

# NLP

- ① Numerical column  $\rightarrow$  Distribution analysis [Boxplot, Histogram]
- ② Categorical column  $\rightarrow$  Frequency analysis [Bar chart, Treemaps]
- ③ Date column  $\rightarrow$  Trend analysis [Line charts, Calendar]
- ④ Location column  $\rightarrow$  Geographical analysis [Maps]
- ⑤ Text data  $\rightarrow$  Bag of word analysis [Word cloud]

## Bag of word analysis

Tokens  $\rightarrow$  Frequency

Corpus — Collection of processed documents

Document — Collection of tokens

Tokens — Collection of words

— Unigram: One word / Token ✓

— Bigrams: Two words / Token

— Trigrams: Three words / Token

— Ngrams: N words / Token

[processed]



Bag of words  $\rightarrow$  Frequency analysis of Token

	Token	Frequency
T <sub>1</sub>	India	50
T <sub>2</sub>	digital	25
⋮	⋮	⋮
T <sub>N</sub>		



Word cloud  $\rightarrow$



Font size  $\approx$  Frequency

Word cloud

D <sub>1</sub>	[ - - - ]
D <sub>2</sub>	[ - - - ]
⋮	[ - - - ]
D <sub>3320</sub>	[ - - - ]

[ - - - ]
[ - - - ]
[ - - - ]
[ - - - ]
[ - - - ]

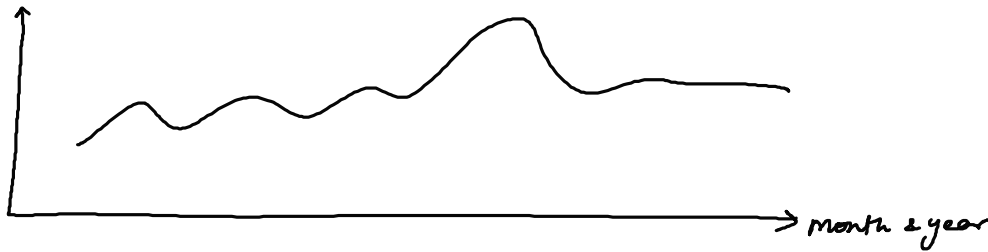
[This product is not good.]

[product good]

## Tent cleaning

- Lower case conversion
- Retained only alphabets & few characters
- Removed common & custom stop words

Hashtag's freq



## IMDB Sentiment

Tent classification → Sentiment analysis

- Rule based
- Tent classification

↳ Ticket system

↳ Intent identification

→ Feature Extraction

→ Stemming

→ Cause of high dimensionality ~ Thousands

→ Vectorization Tent → Number → Document Term matrix

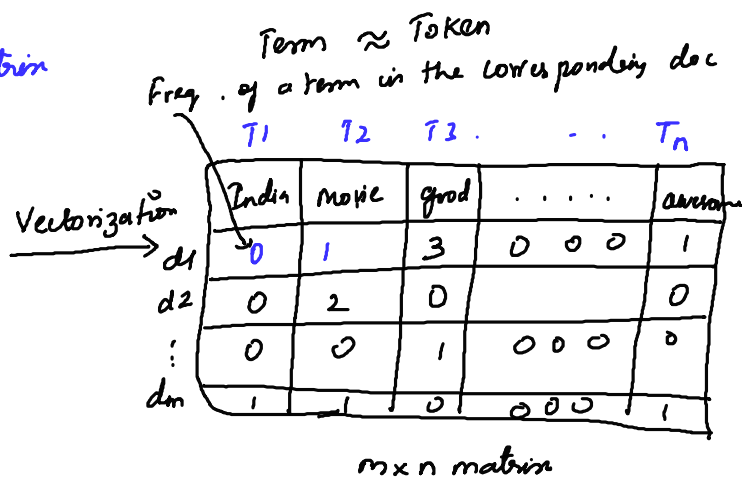
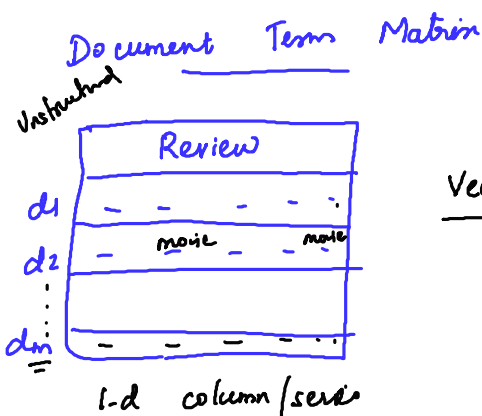
→ TFIDF

## Vectorization

Tent data → Structured numerical features

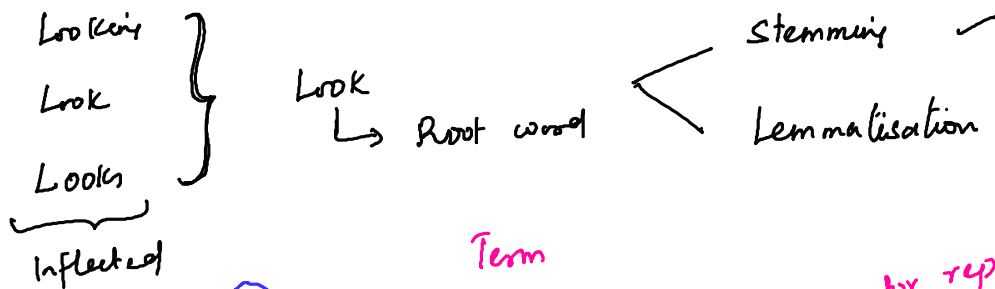
✓ Document Term Matrix

Word Embeddings  
(word2vec)



n → Total no. of unique tokens

m → Total no. of documents



Document

	T1	T2	T3	...	Tn
D1	1	0	0	1	0
D2	0	1	2	0	2
...	0				
Dm	1				

Words: movie, good

Vector representation of doc 1: [1, 0, 0, 1, 0, 1]

Columnwise sum:- Tokens frequency

Row-wise sum:- Document length  
↓  
No. of tokens

Sparse matrix:- 98% 0's

High dimension matrix

Vector representation of the token "movie" [1, 0, 0, 0, 0, 1]

Document Term Matrix - DTM - Doc classification/cluster.

$(DTM)^T \approx$  Term Document Matrix  $\rightarrow$  Word clustering (TDM)

### TF-IDF - Term Frequency - Inverse Document Frequency

	T1	T2	T3	T4
D1	1	0	0	1
D2	0	0	1	1
D3	0	1	1	2

DTM

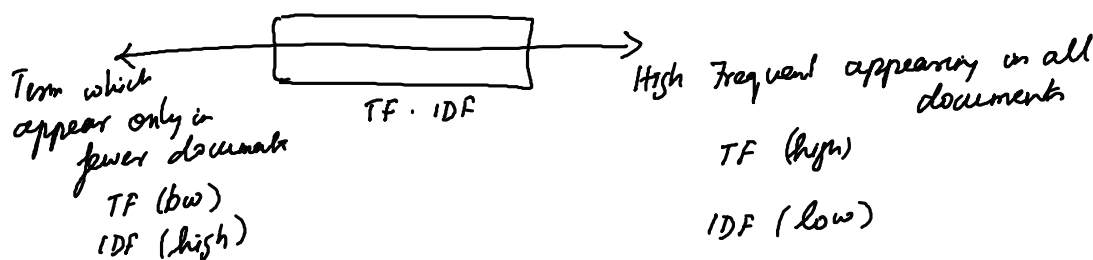
$$T.F(T_i, D_i) = \frac{\text{Frequency of the Term in the doc}}{\text{Total no. of tokens in the doc}}$$

$$T.F(T_1, D_1) = \frac{1}{2} = 0.5$$

$$\text{Inverse Document Frequency } (T_i) = \log \left( \frac{\text{No. of documents in the corpus}}{\text{No. of documents in which the term appears}} \right)$$

$$IDF(T_1) = \log \left( \frac{3}{1} \right) = \underline{\underline{\log(3)}}$$

$$IDF(T_4) = \log \left( \frac{3}{3} \right) = \log(1) = 0$$



$$TF.IDF(T_i, D_i) = \frac{\text{Freq of } T_i \text{ in } D_i}{\text{Total freq. in } D_i} * \log \left( \frac{\text{Total no. of docs in corpus}}{\text{Total no. of docs in which the } T_i \text{ appear}} \right)$$

This movie is good

Unigrams

This  
movie  
is  
good

Bigrams

This movie  
movie is  
is good

Trigram

This movie is  
movie is good