

# FIT5196 Assignment 1

## Task 1

Student name: Ayush Sharma

Due date: 13 Sept 2020

Environment: Python 3.6.0 and Anaconda 4.3.0 (64-bit)

## Introduction

### Steps followed:

- 1.) Import xlsx data into dataframe and Remove empty rows/columns and Remove non-english tweets using langid Combine all tweets for a sheet.
- 2.) Tokenize combined tweet using RegexpTokenizer
- 3.) Generate top 100 bigrams for each sheet using BigramCollocationFinder and Write top 100 bigrams, count for each date to bigram.txt
- 4.) Remove stop words using stopwords file, Remove words appearing in more than 60 sheets and less than 5 sheets, Remove words with length less than 3
- 5.) Stemming all the tokens
- 6.) Get top 100 unigrams for each sheet and count and Write top 100 unigrams, count for each date to unigram.txt
- 7.) Generate top 200 bigrams on complete text from the original tokens
- 8.) Get vocab by taking unique words from all stemmed unigrams Add 200 bigrams to the vocab appending them with underscore Write vocab to vocab.txt
- 9.) Create sparse matrix by with vocab ids of words in for each date and their count and Write this to countVec.txt

### Libraries Used:

- re: This Library is used to tokenize the
- langid: "langid" package, only keeps the tweets that are in English language
- pandas as pd: Load, explore, work on the dataframe that is loaded
- nltk: NLTK is a Natural language text processing library used to perform text analysis
- nltk.tokenize , RegexpTokenizer: This library tokenise the regular expression passed into It.
- from nltk.collocations import \* used to create the bigrams
- from nltk.stem import PorterStemmer: This library is used form stemming
- from itertools import chain: Helps in combing various Lists
- from nltk import bigrams: Generate bigram

- from nltk import FreqDist: gives frequency of the bigram or unigram

## 1.) Importing Libraries

In [1]:



```
1 import re
2 import langid
3 import pandas as pd
4 import nltk
5 from nltk.tokenize import RegexpTokenizer
6 from nltk import bigrams
7 from nltk.util import ngrams
8 from nltk.collocations import *
9 from nltk.stem import PorterStemmer
10 from itertools import chain
11 from nltk import bigrams
12 from nltk import FreqDist
13 import itertools
14 from nltk.collocations import *
15 from collections import Counter
```

## Step 2: Reading and Analysing data

The data contains of 80 sheets of excel that have thousands of tweet text and date along with id

- Read the data
- parsing through the excel files and sheets and removing unwanted NA's and reshaping it
- Applying langid on each text if not english the text will not be appended to the dictionary
- finally creating a dictionary that has date as key and text as a list of values

In [2]:



```

1 # Reading the Excel data
2 data=pd.ExcelFile('30823293.xlsx')
3 # creating blank dict for sheets
4 tweetDict={}
5 # Iterating through each sheet and re shaping it into a better and more usefull format
6 for i in range(len(data.sheet_names)):
7     # parsing through each sheet
8     df=data.parse(i)
9     # removing na values
10    df=df.dropna(0, how = 'all')
11    df=df.dropna(1, how = 'all')
12    # Seting columns names
13    df.columns = ["Text","Id","Date"]
14    # resetting indexes
15    df.index = range(len(df.index))
16    df=df.drop([0])
17    df.index = range(len(df.index))
18    # iterating over the Text column to perform tasks like checking if the type is str
19    # Dropping all those values that are not required
20    for j in df['Text']:
21        if type(j)!=str:
22            nonstrindex=df[df['Text']==j].index
23            df.drop(nonstrindex,inplace=True)
24            df.reset_index(drop = True)
25    # Using the Langid checking each of the Text to remove all non english tweet
26    if langid.classify(str(j))[0]!='en':
27        textindex=df[df['Text']==j].index
28        df.drop(textindex,inplace=True)
29        df.reset_index(drop = True)
30    # Creating a Dictionary with date as key and tweets texts as value
31    tweetDict[data.sheet_names[i]] = '\n'.join(df.Text.values).lower()

```

## Step 3: Tokenizing

- Tokenization refers to creating each words in the text as token by applying the regex `[a-zA-Z]+(?:[- '][a-zA-Z]+)?`
- Tokenizing the data of text and applying regex to create the token
- applying tokenization for to the dictionary that has all the text

In [3]:



```

1 # Creating tokens
2 tokenizer= RegexpTokenizer("[a-zA-Z]+(?:[- '][a-zA-Z]+)?")
3 unigram_token={}
4 # Applying token to the dictionary data
5 for i,j in tweetDict.items():
6     unigram_token[i]=tokenizer.tokenize(j)
7
8

```

## Step 4 Generating 100 Bigram:

- Using the Counter function to calculate the top 100 bigrams by using the package nltk.bigrams that finds each bigram for the text
- Then structuring the bigrams that we have in a way that includes the count of the bigram and formatting according to output file
- Finally having a structured format of a dictionary with key as date and bigrams as a list of values with bigrams within a tuple and a tuple that has the count of each bigram.
- Writing a text file that have the

In [4]:

```
1 # Finding the top 100 bigrams words and counting their occurrences
2 bigramdict={}
3 for i,j in unigram_token.items():
4     bigramword=list(nltk.bigrams(j))
5     bigramcount=Counter(bigramword).most_common(100)
6     bigramdict[i]=bigramcount
```

In [5]:

```
1 # Formatting the bigrams extracted in the required way to write the file
2 final_bi_file_str=str()
3 # iterating within the dict to create the text
4 for i,j in bigramdict.items():
5     final_bi_file_str=final_bi_file_str+str(i)+':'+str(j)+'\n'
```

In [11]:

```
1 # creating and writing the 100bi.txt
2 bigramtxtfile = open("30823293_100bi.txt", "w",encoding='UTF-8')
3 bigramtxtfile.write(str(final_bi_file_str))
4
5 bigramtxtfile.close()
```

## Step 5 Removing stopwords

- Stop words are the words least significant words that are appearing frequently like a, an, the, is etc
- Stop words are removed using the text file provided by the faculty that will help us remove such words.
- we will iterate over our token dictionary by referring to each key which in our case is the date and look for each element within the value and selecting the element of the value to check if it is a stopword or not if it is a stop word then it will be removed from the data.
- after that we have calculated the document frequency of the word by using `FreqDist(list(chain.from_iterable))` to find if a word is in more than 60 documents or if less than 5 documents after finding them these words are removed from the dictionary
- now appending the new dictionary after handling all the removals

In [12]:

```
1 # opening the stopwords file provided to us by faculty
2 stopwords=[]
3 with open('stopwords_en.txt') as f:
4     stopwords=f.read().splitlines()
```

In [13]:



```

1 # filtering and removing the stopwords
2 tokenised_dict={}
3 for i,j in unigram_token.items():
4     tokenised_dict[i]= [x for x in j if x not in stopwords]
5 # Filtering the words with len less than 3 characters within the dictionary
6 for i,j in tokenised_dict.items():
7     tokenised_dict[i]= [x for x in j if len(x)>3]

```

In [14]:



```

1 # removing less favorable word i.e words with < 5 count in the document frequency
2 #and words with >60 in frequency and creating new dict
3 list_final={}
4 for i,j in tokenised_dict.items():
5     list_final = FreqDist(list(chain.from_iterable([set(x) for x in tokenised_dict.values() if i in x])))
6
7 list_remove=[]
8 for i in list_final.keys():
9     if list_final.get(i)>60 or list_final.get(i)<5:
10         list_remove.append(i)
11
12 for i in list_remove:
13     list_final.pop(i)

```

In [15]:



```

1 # appending the above results to new dict
2 tokenised_dict2 = {}
3 for i in tokenised_dict:
4     lst=tokenised_dict.get(i)
5     refined_list=[]
6     for each in lst:
7         if(each in list_final):
8             refined_list.append(each)
9     tokenised_dict2[i]=refined_list

```

## Step 6 Stemming

- We have used PorterStemmer for stemming because we were asked to use it.
- (or 'Porter stemmer') is a process for removing the commoner morphological and inflexional endings from words in English.
- The Porter stemming algorithm removes the commoner morphological and inflexional endings from words in English. like consider, considering, speed, speeding, sleep, sleeping
- remove the words including suffix like ing or e etc thus such words will be converted to similar words after stemming.

In [16]:

```
1 # Stemming the tokens using porterstemmer and passing it into a dictionary
2 tokenised_dict3 = {}
3 ps = PorterStemmer()
4 for i in tokenised_dict2:
5     lst=tokenised_dict2.get(i)
6     refined_list=[]
7     for each in lst:
8         refined_list.append(ps.stem(each))
9     tokenised_dict3[i]=refined_list
```

## Step 6 Generating top 100 Unigrams:

- In this step we have all the umograms after removing unwanted characters and words and stemming the tokens
- now we will count the frequency of all the token for each date
- Finally we will collect and keep the top 100 unigrams that we need to keep
- In the final step we will format the top 100 unigram as per the required output and write a text file for it.

In [17]:

```
1 # Extracted the freq of each remaining tokens after stemming and removing stopwords
2 #and selecting the top 100 unigrams and storing it
3 unigram_dictCount={}
4 for i,j in tokenised_dict3.items():
5     unigram_dictCount[i]=FreqDist((j)).most_common(100)
6
```

In [18]:

```
1 # Formating the unigrams according to the required output to be stored
2 final_uni_file_str=str()
3 for i,j in unigram_dictCount.items():
4     final_uni_file_str=final_uni_file_str+str(i)+':'+str(j)+'\n'
```

In [19]:

```
1 # writting the 100uni.txt files
2 unigramtxtfile = open("30823293_100uni.txt", "w",encoding='UTF-8')
3 unigramtxtfile.write(str(final_uni_file_str))
4 unigramtxtfile.close()
```

## Step 7 Generating top 200 Bigrams:

- with the help of `nltk.collocations.BigramAssocMeasures()` and `nltk.collocations.BigramCollocationFinder.from_words` we will apply this to get the top 200 bigrams using pmi method

In [20]:

```
1 # Generating top 200 bigrams from all tokens within data
2 allbigram=[]
3 bigram_measure=nltk.collocations.BigramAssocMeasures()
4
5 #for i,j in unigaram_token.items():
6 words_bi = list(chain.from_iterable(unigaram_token.values()))
7 finder=nltk.collocations.BigramCollocationFinder.from_words(words_bi)
8 allbigram.append(finder.nbest(bigram_measure.pmi,200))
9
```

## Step 7 Creating Vocab of Words

- now we have 200 bigrams and unigrams in we will format them to keep them in the list and create a formatted list
- we will iterate on the final list and sort it and then put it in a string that will make it formatted as per the requirement
- finally we will write the vocab.txt

In [21]:

```
1 # creating a list of bigrams with each bigram appened to each other as a word for vocab
2 bigramList=[]
3 for i in allbigram:
4     for j in i:
5         if j not in bigramList:
6             bigramList.append(j[0]+'_'+j[1])
7
```

In [22]:

```
1 #Creating a list of unigrams
2 vocabList=[]
3 for i,j in tokenised_dict3.items():
4     for x in j:
5         vocabList.append(x)
6
7
```

In [24]:

```
1 Vocab = []
2 finalVocab = []
3 finalVocab=vocabList+bigramList
4 Vocab=set(finalVocab)
5 Vocab_list = list(Vocab)
6 Vocab_list = sorted(Vocab_list)
```

In [25]:



```
1 string_vocab = ''
2 for i in range(len(Vocab_list)):
3     string_vocab = string_vocab + str(Vocab_list[i]) + ':' + str(i) + '\n'
4
5
6 # writting the vocab.txt files
7 vocabtxtfile = open("30823293_vocab.txt", "w",encoding='UTF-8')
8 vocabtxtfile.write(str(string_vocab))
9 vocabtxtfile.close()
```

In [26]:



```
1 # writting the vocab.txt files
2 vocabtxtfile = open("30823293_vocab.txt", "w",encoding='UTF-8')
3 vocabtxtfile.write(str(string_vocab))
4 vocabtxtfile.close()
```

## Conclusion:

In this task I have created tokens for each date. Then I have carried out steps which includes removal of stopwords, context independent words, context dependent words and words with length less than 3. After this I have carried out stemming for each token. Hence I have created unigrams and vocab files. Also I have generated bigrams which occur the most frequently using collocations and pmi measure techniques. Hence the following files have been generated.

Hence after this task I have gained knowledge on how to work with nltk library and I have understood the concepts of unigrams and bigrams

In [ ]:



1