

Name → Himanshu Gupta

Class/Sem. → 3rd Sem, 4th year

Roll No → BT17CSE093

Branch → CSE

Subject → NLP

Date → 24/11/2020

## 1) Lemmatization & Stemming

(1) Lemmatization → It is the task of determining that two words have the same root, despite their surface differences. It takes into consideration the morphological analysis of words. Below we illustrate the method with examples in both two languages English & ~~Hindi~~ Sanskrit.

Language	Form	Lemma	Morphological information
English	Studies	Study	Third person, singular, present tense.
"	Studying	Study	Present of Verb
"	studied	study	Third person, singular, past tense.
Sanskrit	पढ़त	पढ़	Third person, singular, present tense
Sanskrit	पढ़ते	पढ़	<del>Th<sup>2nd</sup></del> person, <del>एवं ते तात्</del> singular, present
Sanskrit	पढ़त	पढ़	Third person, singular, past tense.

Stemming → It works by cutting off the end or the beginning of the word. It considers a list of common prefixes & suffixes that can be found in an inflected word.

Below we have examples in both language in English & Sanskrit.

form	Suffix	stem
studies	-es	study
studying	-ing	study
studied	-ed	study
पढ़ते	-ते	पढ़ा
पढ़ाये	-ये	पढ़ा
पढ़ा	-ा	पढ़ा

So, from example it is clear that a lemma is a base form of all its inflectional form whereas a stem isn't.

Ans 2

Python code (Remember to install nltk library first)

```

from nltk.tokenize import word_tokenize
from nltk.corpus import wordnet
antonyms = []
text = input().strip()
# tokenization
tokens = word_tokenize(text)
print(tokens)
# antonyms of all words
for word in tokens:
    for syn in wordnet.synsets(word):
        for abc in syn.lemmas():
            if abc.antonyms():
                antonyms.append(abc.antonyms())
print(antonyms)

```

if abc. ontonyms ( ) :

antonyms: append (abc. antonyms ([0]). name())

print( word , set(antonyms) )

卷之三

三

doc-1 My name ~~is~~ Himanshu Gupta, a CSE student at

TTT Nappus

doc-2 → I live in a small town of Jaipur

in Rajastan.

Present I am living in Seminary hills.

do you live in Orange city.

see Mr. Webster today.

my hobbies is reading, dancing, singing,  
playing chess & ~~do~~ and sleeping.

To find total weight of ~~does~~ all words of doe-1

(1) my name = my

$$tf = \frac{1}{10}$$

$$idf = \log_{10}\left(\frac{5}{2}\right)$$

So,

$$tf \cdot idf = tf \times idf = \frac{1}{10} \times 0.398 = 0.0398$$

(2) pa. name

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{1}\right) = 0.0699$$

(3) Himanshu

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{1}\right) = 0.0699$$

(4) Crypta

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{1}\right) = 0.0699$$

(5) a

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{2}\right) = 0.0398$$

(6) CSE

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{1}\right) = 0.0699$$

(7) Student

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{1}\right) = 0.0699$$

(8) act

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{1}\right) = 0.0699$$

(9) IIT

$$tf \cdot idf = \frac{1}{10} \times \log_{10}\left(\frac{5}{1}\right) = 0.0699$$

(o) Nagpur:

$$tf \cdot idf = \frac{f}{f_0} \times \log_{10}\left(\frac{N}{f}\right) = 0.0699$$

Inverted index construction

Term	doc ID	→
my	1	5
name	1	
Himanshu	1	
Crypts	1	
a	1	2
CSE	1	
student	1	
at	1	
IIT	1	
Nagpur	1	

Now

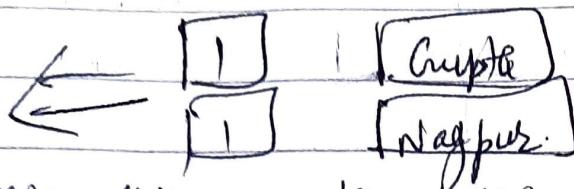
Term doc freq → Posting lists.

Term	doc freq	→	Posting lists
my	2		1 → 5
name	1		1
Himanshu	1		1
Crypts	1		1
a	2		1 → 2
CSE	1		1
student	1		1
at	1		1
IIT	1		1
Nagpur	1		1

## Query Processing

Let Query = Crypta and Nagpur

Locate Crypta in dictionary & Nagpur in Dict.



Merge this posting lists

→ walk through this list simultaneously and merge.



Limitation of Bag-of-words feature in sentiment classification:-

When we consider bag of words method to generate vectors for large documents. then the resultant vectors will be of large dimensions and will contain far too many values resulting in sparse vectors.

Apart from resulting in sparse, Bag of words does a poor job in making sense of text data for example, consider two sentences:-

"I love tea but not coffee."

"I love coffee but not tea".

Now, Bag of words approach will create similar vectors even both sentences have different meaning.

So we can solve this problem by using following techniques.

(i) Positional Indexing :- Creating positional indexing for each word of sentence will solve this problem.

If it is able to distinguish between two sentences/documents

(ii) Converting all words to lower case :-

(iii) Removing ~~stop words~~ Using term frequency (tf)

(iv) Stemming and Lemmatization. E.g. Inverse df (ids)

The process of converting words into numbers are called vectorisation.

#### Step 1 → Tokenizing the Sentence

Sentence 1	Sentence 2	Sentence 3
I am studying in IIT Nagpur.	Nagpur is summer the capital summer of capital of Maharashtra	Nagpur is famous for changes.

#### Step 2 → Creating Dictionary of word frequency

word	frequency
I	1
am	1
studying	1
in	1
IIT	1
Nagpur	3
is	2
the	1

Summers |  
 Capital |  
 of |  
 Maharashtra |  
 famous |  
 for |  
 oranges |

Now, sort word according to frequency.

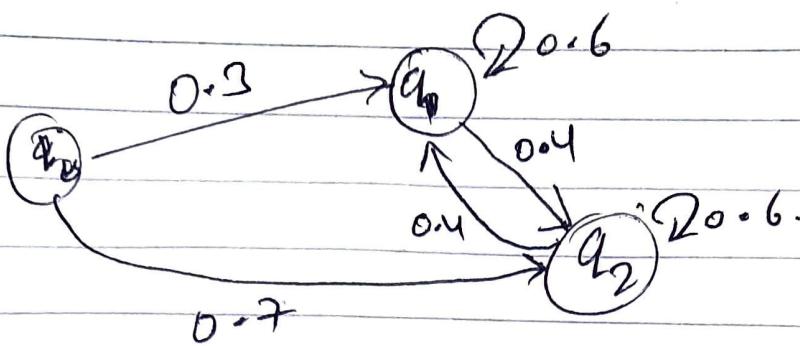
Now, creating vector of word having bit values.

Word	Sentence 1	Sentence 2	Sentence 3
I	1	0	0
am	1	0	0
studying	1	0	0
in	1	0	0
JIIT	1	0	0
Nagpur	1	1	1
is	0	1	1
the	0	1	0
summers	0	1	0
capital	0	1	0
of	0	1	0
Maharashtra	0	1	0
famous	0	0	1
for	0	0	1
oranges	0	0	1

Ans 5)  $Q = \{1, 2\}$ ,  $M = [Q, \Sigma, \delta, Y_1, q_0]$

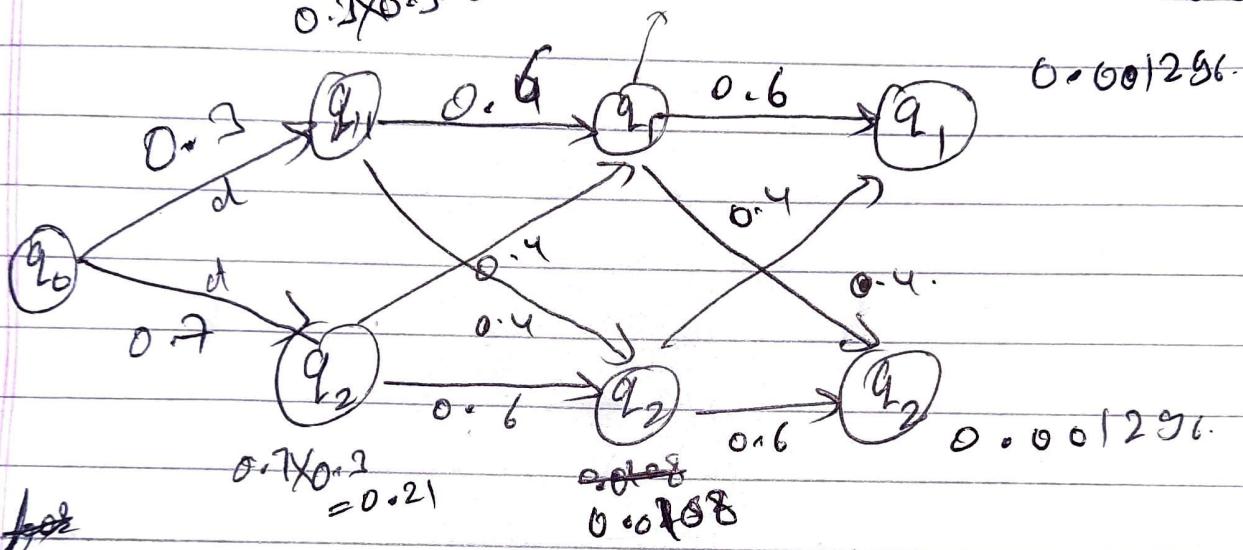
Let start =  $q_0$ ,

Now,



$$0.3 \times 0.3 = 0.09$$

$$\min(0.09 \times 0.6, 0.21 \times 0.4) = \cancel{0.054} = 0.00864$$



for string "dab", State sequence is

$$q_0 \rightarrow q_1 \rightarrow q_2 \rightarrow q_2$$

$$\text{Prob} = 0.09 \times 0.00864 \times 0.001296$$

$$= 1.26 \times 10^{-6}$$