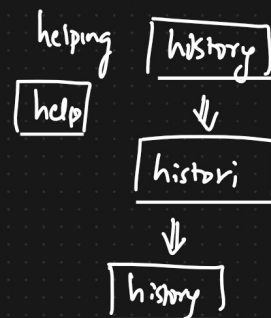


```

graph LR
    1[① DATASET] --> 2[② Text + Prepro-  
cessing → 1]
    2 --> 3[③ Text Processing  
→ 2]

```



↓

Sentiment Analysis

- ① Tokenization
- ② Lowercase the words
- ③ Regular Expression

④

Text \rightarrow Vectors

- ① STEMMING
- ② Lemmatization
- ③ STOPWORDS

- ① One hot Encoding
- ② Bag of words (Bow)
- ③ Tf-idf
- ④ Word2Vec
- ⑤ Avg word2vec

} \Rightarrow ML

} \Rightarrow Deep learning

Topics

- 1) Corpus \rightarrow Paragraph
- 2) Documents \rightarrow Sentences
- 3) Vocabulary \rightarrow unique words
- 4) Words

Vocabulary size = 18 words.

Corpus { "My name is KRISH and I have a interest
in teaching MC, NLP and DL. I am also a
good human".
↓ Tokens

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

(2) I am also a good human.

① One hot Encoding
Text

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	0

D1 The food is good

D2 The food is bad

D3 Pizza is Amazing

Test [Burger is Bad] X

D1 $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}$
 The
 food
 is
 4x7
 good

D2 $\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 & 0 \end{bmatrix}$
 4x7

D3 $\begin{bmatrix} 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$
 3x7

Above we have converted the documents into respective vector values using the vector mapping formed using vocabulary

Advantages

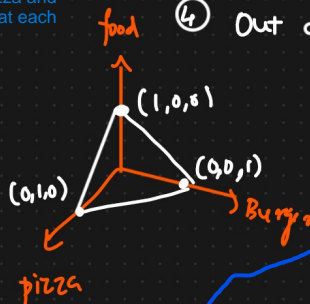
① Easy to Implement with python

[Sklearn OneHotEncoder]

For ex we want to capture semantic/relationship b/w 3 words: food, pizza and burger. If we plot the respective sparse matrix then we can observe that each of these are equidistant from each other representing no or constant relationship

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

Example



During testing or performing prediction over new dataset if we are provided with document containing vocabulary that were present in the set of vocabulary used while training the model then the condition that will arise is called OOV (out of vocabulary)

Disadvantages

① Sparse Matrix → Overfitting

② ML Algorithms → Fixed Size I/p.

③ No semantic is getting captured

④ Out of Vocabulary (OOV)

We are basically creating sparse matrix (matrix with lot of 1 and 0) which will need lot of space in RAM and will lead to overfitting.

Relationship b/w 2 words is not getting captured that's why Semantics not getting captured.

ML algorithm works with fixed Input size but here Input size is not fixed. For Ex; D3 have only 3 words meaning 3 i/p while D1 and D2 have 4 i/p. Hence i/p size is not fixed.

② Bag of Words (BOW)

Dataset

Text

O/P represents the sentiment of the respective sentence or document.
 If 1 then +ve sentiment and if 0 then -ve sentiment.

O/p

Bag of words can be used to resolve the disadvantages raised in One hot encoding approach. By using Small case conversion and Stopwords we are lowering the dimensions.
 In Binary BOW : have only 0 and 1.

① ② ③ good boy

1

lowercase all the words

S1 → good boy boy

She is a good girl

1

⇒

S2 → good girl

Boy And girl are good

1

Stopwords

S3 → boy girl good

Sentences formed by lowering the case followed by removing the stop word of respective Text or Document

Vocabulary

frequency

f₁ f₂ f₃
 [good boy girl]

O/p

Binary BOW

good

3

S1

1

1

0

1

boy

2

S2

1

0

1

1

girl

2

S3

1

1

1

1

In above after converting text into vector using vocabulary mapping we can clearly observe that S1, S2 and S3 have more 1s as compared to 0s. This is the reason why their respective output is also 1

Advantages

- ① Simple and Intuitive (This advantage is common for all the approaches)
- 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

S1 \rightarrow The food is good

S2 \rightarrow The food is not good

\rightarrow

\rightarrow

food	not	good
1	0	1
1	1	1

Unigram (1,1)

Bigram (1,2)

Trigram (1,3)

	food	not	good	+	food not	not good	food good	=	Bigram
S1	1	0	1		0	0	1		
S2	1	1	1		1	1	0		

Unigram

Combination of 2 vocabulary at a time

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]

In TF-IDF, based on the frequency of word we will decide word is important or not.

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

Term frequency



IDF

IDF for good is 0 means it is not important word since it is present in all the sentences.

	S1	S2	S3
good	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$
boy	$\frac{1}{2}$	0	$\frac{1}{3}$
girl	0	$\frac{1}{2}$	$\frac{1}{3}$

Words

good	$\log_c \left(\frac{3}{3} \right) = 0$
boy	$\log_c \left(\frac{3}{2} \right)$
girl	$\log_c \left(\frac{3}{2} \right)$

Final TF-IDF

Final TF-IDF obtained by multiplying respective TF and IDF values. That is first IDF value is multiplied with all the elements of first column in TF so on.

	f1 good	f2 boy	f3 girl	<u>O/P</u>
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

May be final output is representing the sentiment of the respective sentence based on OR logical gate.

Advantages

- ① Intuitive (Simple and intuitive)
- ② Fixed Sized I/p \rightarrow Vocab Size
- ③ Word Importance is getting captured

Disadvantages

- ① DOV
- ② Sparsity is still exists



Word2Vec

Problem of sparsity can be resolved by using Word2Vec technique for the conversion of Text into Vector. This will be discussed in the deep learning sessions.