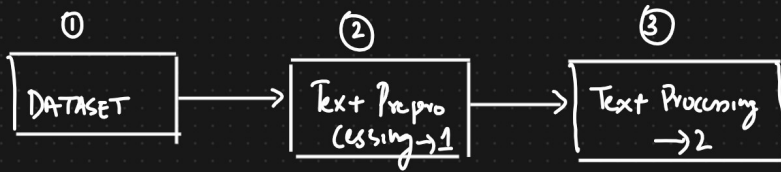


NLP In Machine Learning



the, he, she, of, is, you,

↓

Sentiment Analysis

① Tokenization

② lowercase the words

③ Regular Expression

① STEMMING

② Lemmatization

③ STOPWORDS

④

Text \rightarrow Vectors

① One hot Encoding

② Bag of words (Bow) \Rightarrow NL

③ $\overline{I_f} - I_{df}$

(4) Word2Vec

⑤ Antwort 2 Vec

} \Rightarrow Deep learning

Tokenization

Topics

1) Corpus \longrightarrow Paragraph

2) Documents \rightarrow Sentences

3) Vocabulary → unique words

4) Words

Carpus

"My name is KRISH and I have a interest in teaching ML, NLP and DL. I am also a good human".

↓ Tokens

↓ Tokens

① My name is KRISH and I have an interest in " " " "

② I am also a good human.

Vocabulary size = 18 words.

① One hot Encoding

Text

D1 The food is good

D2 The food is bad

D3 Pizza is Amazing

0/p

1

0

1

Sentiment Analysis Problem

vocabulary
↗

the food is good bad pizza amazing

1 0 0 0 0 0 0

0 1 0 0 0 0

Test [Burger is Bad] X

D1 $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$

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

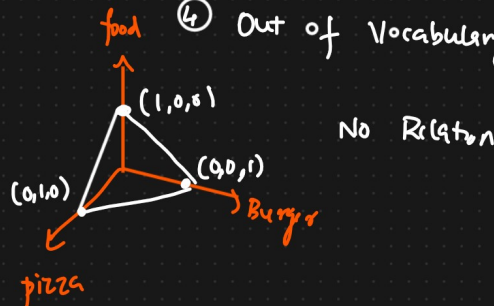
D3 $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$

Advantages

① Easy to Implement with python

[Sklearn OneHotEncoder]

| food | pizza | burger |
|------|-------|--------|
| 1 | 0 | 0 |
| 0 | 1 | 0 |
| 0 | 0 | 1 |



Disadvantages

① Sparse Matrix \rightarrow Overfitting

② ML Algorithms \rightarrow Fixed Size I/p.

③ No semantic is getting captured

④ Out of Vocabulary (OOV)

② Bag of Words (Bow)

Datant

| Text | O/p |
|-----------------------|-----|
| ① ③ ② good boy | 1 |
| She is a good girl | 1 |
| Boy And girl are good | 1 |

downcase all the words \Rightarrow Stop words

S1 \rightarrow good boy boy
S2 \rightarrow good girl
S3 \rightarrow boy girl good

| Vocabulary | frequency | f ₁ [good] | f ₂ boy | f ₃ girl | O/p | Binary Bow |
|------------|-----------|--------------------------|-----------------------|------------------------|-----|------------|
| good | 3 | 1 | 1 | 0 | 1 | |
| boy | 2 | 1 | 0 | 1 | 1 | |
| girl | 2 | 1 | 1 | 1 | 1 | |

Advantages

- ① Simple and Intuitive
- ② Fixed I/P Size \Rightarrow ML Algorithms.

Disadvantages

- ① Sparse Matrix \rightarrow overfitting
- ② Out of Vocabulary (OOV Problem)
- ③ Semantic meaning not there. ✗

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

N-grams

Eg: unigram, bigram, trigram

$\hookrightarrow \downarrow \downarrow$
(1,1)

| | | food | not | good |
|------------------|----------------------|------|-----|------|
| S1 \rightarrow | The food is good | 1 | 0 | 1 |
| S2 \rightarrow | The food is not good | 1 | 1 | 1 |

Bigram (1,2)

Trigram (1,3)

| | food | not | good | food not | not good | food good |
|----|------|-----|------|----------|----------|-----------|
| S1 | 1 | 0 | 1 | 0 | 0 | 1 |
| S2 | 1 | 1 | 1 | 1 | 1 | 0 |

Sklearn \rightarrow n-gram = (1,1) \rightarrow unigrams

(1,2) \rightarrow unigrams, bigrams

(1,3) \rightarrow unigrams, bigrams, trigrams

(2,3) \rightarrow bigram, trigram

(2,2) \rightarrow Bigram

③ TF-IDF [Term frequency - Inverse Document Frequency].

S1 \rightarrow good boy

S2 \rightarrow good girl

S3 \rightarrow boy girl good

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

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

| <u>Term frequency</u> | | | <u>IDF</u> | |
|-----------------------|---------------|---------------|---------------|---|
| | <u>S1</u> | <u>S2</u> | <u>S3</u> | <u>Words</u> |
| good | $\frac{1}{2}$ | $\frac{1}{2}$ | $\frac{1}{3}$ | $\log_c \left(\frac{3}{\frac{1}{3}} \right) = 0$ |
| boy | $\frac{1}{2}$ | 0 | $\frac{1}{3}$ | $\log_c \left(\frac{3}{\frac{1}{2}} \right)$ |
| girl | 0 | $\frac{1}{2}$ | $\frac{1}{3}$ | $\log_c \left(\frac{3}{\frac{1}{2}} \right)$ |

| | <u>Final TF-IDF</u> | | | <u>O/P</u> |
|--------|---------------------|---------------------------|---------------------------|------------|
| | <u>f1</u> good | <u>f2</u> boy | <u>f3</u> girl | |
| Sent-1 | 0 | $\frac{1}{2} \log_c(3/2)$ | 0 | 1 |
| Sent-2 | 0 | 0 | $\frac{1}{2} \log_c(3/2)$ | 1 |
| Sent-3 | 0 | $\frac{1}{3} \log_c(3/2)$ | $\frac{1}{3} \log_c(3/2)$ | 1 |

Advantages

- ① Intuitive
- ② Fixed Sized I/p \rightarrow Vocab Size
- ③ Word Importance is getting Captured

Disadvantages

- ① DOV
- ② Sparsity is still exists



Word2Vec