

Unit-6 (honor)

Page No.

Date

Text processing - It deal with Text data.

Before prediction we need to perform text preprocessing.

Some steps are.

- * Removing punctuations like . , \$ () % @
- * Removing URL
- * Removing stop word
- * lower casing
- * ~~Stop~~ Tokenization
- * Stemming
- * Lemmatization.

Removing punctuations

Import String

~~String~~ fun

```
def remove_punctuation(text):  
    final = "".join([i for i in text if i not  
                      in string.punctuation])  
    return final.
```

lowering case

```
d['a'] = data['b'].apply(lambda x: x.lower())
```

Stop word removal

→ I, am, not, is, g

Page No.

Date

Import nltk

Stopword = nltk.corpus.stopwords.words('english')

Stemming

converting to root base
But meaning may loss.

from nltk.stem.porter import PorterStemmer

post_st = PorterStemmer()

new = ~~old~~ post_st.stem(~~word~~old)

lemmatization

to root base but have meaning

from nltk.stem import WordNetLemmatizer

w_L = WordNetLemmatizer()

new = w_L.lemmatize(old)

Tokenization

- Splitting a string
- text to list of tokens
- text to sentence
→ text to word
- text to RE

Types.

① Word Tokenization

```
from nltk.tokenize import word_tokenize
```

```
text = "hello everyone. welcome"
word_tokenize(text)
```

② Character Tokenization

③ Sentence Tokenization

④ Subword Tokenization
lower → lower-er
Smartst → smart-test

Feature extraction

→ The process of converting textual into number is called vectorization.

→ This feature extraction

→ Easy for computation.

Two main method

① bag-of-words with TF-IDF
bag of word.

— It is representation model of document data.

— Count the no. of word in document.

```
from sklearn.feature_extraction.text import
CountVectorizer
```

disadvantage

① sparse matrix

② can't relate with similar matrix

③ Not relationship

```
Sent = ["This is a very good boy"]
```

```
vector vt = CountVectorizer()
```

```
vt.fit(Sent)
```

```
v = vt.transform(Sent)
```

```
Ans = v.toarray()
```

Ans

TF-IDF

TF - Measure the frequency

$$TF = \frac{\text{count of } d \text{ in } d}{\text{count of no. of word in } d}$$

IDF - Measure the importance of a word

$$IDF = \log \left(\frac{\text{Total no. of document in corpus}}{\text{no. of document containing term}} \right)$$

$$TFIDF = TF \times IDF$$

Here eg. "the" can occur more time But

It has less significance so, TFIDF will detect it. It has less TFIDF value.

② Word Embedding

→ representing numeric input that represent a word in lower-dimensional space.

It allow word with similar meaning to have similar representation.

→ It have feature such as Age, Sports food,

etc. and many

Man →



Page No.

Date

(a) Word2vec

- each word is represented as a vector of 32 or more dimension.
- semantic info stored

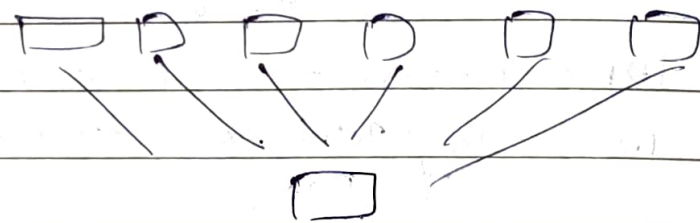
King - Man + Woman = Queen

Two algorithms

(a) Continuous Bag of Word (CBOW)

- predict the target word from the context

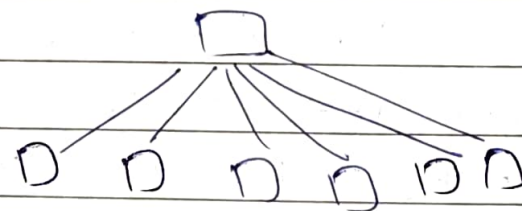
neural network is used



(b) Skip Gram

- predict the context word from target

neural network is used



(b) Glove

derive the relationship b/w the word from Global Statistics.

eg

- 1) I play cricket
- 1 love cricket
- 1 love football

Page No.

Date

	1	play	cricket	football	love
1	0	1	0	0	2
play	1	0	1	0	0
cricket	0	1	0	0	1
football	0	0	1	0	1
love	2	0	1	1	0

Topic Modeling

- discovering abstract "topics" in documents.
- unsupervised classification perform.

Latent Dirichlet Allocation

Aim - find topic.

first step document word matrix

	w_1	w_2	w_3	w_4	w_5	w_6	w_7
D_1	0	1	1	1	0	0	1
D_2		0	0	0	1	1	1
D_3		1	1	0	0	1	0
D_4		0	0	1	1	0	0

document topic matrix

word topic matrix

	k_1	k_2	k_3	...	w_1	w_2	w_3
D_1					k_1		
D_2					k_2		
D_3							

Step 3 Decrement Count by 1 from document topic matrix and same with word topic matrix

Step 4 Calculate

$$P(t_k | d_i) = \frac{n_{ik} + \alpha}{N_i + K\alpha}$$

 n_{ik} = no. of word in i th doc and k th topic

 α = hyper parameter (0.1)

 N_i = total no. of word in i th document

 K = total no. of topic

	T1	T2	T3
Doc 1	2-1	1-1	3-1
Doc 2	2-1	3-1	1-1
Doc 3	1-1		3-1

$$P(t_k | d_i) = \frac{1 + 0.1}{6 - 1 + 3 \times 0.1} = 0.188$$

Step 5 Calculate $P(w_j | T_k)$

	T1	T2	T3
good	12-1	2-1	20-1
train	13-1	10-1	11-1
home	11-1	14-1	25-1
tree	18-1	24-1	

 M_{jk} = no. of time word

 k come in topic K
 β = hyper parameter

 V = to. no. of word

$$P(w_j | t_k) = \frac{M_{jk} + \beta}{\sum_{j \in V} M_{jk} + V\beta} = \frac{11 + 0.1}{47 + 4 \times 0.1} = 0.023$$

Step 6 Calculate $P(w_j | t_k, d_i)$

$$P(w_j | t_k, d_i) = P(t_k | d_i) \times \beta P(w_j | t_k)$$

Step 7 resampling

for a given word w_j and document i find the topic $P(w_j | t_k, d_i)$ is maximum

Step 8 Repeat all step till iteration

Text Similarity.

Used to discover most similar text.

1. Jaccard Similarity.

= ratio of common word to total unique words

= 1 represent high similarity

0 represent low similarity

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

2. Cosine Similarity.

= measure cosine of angle b/w two vectors

= Vector can be bag of words, TF-IDF, or Doc2Vec

$$\text{similarity}(A, B) = \frac{A \cdot B}{\|A\| \times \|B\|}$$