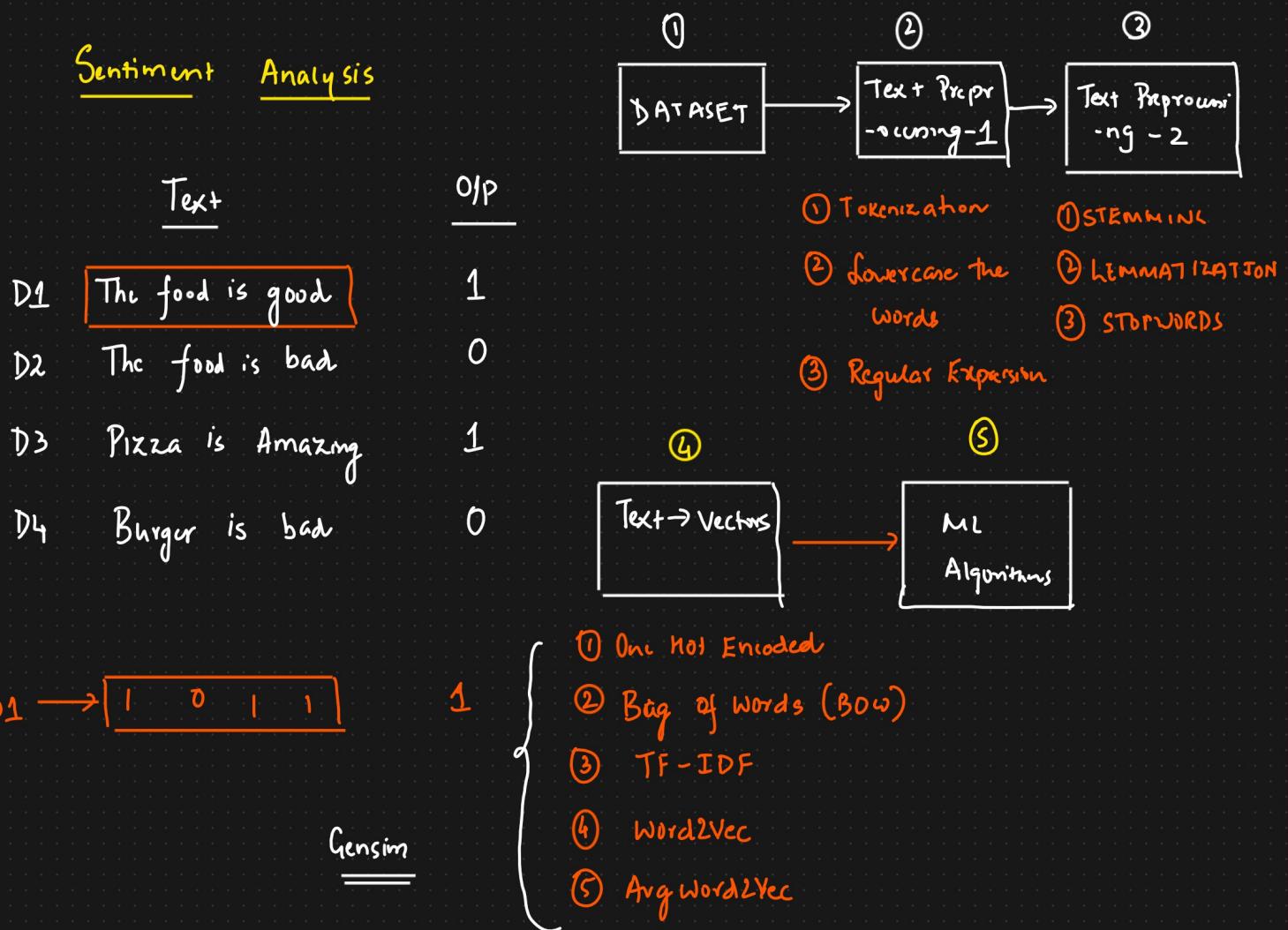
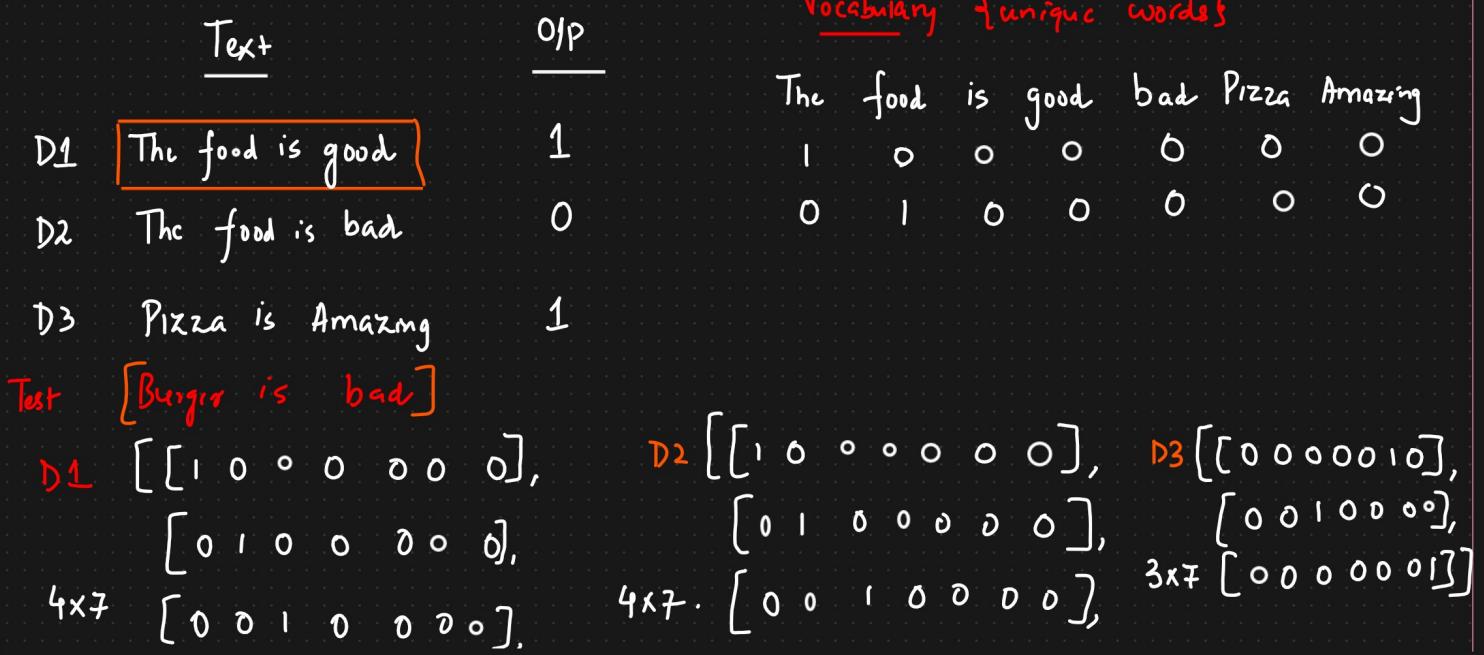


Text Preprocessing → What We have Learnt?



① One Hot Encoding

Vocabulary size: 7



$$\begin{bmatrix} 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix} \quad \begin{bmatrix} 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$$

$\left\{ \text{50K vocabulary size} \right\}$ —
Disadvantages

Advantages

① Easy to implement with python

[*sklearn OneHotEncoder, pd.get_dummies()*]

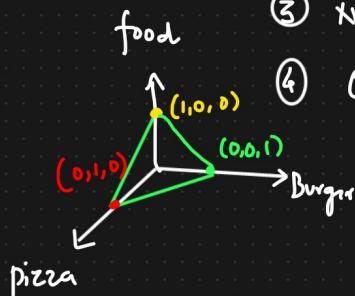
① Sparse matrix \rightarrow Overfitting

② ML Algorithm \rightarrow Fixed Size IIP

③ No semantic meaning is getting captured

④ Out of Vocabulary (OOV).

food	pizza	burger
1	0	0
0	1	0
0	0	1



② Bag of Words

Dataset

Text	O/P		
He is a good boy	1	Counts all the words case	S1 \rightarrow good boy
She is a good girl	1	\Rightarrow	S2 \rightarrow good girl good
Boy and girl are good	1	Stopwords	S3 \rightarrow Boy girl good [School] [Test]

Vocabulary	frequency		[good boy girl]	O/P
good	3	\downarrow	S1 [1 1 0]	1
boy	2	\Rightarrow	S2 [1 0 1]	1
girl	2	\downarrow	S3 [1 1 1]	1

Binary BOW and BOW

{ 1 and 0 }

{ Count will get updated
based on frequency }

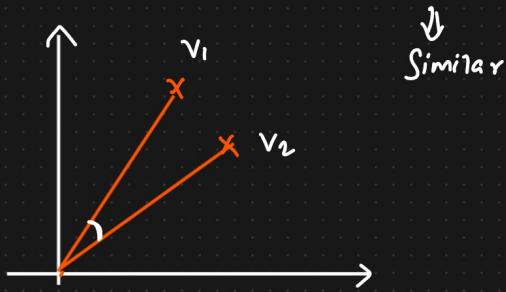
Advantages

- ① Simple and Intuitive
- ② Fixed Sized I/p \rightarrow ML Algorithms

Disadvantages

- ① Sparse matrix or array \rightarrow Overfitting
- ② Ordering of the word is getting changed
- ③ Out of Vocabulary (OOV).
- ④ Semantic meaning is still not captured.

{ The food is good \rightarrow $[1 1 1 0 1]$ $\rightarrow v_1$
 The food is not good \rightarrow $[1 1 1 1 1]$ $\rightarrow v_2$



④ N-grams Eg: bigrams, trigrams

$S_1 \rightarrow$ The food is good

$S_2 \rightarrow$ The food is not good

Bigram

food not good

1 0 1

1 1 1

$[$ food not good food good food not not good $]$

S_1 1 0 1 1 0 0 }

S_2 1 1 1 0 1 1 }

S_k learn \rightarrow n-grams = (1,1) \rightarrow unigrams

= (1,2) \rightarrow unigram, bigram

= (1,3) \rightarrow unigram, bigram, trigram

= (2,3) \rightarrow Bigram, trigram.

④ TF-IDF [Term Frequency - Inverse Document Frequency]

$S_1 \rightarrow \text{good boy}$

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

$S_2 \rightarrow \text{good girl}$

$S_3 \rightarrow \text{boy girl good}$

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

Term Frequency * IDF

	S_1	S_2	S_3	Words	IDF
good	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	good	$\log_e(3/3) = 0$
boy	$\frac{1}{2}$	0	$\frac{1}{3}$	boy	$\log_e(3/2)$
girl	0	$\frac{1}{2}$	$\frac{1}{3}$	girl	$\log_e(3/2)$

Final TF-IDF

	good	boy	girl	<u>Op</u>	<u>Bow</u>
Sent 1	0	$\frac{1}{2} \times \log_e(3/2)$	0	1	1 0
Sent 2	0	0	$\frac{1}{2} \log_e(3/2)$	1	0 1
Sent 3	0	$\frac{1}{3} \log_e(3/2)$	$\frac{1}{3} \log_e(3/2)$	1	1 1

Advantages

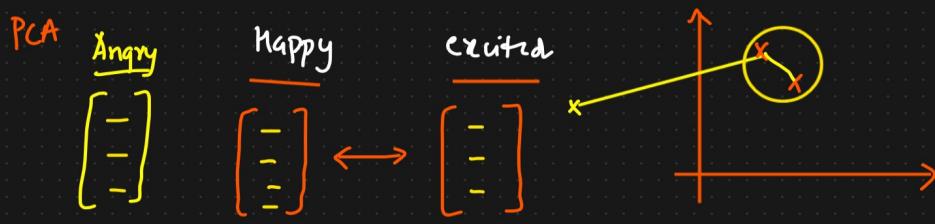
- ① Intuitive
- ② Fixed Size \rightarrow Vocab size
- ③ Word Importance is getting captured

Disadvantages

- ① Sparsity still exists
- ② OOV

Word Embeddings [Wikipedia]

In natural language processing (NLP), word embedding is a term used for the representation of words for text analysis, typically in the form of a real-valued vector that encodes the meaning of the word such that the words that are closer in the vector space are expected to be similar in meaning.



Words \rightarrow vectors

Sentence \rightarrow vectors

Word Embeddings

[ANN]

Count or frequency

- 1) One
- 2) Bow
- 3) TF-IDF

Deep Learning Trained Model

↓

Word2Vec

Cbow Skipgram

[Continuous BAG OF WORDS]

Word2Vec \rightarrow Feature Representation

↑ Google

Word2vec is a technique for natural language processing published in 2013. The word2vec algorithm uses a neural network model to learn word associations from a large corpus of text. Once trained, such a model can detect synonymous words or suggest additional words for a partial sentence. As the name implies, word2vec represents each distinct word with a particular list of numbers called a vector.

Vocabulary \rightarrow Unique words \rightarrow Corpus.

	Boy	Girl	KING	QUEEN	Apple	Mango
--	-----	------	------	-------	-------	-------

Gender

Gender	-1	1	-0.92	0.93	0.01	0.05
Feature Representation	0.01	0.02	0.95	0.96	-0.02	0.02
Royal	0.03	0.02	0.75	0.68	0.95	0.96
Age	-	-	-	-	0.91	0.92
Food	-	-	-	-	-	-

300 dimension !

$$\begin{bmatrix} - \\ - \\ - \\ - \\ - \\ - \end{bmatrix}$$

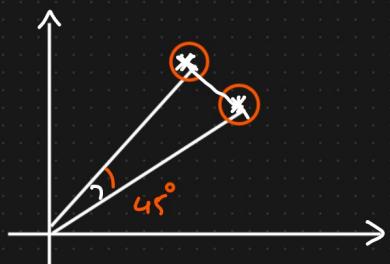
$$[\text{KING} - \text{BOY} + \text{QUEEN} = \text{GIRL}]$$

$$\text{KING} [0.95, 0.96] \quad \text{Man} [0.95, 0.98]$$

$$\text{QUEEN} [-0.96, 0.95] \quad \text{Women} [-0.94, -0.96]$$

$$\text{KING} - \text{MAN} + \text{QUEEN} = \text{WOMEN}$$

Cosine Similarity



$$\text{Distance} = 1 - \text{Cosine Similarity}$$

$$\text{Cosine-Sim} = \cos 45^\circ = \frac{1}{\sqrt{2}} = 0.7071$$

$$\text{Distance} = 1 - 0.7071 \\ \hookrightarrow = 0.29$$

$$\boxed{\sqrt{1-1}=0}$$

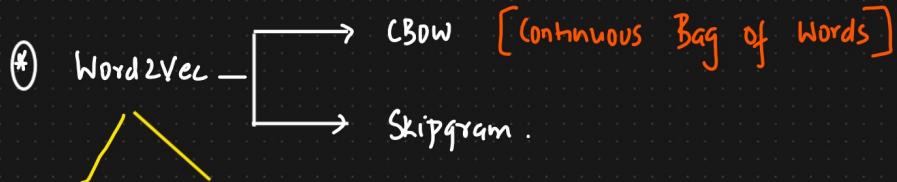
$$\text{Distance} = 1 - 0 \\ = 1$$



$$\text{Distance} = 1 - \cos 0^\circ \\ = 1 - 1 \\ = 0\%$$



ANN, Loss, Optimizers



Pretained Train A

Model Model For Scratch

② Skipgram - Word2Vec

Window Size = 5

→ O/p
[Interon, Company, Related, To]

→ Company, IS, To, DATA

→ Is, Related, DATA, SCIENCE

10

15

Related

Diagram illustrating a neural network layer:

- Inputs (IS):** A vertical vector of 10 inputs, labeled 0, 0, 1, 0, 0, 0, 0, 0, 0, 0.
- Weights:** A matrix labeled 7×5 representing the weight connections between the input layer and the current layer.
- Randomly Weights:** A label indicating the initial state of the weights.
- IP layer:** A label indicating the layer type.
- ANN:** A label indicating the overall system type.

5x7

Two rectangles drawn in yellow ink on a grid background. The top rectangle is oriented vertically and has a horizontal line extending from its top edge. The bottom rectangle is also oriented vertically.

for function ↓

Improve

CBow or Skipgram

① Generating the Training Data

When Should we apply CBOW or SkipGram.

2) Increase the window size
-vector dimension is also increasing.

Small Dataset \rightarrow CBOW
Huge Dataset \rightarrow SkipGram

Google Word2vec

Gunsim

3 billion words → Google News

feature representation of 300 dimension vectors]

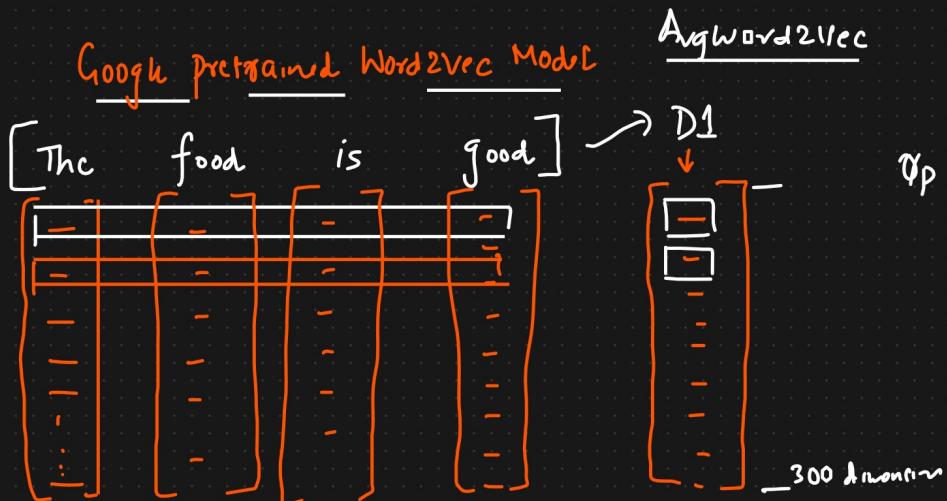
Cricket → [- - - - - - - - - -].

Advantages of Word2Vec

- f) Sparse Matrix \longrightarrow Dense Matrix
- f) Semantic Info is getting captured $\begin{bmatrix} \text{Moust, good} \end{bmatrix}$
- f) Vocabulary Size \longrightarrow Fixed set of dimension vectors
Google Word2Vec $\begin{bmatrix} 300 \text{ dimension} \end{bmatrix}$
- f) OOV is also solved

⑥ Avg Word2Vec

<u>Text</u>	<u>O/P</u>	<u>Avg Word2Vec</u>	<u>O/P</u>
D1 The food is good	1	$\begin{bmatrix} \cdot & \cdot \end{bmatrix}$	$\frac{1}{1}$
D2 The food is bad	0	$\begin{bmatrix} \cdot \end{bmatrix}$	$\begin{bmatrix} \cdot \end{bmatrix}$
D3 Pizza is Amazing	1	$\begin{bmatrix} \cdot \end{bmatrix}$	$\begin{bmatrix} \cdot \end{bmatrix}$



⑦ Gensim, Glove

\hookrightarrow pretrained Google Word2Vec

