

Aim : Text Analytics

1. Extract Sample document and apply following document preprocessing methods: Tokenization, POS Tagging, stop words removal, Stemming and Lemmatization.
2. Create representation of document by calculating Term Frequency and Inverse Document Frequency.

Code :

```
In [4]: 1 import nltk
```

```
In [5]: 1 nltk.download('punkt')
```

```
[nltk_data] Downloading package punkt to
[nltk_data] C:\Users\Welcome\AppData\Roaming\nltk_data...
[nltk_data] Package punkt is already up-to-date!
```

Out[5]: True

```
In [6]: 1 nltk.download('stopwords')
```

```
[nltk_data] Downloading package stopwords to
[nltk_data] C:\Users\Welcome\AppData\Roaming\nltk_data...
[nltk_data] Unzipping corpora\stopwords.zip.
```

Out[6]: True

```
In [7]: 1 nltk.download('wordnet')
```

```
[nltk_data] Downloading package wordnet to
[nltk_data] C:\Users\Welcome\AppData\Roaming\nltk_data...
```

Out[7]: True

```
In [9]: 1 nltk.download('averaged_perceptron_tagger')
```

```
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] C:\Users\Welcome\AppData\Roaming\nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
```

Out[9]: True

```
In [10]: 1 text= "Tokenization is the first step in text analytics. The process of breaking down a text p
2 text
```

Out[10]: 'Tokenization is the first step in text analytics. The process of breaking down a text paragraph i
nto smaller chunkssuch as words or sentences is called Tokenization.'

```
In [11]: 1 from nltk.tokenize import sent_tokenize
2 tokenized_text= sent_tokenize(text)
3 print(tokenized_text)
```

```
['Tokenization is the first step in text analytics.', 'The process of breaking down a text paragra  
ph into smaller chunkssuch as words or sentences is called Tokenization.']
```

```
In [12]: 1 from nltk.tokenize import word_tokenize
2 tokenized_word=word_tokenize(text)
3 print(tokenized_word)
```

```
['Tokenization', 'is', 'the', 'first', 'step', 'in', 'text', 'analytics', '.', 'The', 'process',  
'of', 'breaking', 'down', 'a', 'text', 'paragraph', 'into', 'smaller', 'chunkssuch', 'as', 'word  
s', 'or', 'sentences', 'is', 'called', 'Tokenization', '.']
```

```
In [17]: 1 import re
2 from nltk.corpus import stopwords
3 stop_words=set(stopwords.words("english"))
4 print(stop_words)
```

{'aren', 'you're', 'isn', 'didn', 'and', 'those', 'we'll', 'they're', 'weren't', 'about', 'needn', 'as', 'they'd', 'below', 'should've', 'am', 'we', 've', 'very', 'each', 'ma', 'under', 'a', 'are n't', 'yourselves', 'herself', 'these', 'wasn', 'won't', 'didn't', 'why', 'you'd', 'don't', 'i've', 'i'll', 'i'm', 'mustn', 'own', 'they'll', 'won', 'him', 'shan', 'she', 'that'll', 'them', 'which', 'your', 'for', 'it', 'it'd', 'mustn't', 'but', 'ain', 'because', 'now', 'at', 'they've', 'some', 'she'd', 'she'll', 'after', 'ours', 'we've', 'only', 'you'll', 'of', 'myself', 'who', 'weren', 'my', 'doing', 'haven', 'he', 'me', 'over', 'shouldn't', 'just', 'i'd', 'it's', 'd', 'did', 'don', 'before', 'if', 'have', 're', 'couldn', 'our', 'is', 'there', 'it'll', 'hadn't', 'his', 'by', 'doesn't', 'how', 'same', 'we'd', 'he'll', 'o', 'no', 'do', 'most', 'was', 'you've', 'so', 'isn't', 'hasn't', 'can', 'more', 'being', 'needn't', 'out', 'than', 'were', 'during', 'its', 'above', 'are', 'you', 'both', 'with', 'not', 'been', 'further', 'here', 'in', 'their', 'to', 'where', 'down', 'mightn't', 'has', 'mightn', 'll', 'having', 'theirs', 'we're', 'couldn't', 'into', 'an', 'yourself', 's', 'themselves', 'when', 'wouldn', 'she's', 'itself', 'up', 'while', 'be', 'from', 'whom', 'will', 'shouldn', 'had', 'against', 'he'd', 'wouldn't', 'such', 'again', 'the', 'through', 'off', 'y', 'this', 'other', 'hadn', 'hers', 'they', 'yours', 'on', 'haven't', 'i', 'doesn', 'few', 'he r', 'himself', 'all', 'then', 'that', 'wasn't', 'or', 'he's', 'should', 't', 'between', 'm', 'until', 'ourselves', 'nor', 'what', 'once', 'shan't', 'too', 'hasn', 'any', 'does'}

```
In [18]: 1 text= "How to remove stop words with NLTK library inPython?"
2 text= re.sub('[^a-zA-Z]', ' ',text)
3 tokens = word_tokenize(text.lower())
4 filtered_text=[]
5 for w in tokens:
6     if w not in stop_words:
7         filtered_text.append(w)
8 print("Tokenized Sentence:",tokens)
9 print("Filterd Sentence:",filtered_text)
```

Tokenized Sentence: ['how', 'to', 'remove', 'stop', 'words', 'with', 'nltk', 'library', 'inpytho n']
Filterd Sentence: ['remove', 'stop', 'words', 'nltk', 'library', 'inpython']

```
In [20]: 1 from nltk.stem import PorterStemmer
2 e_words= ["wait", "waiting", "waited","waits"]
3 ps =PorterStemmer()
4 for w in e_words:
5     rootWord=ps.stem(w)
6 print(rootWord)
```

wait

```
In [25]: 1 from nltk.stem import WordNetLemmatizer
2 wordnet_lemmatizer = WordNetLemmatizer()
3 text = "studies studying cries cry"
4 tokenization = nltk.word_tokenize(text)
5 for w in tokenization:
6     print("Lemma for {} is {}".format(w,wordnet_lemmatizer.lemmatize(w)))
```

Lemma for studies is study
Lemma for studying is studying
Lemma for cries is cry
Lemma for cry is cry

```
In [29]: 1 import nltk
2 from nltk.tokenize import word_tokenize
3 data="The pink sweater fit herperfectly"
4 words=word_tokenize(data)
5 print(data)
6 print(words)
```

The pink sweater fit herperfectly
['The', 'pink', 'sweater', 'fit', 'herperfectly']

```
In [30]: 1 for word in words:
2         print(nltk.pos_tag([word]))
3
```

```
[('The', 'DT')]
[('pink', 'NN')]
[('sweater', 'NN')]
[('fit', 'NN')]
[('herperfectly', 'RB')]
```

```
In [31]: 1 import pandas as pd
2         from sklearn.feature_extraction.text import TfidfVectorizer
```

```
In [32]: 1 documentA = 'Jupiter is the largest Planet'
2         documentB = 'Mars is the fourth planet from the Sun'
3         print(documentA)
4         print(documentB)
```

```
Jupiter is the largest Planet
Mars is the fourth planet from the Sun
```

```
In [34]: 1 bagOfWordsA = documentA.split(' ')
2         bagOfWordsB = documentB.split(' ')
3         print(bagOfWordsA)
4         print(bagOfWordsB)
```

```
['Jupiter', 'is', 'the', 'largest', 'Planet']
['Mars', 'is', 'the', 'fourth', 'planet', 'from', 'the', 'Sun']
```

```
In [36]: 1 uniqueWords = set(bagOfWordsA).union(set(bagOfWordsB))
```

```
In [37]: 1 numOfWordsA = dict.fromkeys(uniqueWords, 0)
2         for word in bagOfWordsA:
3             numOfWordsA[word] += 1
4         numOfWordsB = dict.fromkeys(uniqueWords, 0)
5         for word in bagOfWordsB:
6             numOfWordsB[word] += 1
```

```
In [47]: 1 def computeTF(wordDict, bagOfWords):
2         tfDict = {}
3         bagOfWordsCount = len(bagOfWords)
4         for word, count in wordDict.items():
5             tfDict[word] = count / float(bagOfWordsCount)
6         return tfDict
7         tfA = computeTF(numOfWordsA, bagOfWordsA)
8         tfB = computeTF(numOfWordsB, bagOfWordsB)
9
```

```
In [48]: 1 def computeIDF(documents):
2         import math
3         N = len(documents)
4         idfDict = dict.fromkeys(documents[0].keys(), 0)
5         for document in documents:
6             for word, val in document.items():
7                 if val > 0:
8                     idfDict[word] += 1
9         for word, val in idfDict.items():
10            idfDict[word] = math.log(N / float(val))
11         return idfDict
```

```
In [49]: 1 idfs = computeIDF([numOfWordsA, numOfWordsB])
2 idfs
```

```
Out[49]: {'Jupiter': 0.6931471805599453,
'planet': 0.6931471805599453,
'largest': 0.6931471805599453,
'is': 0.0,
'from': 0.6931471805599453,
'fourth': 0.6931471805599453,
'Planet': 0.6931471805599453,
'Mars': 0.6931471805599453,
'Sun': 0.6931471805599453,
'the': 0.0}
```

```
In [50]: 1 def computeTFIDF(tfBagOfWords, idfs):
2     tfidf = {}
3     for word, val in tfBagOfWords.items():
4         tfidf[word] = val * idfs[word]
5     return tfidf
6 tfidfA = computeTFIDF(tfA,idfs)
7 tfidfB = computeTFIDF(tfB,idfs)
8 df = pd.DataFrame([tfidfA,tfidfB])
9 df
10
```

```
Out[50]:
```

	Jupiter	planet	largest	is	from	fourth	Planet	Mars	Sun	the
0	0.138629	0.000000	0.138629	0.0	0.000000	0.000000	0.138629	0.000000	0.000000	0.0
1	0.000000	0.086643	0.000000	0.0	0.086643	0.086643	0.000000	0.086643	0.086643	0.0

Name : Mohan Kadambande

Roll No. : 13212 (TECO-b1)

```
In [ ]: 1
```