# Day 2 - NLP For ML And DL

## Agenda

{NLTK}

① Text preprocessing → Words → Vectors

    a) OHE (One hot Encoding)
    b) Bag of Words (BOW)    ⇒ practical Implementation
                                         ⇒ ngrams
    c) TF-IDF (Term Frequency - Inverse Document Frequency)

    d) Word2Vec

② Quiz → Live → 5000Rs INR

             1st Prize → 2000 Rs INR
             2nd Prize → 1500 Rs INR
             3rd Prize → 1500 Rs INR

## Basic Terminologies Used In NLP

Sentiment Analysis

① CORPUS ✓→ Paragraph → [D1, D2, D3, D4]

② Documents ✓→ Sentence     Text     O/P     → Dictionary Book

③ Vocabulary ✓⇒ 10K unique words

               → D1   The food is good    1     10K unique words

④ Word  ⇒ word
             → D2   The food is bad     0

             → D3   Pizza is amazing   1

             → D4   Burger is bad    0

                   ↑
              DATASET



DATASET → Text -1 Preprocessing → Text Preprocessing -2 → Words → vectors

Tokenization,        ① STEMMING    ① BOW ?

## ① One hot Encoding

Paragraph

Yocabulary

```
    ┌→ A man eat food ⎫
    │    =   =   =   =  ⎬ ⇒ CORPUS
    →  Cat Eat food  by ⎭→ size
    →  People Watch KRISH YT]
       =      =     =   =
```

9
⇓
food ⤴

A    man    eat    Cat    People    Watch    KRISH   YT

Out of vocabulary

⇒ CANNOT TRAIN ⎫
     THE MODEL  ⎬

$$D1 \rightarrow \begin{bmatrix} [1 & 0 & 0 & 0], \\ [0 & 1 & 0 & 0], \\ [0 & 0 & 1 & 0], \\ [0 & 0 & 0 & 1] \end{bmatrix}$$

$$D2 - \begin{bmatrix} [1 & 0 & 0] \\ [0 & 1 & 0], \\ [0 & 0 & 1] \end{bmatrix}$$

### Advantages

① Simple to Implement ⎫
② Intuitive        ⎬

### Disadvantage

① Sparse Matrix ✓    [Extra Test data]

② OOV {out of Vocabulary} ✓

③ Not fixed size ✓

④ Semantic meaning between

word is not captured

## ② Bag of Words

→ Stop words
lower all the words
case

D1 → ⊙He ⊙is ⊙a good boy    D1 → good   boy   good

D2 → ⊙She ⊙is ⊙a good girl ✗   D2 → good   girl

D3 → Boy ⊙nd girl ⊙re good   D3 → Boy    girl   good

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

Bow => Binary Bow

|  | $f_1$ good | $f_2$ boy | $f_3$ girl |
|---|---|---|---|
| Doc 1 → | 1 | 1 | 0 |
| Doc 2 | 1 | 0 | 1 |
| Doc 3 → | 1 | 1 | 1 |

O/p => Assumptions



{ Euclidian Distance
  Cosine Similarity .

## Advantages

① Simple and Intuitve

Cosine Similarity



$P_1$
(0,1)
90°
Similar
$P_2$
(1,0)



Minions

Avengers

IRON

Caphre the semantic Info

## Disadvantages

① Sparsity

② OOV

③ Ordering of the words.

④ Semantic meaning Not able to
Cos 45 = 0.53  } Capture .

1 - 0.53 = Cos - similarity

0.47

Cos 90 = 0

1 - 0 = [1]

Cos 0 = 1

1 - 1 = 0

# Ngrams $\Rightarrow$ Bigrams, Trigrams, — N grams

| | f1 good | f2 boy | f3 girl | f4 good boy | f5 good girl |
|---|---|---|---|---|---|
| Sent 1 | 1 | 1 | 0 | 1 | 0 |
| Sent 2 | 1 | 0 | 1 | 0 | 1 |
| Sent 3 | 1 | 1 | 1 | 0 | 0 |

KRISH   EATS   FOOD

practically

BI-GRAMS? 2 Bigrams

KRISH EATS        EATS FOOD

TRIGRAMS $\Rightarrow$ 3 Trigrams

I am not feeling Well

I am not        Am Not feeling        Not feeling well

ngrams

KRISH IS NOT FEELING WELL        (1,3)

| f1 | f2 | f3 | f4 | F5 | KRISH IS |
|---|---|---|---|---|---|

| KRISH | IS | NOT | FEELING | WELL |
|---|---|---|---|---|