NAME : G.SRIKAR

**REGISTRATION NO.:** 19BCE7081

**SUBJECT** : Advanced Data Analytics

(CSE-4029)

**SLOT** : L21 + L22

\*\* Lab Assignment-4 \*\*

# **Twitter Data Analysis**

**<u>AIM</u>**: Ukraine Russia War Twitter data Sentiment Analysis using Python

## Reading & Understanding the Data:

### **PYTHON-CODE**:-

import pandas as pd

import seaborn as sns

import matplotlib.pyplot as plt

 $from\ nltk. sentiment. vader\ import\ SentimentIntensity Analyzer$ 

 $from\ wordcloud\ import\ WordCloud,\ STOPWORDS,\ ImageColorGenerator$ 

import nltk

import re

from nltk.corpus import stopwords

import string

data = pd.read\_csv("filename.csv")

print(data.head())

```
warnings.warn("The twython library has not been installed.
0 1508562464916582400 1508532968867909640 2022-03-28 21:51:34 UTC
  1508562460445454338 1508562460445454338 2022-03-28 21:51:32 UTC 1508562459476561926 1508460376924557316 2022-03-28 21:51:32 UTC 1508562459149451265 1508460376924557316 2022-03-28 21:51:32 UTC
   2022-03-28 21:51:33
2022-03-28 21:51:32
   2022-03-28 21:51:32
                                                                           canadageneva
                                      name place ... geo source user_rt_id user_rt \
oelho NaN ... NaN NaN NaN NaN
                         Gianluca Rubeo NaN ... NaN
                                                                    NaN
                                                                                   NaN
                                                                                              NaN
  MIKA UL YARMYAOHUH HAOLULYAOHUH
                    Canada in Geneva 🍁
                                                                    NaN
                                                                                              NaN
           NaN [{'screen name': 'CarlosPortSPFC', 'name': 'Ca...
           NaN
           NaN
                [{'screen_name': 'Amb_Ulyanov', 'name': 'Mikha...
                                                                                                NaN
```

```
translate trans_src trans_dest

0  NaN  NaN  NaN

1  NaN  NaN  NaN

2  NaN  NaN  NaN

3  NaN  NaN  NaN

4  NaN  NaN  NaN

[5 rows x 36 columns]
```

## Printing all the column names of the dataset:

**PYTHON-CODE**:-

print(data.columns)

**OUTPUT**:-

Need three columns for this task (username, tweet, and language); I will only select these columns and move forward:

```
↓ PYTHON-CODE:
data = data[["username", "tweet", "language"]]
```

Cheking whether any of these columns contains any null values or not:

**PYTHON-CODE**:

data.isnull().sum()

**OUTPUT**:-

```
data.isnull().sum()

username    0
tweet     0
language    0
dtype: int64
```

So none of the columns has null values, checking for how many tweets are posted in each language:

```
PYTHON-CODE:-
```

data["language"].value\_counts()

**OUTPUT**:-

```
data["language"].value counts()
       8761
en
it
        382
        324
pt
       194
und
         53
ru
         39
in
fr
         38
es
         31
de
         19
         19
ca
ja
         18
         16
tr
nl
         15
         13
pl
         13
ar
          10
hi
CS
tl
           6
et
           6
uk
fi
           5
           4
су
ro
zh
```

Most of the tweets are in English. Preparing this data for the task of sentiment analysis. Here I will remove all the links, punctuation, symbols and other language errors from the tweets:

```
# PYTHON-CODE:
nltk.download('stopwords')
stemmer = nltk.SnowballStemmer("english")
stopword=set(stopwords.words('english'))

def clean(text):
    text = str(text).lower()
    text = re.sub('\[.*?\]', ", text)
    text = re.sub('https?://\S+|www\.\\S+', ", text)
    text = re.sub('<.*?>+', ", text)
    text = re.sub('\[.*?\]' % re.escape(string.punctuation), ", text)
    text = re.sub('\[.*]', ", text)
    text = re.sub('\[.*]', ", text)
    text = re.sub('\[.*]', ", text)
    text = [word for word in text.split(' ') if word not in stopword]
    text=" ".join(text)
    text = [stemmer.stem(word) for word in text.split(' ')]
```

```
text=" ".join(text)
return text
data["tweet"] = data["tweet"].apply(clean)
# OUTPUT:-
```

```
nltk.download('stopwords')
stemmer = nltk.SnowballStemmer("english")
stopword=set(stopwords.words('english'))

def clean(text):
    text = str(text).lower()
    text = re.sub('\[.*?\]', '', text)
    text = re.sub('https?://\S+|www\.\S+', '', text)
    text = re.sub('\[.*?\]', '', text)
    text = re.sub('\[.*?\]', '', text)
    text = re.sub('\[.*?\]', '', text)
    text = re.sub('\[.*]\]', '', text)
    text = re.sub('\[.*]\]', '', text)
    text = re.sub('\[.*]\]', '', text)
    text = [word for word in text.split(' ') if word not in stopword]
    text=" ".join(text)
    text = [stemmer.stem(word) for word in text.split(' ')]
    text=" ".join(text)
    return text
data["tweet"] = data["tweet"].apply(clean)

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data] Package stopwords is already up-to-date!
```

looking at the wordcloud of the tweets, which will show the most frequently used words in the tweets by people sharing their feelings and updates about the Ukraine and Russia war.

```
text = " ".join(i for i in data.tweet)
stopwords = set(STOPWORDS)
wordcloud = WordCloud(stopwords=stopwords,
background_color="white").generate(text)
plt.figure( figsize=(15,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```

## **PYTHON-CODE**:

```
nltk.download('vader_lexicon')
sentiments = SentimentIntensityAnalyzer()
data["Positive"] = [sentiments.polarity_scores(i)["pos"] for i in data["tweet"]]
data["Negative"] = [sentiments.polarity_scores(i)["neg"] for i in data["tweet"]]
data["Neutral"] = [sentiments.polarity_scores(i)["neu"] for i in data["tweet"]]
data = data[["tweet", "Positive", "Negative", "Neutral"]]
print(data.head())
```

# Looking at the most frequent words used by people with positive sentiments:

## **PYTHON-CODE**:

```
positive =' '.join([i for i in data['tweet'][data['Positive'] > data["Negative"]]])
stopwords = set(STOPWORDS)
wordcloud = WordCloud(stopwords=stopwords,
background_color="white").generate(positive)
plt.figure( figsize=(15,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```

```
positive = ' '.join([i for i in data['tweet'] [data['Positive'] > data["Negative"]]])

stopwords = set(STOFMORDS)

wordcloud = Wordcloud(stopwords=stopwords, background_color="white").generate(positive)

plt.figure(figsize=(15,10))

plt.mishow(wordcloud, interpolation='bilinear')

plt.axis("off")

plt.show()

wordcloud;

plt.figure(figsize=(15,10))

plt.mishow(wordcloud, interpolation='bilinear')

plt.show()

wordcloud;

plt.show()

mani year

man
```

# Looking at the most frequent words used by people with negative sentiments:

## **PYTHON-CODE**:

```
negative =' '.join([i for i in data['tweet'][data['Negative'] > data["Positive"]]])
stopwords = set(STOPWORDS)
wordcloud = WordCloud(stopwords=stopwords,
background_color="white").generate(negative)
plt.figure( figsize=(15,10))
plt.imshow(wordcloud, interpolation='bilinear')
plt.axis("off")
plt.show()
```