

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

Slip size → first

11

Up size  $\rightarrow$  first  
= Test data  $\rightarrow$  [Good : Karen ✓]  $\rightarrow$  New words

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

don't → Stopwords

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

{ Binary = True }

Vocabulary

Frequency

Bow

Vocabulary

~~~~~  
O/P

{ good 3 }  
 { boy 2 }  
 { girl 2 }

[ good boy girl ]

|    |   |   |   |   |
|----|---|---|---|---|
| D1 | 1 | 1 | 0 | - |
| D2 | 1 | 0 | 1 | - |
| D3 | 1 | 1 | 1 | - |

1/D ⇒ fixed

Advantages

Disadvantages

- ① Simple & Intuitive  
 ② Input fixed issue is resolved

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

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

Captured Not Captured

{ [ 1 0 1 1 ] }  
 { [ 1 1 1 1 ] }

Vocab

like      freq

2

⇒ ORDERING

④ N-Grams

pizza  
 don't

1 (1,2) ⇒ Bigram along with unigram

don't pizza

He is intelligent like pizza don't like don't pizza don't don't like pizzas like  
 This is not a { D1 | | | 0 1 }      D2 | | | 1 1 1 0 0 0 0 0 0 0 0 0

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.