

11) One Hot encoding:-

	Text	O/P
D1	The food is good	1
D2	The food is bad	0
D3	The Pizza is Amazing	1

12) → Vocabulary:-

The, food, is, good, bad, Pizza, Amazing

<3> D1 $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$, vector representation for D1.
 $\begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$,
 $\begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix}$,
 $\begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$

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

D3 $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 \end{bmatrix}$,
 $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$,
 $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$,
 ~~$\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$~~ D3

OHE advantages :-

- Easy to implement.

OHE Disadvantages :-

- Sparse matrix. → Overfitting.
- For ML Algorithm fix size of ZIP but we don't get it.
- No semantic meaning is getting capture.
- out of vocabulary. (~~out~~ OOV)

12) Bags of word

Dataset

Text	O/P
He is a good boy	1
She is a good girl	1
Boy & girl are good	1

Step 1 → lower the text

Step 2 → remove stopwords

T1			T2			
	Vocabulary	freq		good	boy	girl
	good	3	S1	1	1	0
	boy	2	S2	1	0	1
	girl	2	S3	1	1	1

- 1) Binary Bag of words contain only 1 & 0
- 2) normal Bag of words contain counts the freq.

Date
Page

Advantages of BOW

- Simple & intuitive
- fixed size $n \times p$ for ML Algo

Disadvantage of BOW

- Sparse matrix problem NLP \rightarrow Overfitting
- ordering of words getting changed.
- out of vocabulary (OOV)
- Semantic meaning is still not catching.

s_1 The food is good

s_2 The food is not good.

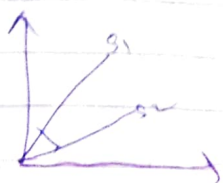
[vectors

s_1 [The food is ~~not~~ good not]

s_2 [1 1 1 1 1]

→ Semantic meaning or cosine similarity show near

→ Using PCA we convert this in 2D plot



← showing near to each other.

13) n-grams

$S_1 \rightarrow$ The food is good

$S_2 \rightarrow$ The food is not good

Step 1: vocabulary & remove stopwords

	food	not	good
$S_1 \rightarrow$	1	0	1
$S_2 \rightarrow$	1	1	1

Step 2: Bigram using

	food	not	good	foodgood	foodnot	notgood
S_1	1	0	1	1	0	0
S_2	1	1	1	0	1	1

SKlearn \rightarrow n-gram $\rightarrow (1,1) \rightarrow$ unigram

$(1,2) \rightarrow$ uni, bigram, ~~trigram~~

$(1,3) \rightarrow$ uni, bi, tri

$(2,3) \rightarrow$ bi, trigram

b) TF-IDF (Term freqⁿ - inverse term freqⁿ)

S₁ → good boy
S₂ → good girl
S₃ → boy girl good

$$\text{Term freq}^n = \frac{\text{no. of rep. of words in sentence}}{\text{no. of words in sentence}}$$

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

Step 1

	Term freq ⁿ				IDF		
	S ₁	S ₂	S ₃		S₁	S₂	S₃
good	1/2	1/2	1/3	good	$\log_e (3/3) = 0$		
boy	1/2	0	1/3	boy	$\log_e (3/2) =$		
girl	1/2	1/2	1/3	girl	$\log_e (3/2) =$		

Step 2

Final TF-IDF

	good	boy	girl	girl
S ₁	0	$\frac{1}{2} \times (\log_e (3/2))$	0	
S ₂	0	0	$\frac{1}{2} \times \log_e (3/2)$	
S ₃	0	$\frac{1}{3} \times \log_e (3/2)$	$\frac{1}{3} \times \log_e (3/2)$	

Note: If a word present in every document the importance is less.

Date _____
Page _____

Advantages of TF-IDF

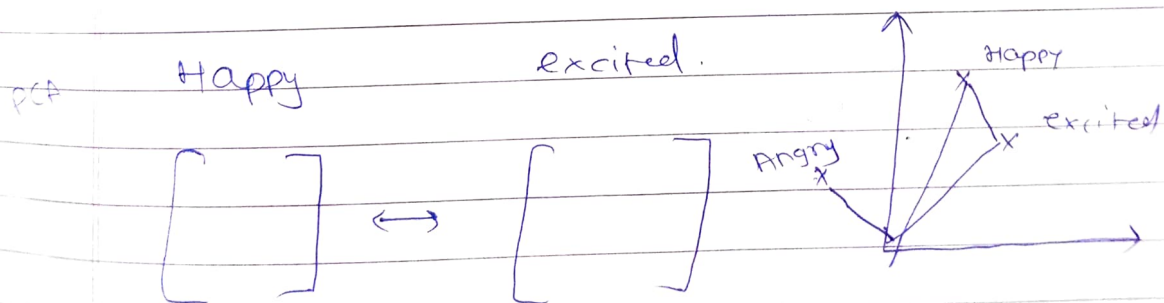
- Simple implementation
- Fixed size \rightarrow vocab size
- Word importance is captured.

disadvantages

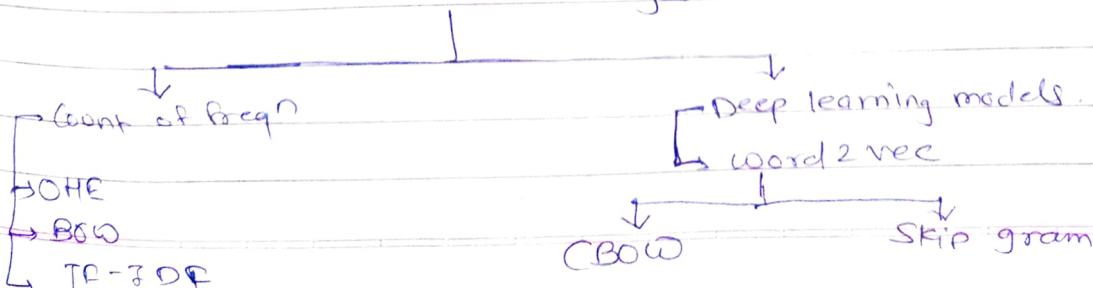
- sparse still exist - overfitting
- OOV (out of vocabulary).

15) Word embeddings

Word embeddings are like 'numbered maps' that represent the meaning of words. They turn words into simple number so that computers can understand and work with them.



Word embeddings



alt
Note:- If a word occurs in every document its importance is less.

Page
Page

Advantages of TF-IDF

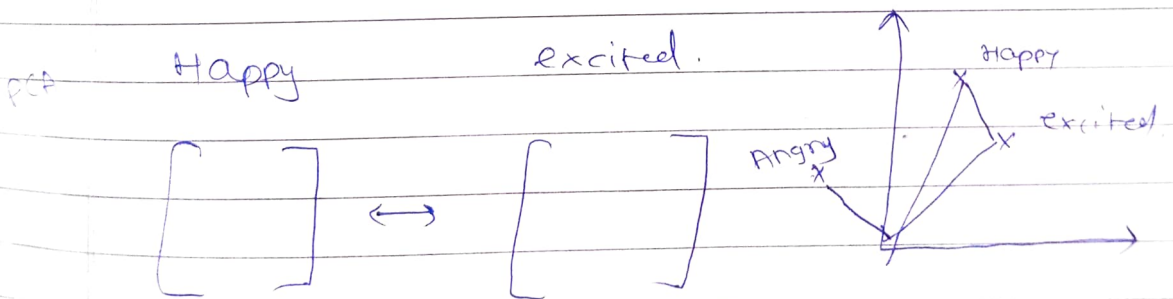
- Simple implementation
- Fixed size \rightarrow vocab size
- Word importance is captured.

disadvantages

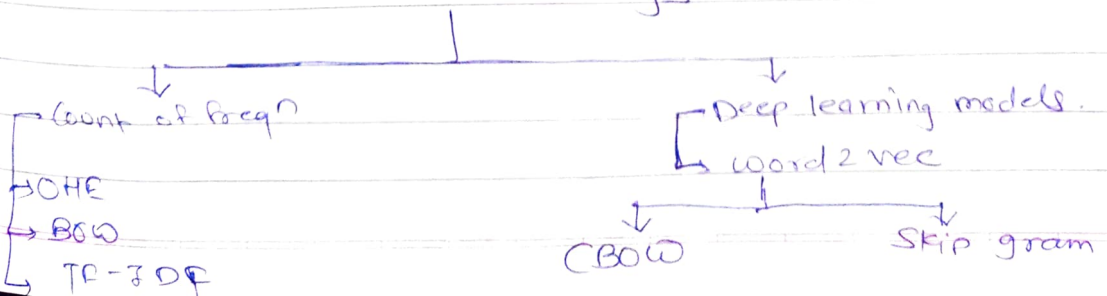
- Sparse still exist - overfitting
- OOV (out of vocabulary).

15) Word embeddings

Word embeddings are like "numbered maps" that represent the meaning of words. They turn words into simple number so that computers can understand and work with them.



Word embeddings

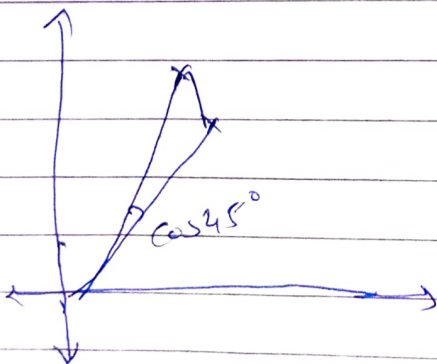


word holonymy \rightarrow unique words \rightarrow corpus

feature representation	words				
\downarrow	Boy	girl	king	queen	Apple
Gender	-1	1	-0.93	0.94	0.1
Royal	0.01	0.02	0.94	0.95	0.1
Age	0.03	0.04	0.75	0.65	0.95
food	0.01	0.01	0.2	0.1	0.99
-	-	-	-	-	-
-	-	-	-	-	-
-	-	-	-	-	-
-	-	-	-	-	-

$$[\text{king} - \text{Boy} + \text{queen}] = \text{girl}$$

✧ Cosine similarity :-



$$\text{Distance} = 1 - \text{cosine sim}$$

$$\text{cosine} = 1 - \cos 45^\circ$$

$$= 1 - 0.70$$

$$= 0.29$$

movie recommendation

