

Bag of Words: (BOW) \Rightarrow frequency of words

Tom is travelling in bus

John is travelling in bus

Emma took train and bus

D1 \Rightarrow Tom travelling bus.

D2 \Rightarrow John travelling bus.

D3 \Rightarrow Emma took train bus.

Vocabulary:

Tom - 1

travelling - 2

bus - 3

John - 1

Emma - 1

took - 1

train - 1

[Words are arranged in descending order of frequency and kept in the bag.] (X)

[bus, travelling, tom, john, Emma, took, train]

[Bus travelling tom John Emma took train] \rightarrow (look)

1 1 1 0 0 0 0 \Rightarrow D-1

1 1 0 1 0 0 0 \Rightarrow D-2

1 0 0 0 1 1 1 \rightarrow D-3

New sentence \rightarrow Tom took bus (then) train (then) bus

2 0 1 0 0 1 1



\rightarrow oov \leftarrow

Disadvantage:

- * sparse matrix
- * words are shuffled based on frequency.
- * semantic meaning is unclear.
- * OOV.

Advantage:

- * No need of fixed length document.
- * Easy implementation & cost effective.

⇒ TF-IDF :

↳ [Term Frequency - Inverse Document Frequency]

⊗ Term Frequency: Focused on words in each sentence

$$TF : \frac{\# \text{ repetition of words in sentence}}{\# \text{ words in the sentence}}$$

⊗ Inverse document Frequency: Focused on the sentences in the corpus.

$$IDF : \log_e \left(\frac{\# \text{ sentences}}{\# \text{ sentences containing the word}} \right)$$

D-1 : Good boy

D-2 : Good girl

D-3 : Boy girl good

[Good, boy, girl] \Rightarrow vocabulary

Term Frequency — step:1

D1 D2 D3

Good $\frac{1}{2}$ $\frac{1}{2}$ $\frac{1}{3}$

boy $\frac{1}{2}$ 0 $\frac{1}{3}$

girl 0 $\frac{1}{2}$ $\frac{1}{3}$

IDF : step:2

Good $\rightarrow \log_e\left(\frac{3}{2}\right)$

boy $\rightarrow \log_e\left(\frac{3}{2}\right)$

girl $\rightarrow \log_e\left(\frac{3}{2}\right)$

step-3

D-1

0

$\frac{1}{2} \times \log_e\left(\frac{3}{2}\right)$

0

D-2

0

0

$\frac{1}{2} \times \log_e\left(\frac{3}{2}\right)$

D-3

0

$\frac{1}{3} \times \log_e\left(\frac{3}{2}\right)$

$\frac{1}{3} \times \log_e\left(\frac{3}{2}\right)$

Disadvantage

- OOV
- sparse matrix

Advantage:

- word importance is captured.
- length of sentences can be dynamic.
- Easy to implement.

Word Frequency and encoding:

- one hot encoding
- Bag of Words (Bow)
- TF-IDF

Common Issue

- sparse matrix
- oov

To fix the above issue Google came up with the idea of deep learning trained library in 2013

↳ ANN



⇒ Word2Vec - CBOW (continuous bag of words)

vocabulary \Rightarrow Boy Girl king Queen Mango Apple Women

Gender 1 -1 0.98 -0.97 0.01 0.02 -0.99

Royal 0.01 0.03 0.99 0.98 0.02 -0.01 0.01

Age 0.9 0.91 0.9 0.89 0.4 0.45 0.86

Food 0.01 0.01 0.02 0.02 0.99 0.96 0.01

\Downarrow

dimensions

\Rightarrow Gender [king - boy + Queen] = Women

\rightarrow Gender [0.98 - 1 - 0.97] \Rightarrow 0.99 (cov is fixed)

	Mango	Apple	Orange	<u>Grapes</u>
Gender	0.01	-0.03	0.02	0.06
Food	0.97	0.96	0.94	0.92

→ Gender [Mango - Apple + orange]

Gender [0.01 + 0.03 + 0.02] → 0.06 → Grapes do not belong to Gender

→ Food [Mango - Apple + orange]

Food [0.97 - 0.96 + 0.94] → 0.92 → Grapes belong to Food