

```
from zipfile import zipfile  
import glob
```

```
import pandas as pd.
```

```
import nltk
```

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```
from sklearn.metrics.pairwise import linear_kernel
```

```
from nltk.corpus import stopwords
```

```
import warnings.
```

```
warnings.filterwarnings('ignore')
```

Exe :- 1

```
① file_name = "movies.zip"  
with zipfile(file_name, 'r') as zip:  
    zip.printdir()
```

```
files = [file for file in glob.glob("movies/*")]  
files.
```

```
nltk.download('punkt')
```

```
nltk.download('stopwords')
```

```
stop_words = set(stopwords.words('english'))
```

```
tokenizer = nltk.tokenize.WhitespaceTokenizer()
```

```
from nltk.stem import PorterStemmer.
```

```
ps = PorterStemmer()
```

```
from nltk.stem import LancasterStemmer.
```

```
ls = LancasterStemmer()
```

```
from nltk.stem import WordNetLemmatizer
```

```
lemmatizer = WordNetLemmatizer()
```

## Natural Language Processing Lab

### Lab5. Stemming and Lemmatization on Movie Dataset

#### EXERCISE-1

The file movie.zip contains 20 files about various movies. For each of the files in movies.zip, you will have to do the following:

- How many sentences in each file?
- How many tokens in each file?
- How many tokens excluding stop words in each file?
- How many unique stems (ie., stemming) in each file? (Use PorterStemmer)
- How many unique stems (ie., stemming) in each file? (Use LancasterStemmer)
- How many unique words (ie., lemmatization) in each file? (Use WordNetLemmatizer)
- Pretty Printing: Print the details of A to E in the following order

File Name	Sentences	Tokens	Tokens-Only	StemsPorter	StemsLancaster	Lemmas
-----	-----	-----	-----	-----	-----	-----

#### EXERCISE-2

In this exercise, you will build your Term-Document Matrix for this movie collection of 20 movies. In order to improve the similarity search experience, you will use only lemmatized terms for creating the matrix.

##### Step-1

For each movie:

- Tokenize terms and build list of tokens
- Find lemmatized words from the tokens

##### Step-2

Build Term-Document matrix using TfidfVectorizer

##### Step-3

Take vectors of any two movies and compute cosine similarity

#### EXERCISE-3

Will lemmatized matrix help to achieve better similarity search or not?. Please comment.

for file in files:

with open (file, 'r', encoding='cp1252') as f:

Contents = f.readlines()

print (Contents)

print ("\*\*\*\*\*")

print ("")

A. How many sentences in each file?

B. How many tokens in each file?

C. How many tokens excluding stop words in each file?

files = [file for file in glob.glob('movies/\*')]

for file in files:

with open (file, 'r', encoding='cp1252') as f:

for row in Contents:

sent\_text = nltk.sent\_tokenize(row)

print("sentence tokenize", len(sent\_text))

for row in contents

words = nltk.word\_tokenize(row)

print ("Word tokenize", len(words))

filtered\_sentences = [w for w in words if not w in stop\_words]

print ("stopwords", len(filtered\_sentences))

print ("\*\*\*\*\*")

D) def port\_stem\_sentence(sentence):

tok = tokenizer.tokenize(sentence)

filtered\_sentence = [w for w in tok if not w in stop\_words]

stem\_sentence = []

for word in filtered\_sentence:

stem\_sentence.append(nltk.stem.porter.stem(word))

## NOTES

for files .in files:

with open (file, 'r', encoding='cp1252') as f:

contents = f.readline()

print ("porter-stemming")

print (port-stem Sentence (contents))

print ("\*\*\*\*\*")

E). for lan-stem Sentence (sentence):

tok = tokenizer.tokenize (sentence).

filtered\_sentence = [w for w in tok if not w in stop-words]

stem\_sentence = []

for word in filtered\_sentence:

stem\_sentence.append (l.s. stem (word))

return lan (stem\_sentence)

for file .in files:

with open (file, 'r', encoding='cp1252') as f:

contents = f.readline()

print ("Lancaster-stemming")

print (lan-stem Sentence (contents))

print ("\*\*\*\*\*")

F) def lemmSentence (sentence):

tok = tokenizer.tokenize (sentence).

filtered\_sentence = [w for w in tok if not w in stop-words]

lemm\_sentences = []

for word in filtered\_sentence:

lemm\_sentence.append (lemmatizer.lemmatize (word))

return lan (lemm\_sentence)

for file in files:

with open (file, 'r', encoding='utf-8') as f:

contents = f.readline()

print ("lemmatization")

print (lemmSentence (contents))

print ("\*\*\*\*\*")

Exercise-2

tok = []

for file in files:

with open (file, 'r', encoding='utf-8') as f:

contents = f.read()

let = tokenizer.tokenize(contents)

tok.append(let)

tok.

~~tok~~ - lem = []

for i in tok:

for j in i:

tok-lem = lemmatizer.lemmatize(j)

tok-lem.append(tok-lem)

tok-lem

Step:2

for file in files:

with open (file, 'r', encoding='utf-8') as f:

contents = f.read()

tok = tokenizer.tokenize(contents)

filtered\_sentence = [w for w in tok if not w in stop\_words]

tfidf = TfidfVectorizer(min\_df=2, max\_df=0.5,

features = tfidf.fit\_transform(filtered\_sentence)

df = pd.DataFrame(features.toarray(), columns=tfidf.get\_feature\_names())

print(df) print("\*\*\*\*\*")

## NOTES

Example:

with open('files[10]', 'r', encoding='utf-8') as f:

contents = f.read()

tok = Tokenizer.tokenize(contents)

filtered\_sentence = [w for w in tok if not w in stop\_words]

tfidf = TfidfVectorizer(min\_df=2, max\_df=0.5,

ngram\_range=(1,2))  
movie\_1 = tfidf.fit\_transform(filtered\_sentence)  
print(movie\_1)

with open('files[10]', 'r', encoding='utf-8') as f:

contents = f.read()

tok = Tokenizer.tokenize(contents)

filtered\_sentence = [w for w in tok if not w in stop\_words]

tfidf = TfidfVectorizer(min\_df=2, max\_df=0.5,

ngram\_range=(1,2))  
movie\_2 = tfidf.fit\_transform(filtered\_sentence)  
print(movie\_2)

doc1 = movie\_1[0:10]

doc2 = movie\_2[1:]

score = linear\_kernel(doc1, doc2)

print(score)