# NLP
## Natural Language Processing

① Roadmap

why NLP?
→ Where text is data, we need the machine
to understand

Roadmap

→ ① Text Preprocessing → Bag of words, Tokenization, stemming, TF-IDF, Word to Vec, stopwords, Lemitization

② Text Preprocessing2- Unigrams, Bigrams

③ Text prepro cy2 - Avg Word2Vec, gensim

③ ML use cases - Chat Bots, etc.

④ RNN, LSTM RNN, GRU RNN

⑤ Advanced Preprocessing - Word embedding

⑥ Bidirectional LSTM, Encoders, Decoders, Attention model

⑦ Transformers Bert.

Library

NLTK
SPacy    } ML
Text Blob

Tensorflow } D L
Pytorch
Hugging face

① Tokenization

    Connecting sentence into words

    ② stopwords :- to, of, a, an, the, etc. :

    ③ stemming :- Process of Reducing words to their bare word stem

       Disadvantage
       → The bare word may not have any meaning

       eg:- historical → histori       finally    }
              history             final     } → fi·al
                                  finalized

    Advantages
    → It is really fast

    ④ Lemmatisation —→ This has dictionary of words

       history
              → history
       historical

                        Disadvantage
    Advantage                    ① It is slow
    ① Get a meaningfull word

    Usecase                 Lemmtion
      stemming             ① Text generation
      ↳ spam classification    ② Language translation
      ↳ Reviews classification   ③ chatbot

    step 2 :- Words to Vectors

    ① Bag of words
    ② TF-IDF

## Terms

1) Corpus → $D_1, D_2, D_3, D_4$   content → Para

2) Documents → suture

3) Vocabulary → unique words

4) words

| | example | Text | O/P |
|---|---|---|---|
| | D1 | The food is good | 1 |
| | D2 | food is bad | 0 |
| | D3 | Pizza is nyce | 1 |
| | D4 | bryni is bad | 0 |

Dataset → Text Pupuais 1 → Text Pupuaiig 2 → words into verbs

Text Pupuais 1:
lowering ↓
Tokenization / lowering
lemahzation
Steming

words into verbs:
1) BOW
2) TFIDF
3) Word2vec

---

## 1) One hot Encoding

Corpus :- A man eat food
→ Cat eat food
→ People watch KushTube YT

Vocabulary
A man eat Cat PPl
watch Kush YT

$D_1 \rightarrow$
$$\begin{bmatrix} 1 & 0 & 0 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 0 & 1 & 0 & 0 \end{bmatrix} \quad \text{features}$$
$$\begin{bmatrix} 0 & 0 & 1 & 0 \end{bmatrix}$$
$$\begin{bmatrix} 0 & 0 & 0 & 1 \end{bmatrix}$$

D1 $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$
$\begin{bmatrix} 0 & 1 & 0 & 0 & - & - & - \end{bmatrix}$
$\begin{bmatrix} 0 & 0 & 1 & 0 & . & . & . \end{bmatrix}$

---

### Issues

1) Spause Matrix

2) features will change on size on sentence

3) Semantic meaning within the word is not capture

4) Out of vocabulary

### Advantages

1) simple to Implement

2) Intutive

② **Bag of words**

D1 :- He is a good boy
D2 :- She is a good girl
D3 :- Boys and girls are good.

↳ Applying stopwords

D1 :- good boy
D2 :- good girl
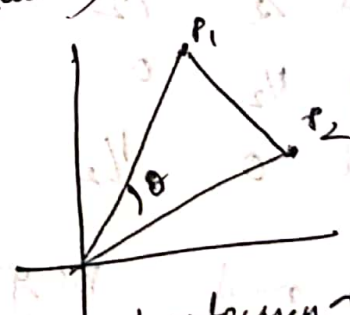D3 :- Boys girls good

| Vocabulary | Frequency |
|---|---|
| good | 3 |
| boy | 2 |
| girl | 2 |

→ order will be based on frequency of o/p

|  | d1 | d2 | d3 |
|---|---|---|---|
|  | good | boy | girl |
| D1 | 1 | 1 | 0 |
| D2 | 1 | 0 | 1 |
| D3 | 1 | 1 | 1 |

**Advantages and disadvantages**

↳ simple and intuitive

→ Sparcity

→ Out of Vocabulary

↳ ordeing of the words has completely changed

**Cosine Similarity** (Application of cosine Rule)

| 1 − Cos θ = Cosine Similarity |

In order to capture semantic info

→ to follow the same order

→ N grams :- Bigrams, Trigrams

→ In count vectorizer → byprocedure

good(d1) boy(d2) girl(d3)   d4 good boy   d5 good boy

|  | good(d1) | boy(d2) | girl(d3) | d4 good boy | d5 good boy |
|---|---|---|---|---|---|
| D1 :- | 1 | 1 | 0 | 1 | 1 |
| D2 :- | 1 | 0 | 1 | 0 | 0 |
| D3 | 1 | 1 | 1 | 0 | 0 |

Anish Eats food

Bigrams: Anish Eats
:- Eats food

## ③ TF-IDF

Term frequency and Inverse Document frequency

D1 :- good boy

D2 :- good girl

D3 :- boy girl good

> { More weightage to rare words }

$$\text{Term Frequency} = \frac{\text{no. of repetitions of words in sentence}}{\text{no. of words in sentence}}$$

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

$$\boxed{TF\text{-}IDF = TF \times IDF}$$

ex:-

| TF → For sentence | | | |
|---|---|---|---|
| | D1 | D2 | D3 |
| good | ½ | ½ | ⅓ |
| boy | ½ | 0 | ⅓ |
| girl | 0 | ½ | ⅓ |

IDF → for words

good  $\ln\left(\frac{3}{3}\right) \Rightarrow 0$

boy  $\ln\left(\frac{3}{2}\right)$

girl  $\ln\left(\frac{3}{2}\right)$

good ——— boy ——— girl

D

| | good | boy | girl |
|---|---|---|---|
| D1 | 0 | $\frac{1}{2}\ln\frac{3}{2}$ | 0 |
| D2 | 0 | 0 | $\frac{1}{2}\ln\frac{3}{2}$ |
| D3 | 0 | $\frac{1}{3}\ln\left(\frac{3}{2}\right)$ | $\frac{1}{3}\ln\left(\frac{3}{2}\right)$ |

Advantages
1. Intuitive
2. Word Importance is getting
   capture

Disadvantage
1. Sparcity
2. Out of vocabulary

2. Word 2 vec
   ( Feature Representation :
   - limited dimensions
   2. sparcity is reduced
   3. Semantic meaning is maintained

1. Word Embedding → CBOW, skipgram
2. Word 2 vec

Word Embedding
├── Frequency
│   1. One hot Encody ✓
│   2. Bag of words ✓
│   3. TF-IDF ✓
└── Deep learning trained model
    1. word 2 vec
       ├── CBOW (Continuous BOW)
       └── skipgrams

|       | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ |
|-------|-------|-------|-------|-------|-------|-------|
|       | Boy   | girl  | king  | queen | Apple | mango |
| center|       |       | -0,92 | +0,93 | 0     | 0.4   |
| male  | -1    | 1     |       |       |       |       |
| Royal | 0.01  | 0.02  | 0.95  | 0.96  | -0.02 |       |

* Word 2 vec
1. CBOW (continuous Bag of words)

CORPUS :- "KRISH CHANNEL IS RELATED TO DATA science
WINDOW SIZE :- 5

center word is O/p, target &
context word

TRAINING DATA

| Independent feature | O/P |
|---------------------|-----|
| 1. Krish, CHANNEL RELATED | IS |
|           to        | TO  |
| 2. CHANNEL, IS, TO, DATA | to |
| 3. IS RELATED DATA science | to |

Bow. (Hence this is called continuon Bow

Kush    1 1 0 0 0 0 0 0

chaul  0    1 0  -

IS         0 0 1

RIGHT Bo  0 0 0 1 0 ...

ANN                    window size 25              softmax
Word 1                                             Layer
                                                   bush
      → Neurons                                    [          ]
                                                        Vector

Word                                               

Loss = $(y - \tilde{y})$

skip gram
        I/P           O/P
        IS            other
        related       column
        to