

# PHASE - 2

## TEXT - TO - VECTORS

Text → Vectors means converting words or sentences into numerical form so, ML / DL understands them.

### ONE HOT ENCODING :-

One hot Encoding is a technique to convert categorical text data into numerical vectors for ML models.

What it does :-

→ Each category becomes a binary vector.  
(0s and 1s)

Example :-

Categories :-

['Red', 'Blue', 'Green']

One-hot vectors :-

Red → [1, 0, 0]

Blue → [0, 1, 0]

Green → [0, 0, 1]

**Note :-** In NLP, we use Unique Vocabulary.

Unique Vocabulary means the set of distinct words present in a text or corpus.

Example :- Text

D1 The food is good

D2 The food is bad

{ Give, Text, D1 and D2, are the reviews of customers here we use unique vocabulary to find vectors for our NLP model. }

So, we will create a unique Vocabulary.

	good	bad			
1	0	0	0	0	0
0	1	0	0	0	0
0	0	1	0	0	0
0	0	0	1	0	0
0	0	0	0	0	1

which means,

The  $\rightarrow [1, 0, 0, 0, 0]$

food  $\rightarrow [0, 1, 0, 0, 0]$

is  $\rightarrow [0, 0, 1, 0, 0]$

good  $\rightarrow [0, 0, 0, 1, 0]$

bad  $\rightarrow [0, 0, 0, 0, 1]$

For D1, The food is good

Our vectors are,

D1  $[ [1, 0, 0, 0, 0]$   
 $[0, 1, 0, 0, 0]$   
 $[0, 0, 1, 0, 0]$   
 $[0, 0, 0, 1, 0] ]$

For D2, The food is bad

Our vectors are,

D2  $[ [1, 0, 0, 0, 0]$   
 $[0, 1, 0, 0, 0]$   
 $[0, 0, 1, 0, 0]$   
 $[0, 0, 0, 0, 1] ]$

## ONE-HOT ENCODING

### Advantage

- Easy to implement with Python  
[Sklearn, OnehotEncoder, pd.get\_dummies]

### Disadvantages

- Sparse matrix - Overfitting (matrix, which most values 0)  
 $(1, 0, 0, 0, 0)$
- Ml Aglo. → Fixed size I/P  
→ all input must be same size
- is getting captured  
X does not know the word's meaning,
- X different words can have same meaning
- Out of vocabulary  
(if new I/P is given its not in vocabulary)

# BAG OF WORDS :- (Bow)

Bag of Words (Bow) is a simple text  $\rightarrow$  vector technique that represents text by word frequency, ignoring grammar and word order.

## # How Bow works :-

- $\Rightarrow$  Collect Unique vocabulary
- $\Rightarrow$  Collect word occurrences
- $\Rightarrow$  Create fixed length vector.

## Example :-

### Step 1

#### Text

HP is a good boy lower all the words S1  $\rightarrow$  good boy  
She is a good girl  $\Rightarrow$  S2  $\rightarrow$  good girl  
Boy and girl are good Stopwords S3  $\rightarrow$  Boy girl good

#### Note:-

Here, we are applying .lower() and after removing Stopwords we are getting new S1, S2 and S3

What is { .lower()  $\rightarrow$  Boy  $\rightarrow$  boy }  
happening, { Stopwords  $\rightarrow$  is, a, and got removed }

Step-2  $\Rightarrow$

Now we will be counting frequency of unique words. {Vocabulary}

Vocabulary	Frequency	
good	3	max
boy	2	
girl	2	min

Note  $\rightarrow$  Frequency of unique words {vocabulary} is always in descending order, so we can also say that, the word which came maximum times appeared on top.

Step 3  $\Rightarrow$

Vocabulary	Frequency	good	boy	girl
good	3	S1 [1	1	0]
boy	2	$\Rightarrow$ S2 [1	0	1]
girl	2	S3 [1	1	1]

Creating Vector of  $S_1, S_2$ , and  $S_3$  using vocabulary or Unique vocabulary

## Binary Bow

→ { 1s and 0s }

## Bow

{ count will get updated based on frequency }

→ Example :-

"I love love Python"

created

Vocabulary

[ "I" , "Love" , "Python" ]

→ Example :-

"I love love Python"

vocabulary

[ 'I' , 'Love' , 'Python' ]

Vector will Be

[ 1, 1, 1 ]

Vector will be

[ 1, 2, 1 ]

'2' because of count of  
'love'

## BAG OF WORD (BOW)

### Advantages

- Simple and Intuitive
- Fixed Size Input doesn't matter because it can handle any size of sentence.

### Disadvantages

- Sparse matrix or array
- Ordering of word is getting changed
- Out of vocabulary (OOV)
- Semantic meaning is still not captured.

# TF-IDF [TERM FREQUENCY - INVERSE DOCUMENT FREQUENCY]

∅ TF-IDF tells which words are important in a document.

Taking previous Example

S<sub>1</sub> → good boy  
S<sub>2</sub> → good girl  
S<sub>3</sub> → ~~good~~ boy girl good

Term Frequency :-

How many time a word appears in one document.

Term Frequency (TF) =  $\frac{\text{No. of rep of words in sentence}}{\text{No. of words in sentence}}$

So, we are calculating TF for S<sub>1</sub>, S<sub>2</sub> and S<sub>3</sub>

	S <sub>1</sub>	S <sub>2</sub>	S <sub>3</sub>
good	1/2	1/2	1/3
boy	1/2	0	1/3
girl	0	1/2	1/3

## IDF (Inverse Document frequency) :-

IDF tells how important a word is, across all document

⇒ Rare word = More Important

⇒ Common word = Less Important

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

Now, calculating IDF for S1, S2 and S3

<u>Words</u>	<u>IDF</u>	(S1, S2, S3) no. of sentence <u>good</u> is present in all sentences
good	$\log_e (3/3) = 0 \#$	
boy	$\log_e (3/2)$	
girl	$\log_e (3/2)$	

\* Note :- To find TF-IDF we basically multiply the

$$TF-IDF = TF * IDF$$

### Term Frequency (TF)

	S1	S2	S3		IDF words	IDF
good	1/2	1/2	1/3	x	good	$\log_e(3/3) = 0$
boy	1/2	0	1/3	x	boy	$\log_e(3/2)$
girl	0	1/2	1/3	x	girl	$\log_e(3/2)$

### Final TF-IDF

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

	good	boy	girl
Sent 1	0	$1/2 * \log_e(3/2)$	0
Sent 2	0	0	$1/2 * \log_e(3/2)$
Sent 3	0	$1/3 * \log_e(3/2)$	$1/3 * \log_e(3/2)$

Here, Multiplying Sentence 1 ( $s_1$ )'s TF with IDF value?

$$\begin{matrix} \text{good} & \text{good} \\ \text{boy} & \text{boy} \\ \text{girl} & \text{girl} \end{matrix} *$$

which means

$$\Rightarrow s_1(\text{good}) * \text{IDF(good)}$$

$$\Rightarrow \frac{1}{2} * 0$$

$$\Rightarrow \Rightarrow 0$$

### TF-IDF

#### Advantages

→ Intuitive

→ Fixed size → Vocab size

→ Word Imp is getting captured.

#### Disadvantages

→ Sparsity still exists

→ Out of vocabulary (OOV)