# A special thanks to all my mentors for helping me constantly to progress technically

# Jupyter notebook prepared, arranged and executed by Karthi Balasundaram , sentimentally analysing Russian Ukraine War using real tweet data from twitter.

In [3]:
```python
# installing natural language toolkit(nltk)
!pip install nltk
```

```
Defaulting to user installation because normal site-packages is not writeable
Collecting nltk
  Downloading nltk-3.7-py3-none-any.whl (1.5 MB)
     |████████████████████████████████| 1.5 MB 2.4 MB/s
Requirement already satisfied: click in ./Library/Python/3.9/lib/python/site-p
ackages (from nltk) (8.0.3)
Collecting regex>=2021.8.3
  Downloading regex-2022.3.15-cp39-cp39-macosx_10_9_x86_64.whl (288 kB)
     |████████████████████████████████| 288 kB 2.6 MB/s
Requirement already satisfied: tqdm in /Library/Frameworks/Python.framework/Ve
rsions/3.9/lib/python3.9/site-packages (from nltk) (4.62.3)
Requirement already satisfied: joblib in /Library/Frameworks/Python.framework/
Versions/3.9/lib/python3.9/site-packages (from nltk) (1.0.1)
Installing collected packages: regex, nltk
  WARNING: The script nltk is installed in '/Users/karthibalasundaram/Library/
Python/3.9/bin' which is not on PATH.
  Consider adding this directory to PATH or, if you prefer to suppress this wa
rning, use --no-warn-script-location.
Successfully installed nltk-3.7 regex-2022.3.15
WARNING: You are using pip version 21.3.1; however, version 22.0.4 is availabl
e.
You should consider upgrading via the '/Library/Frameworks/Python.framework/Ve
rsions/3.9/bin/python3.9 -m pip install --upgrade pip' command.
```

In [9]:
```python
# installing openpyxl (a python library to read/write excel files)
!pip install openpyxl
```

```
Defaulting to user installation because normal site-packages is not writeable
Collecting openpyxl
  Downloading openpyxl-3.0.9-py2.py3-none-any.whl (242 kB)
     |████████████████████████████████| 242 kB 1.7 MB/s
Collecting et-xmlfile
  Downloading et_xmlfile-1.1.0-py3-none-any.whl (4.7 kB)
Installing collected packages: et-xmlfile, openpyxl
Successfully installed et-xmlfile-1.1.0 openpyxl-3.0.9
WARNING: You are using pip version 21.3.1; however, version 22.0.4 is availabl
e.
You should consider upgrading via the '/Library/Frameworks/Python.framework/Ve
rsions/3.9/bin/python3.9 -m pip install --upgrade pip' command.
```

In [4]:
```python
#importing other default and necessary libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import string
import re
import nltk
from nltk.util import pr
from nltk.corpus import stopwords
import warnings
warnings.filterwarnings('ignore')
stemmer = nltk.SnowballStemmer("english")
nltk.download('stopwords')
stopword=set(stopwords.words('english'))
```

```
[nltk_data] Downloading package stopwords to
[nltk_data]     /Users/karthibalasundaram/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
```

In [111…
```python
#reading the excel file using pandas library
data = pd.read_excel("/Users/karthibalasundaram/Downloads/Russia_Ukraine_war/
```

In [112…
```python
#the below line calls last 5 rows from the excel
data.tail()
```

Out[112…

| | id | conversation_id | created_at | date | time | timezone | user_id |
|---|---|---|---|---|---|---|---|
| **10009** | 1.504308e+18 | 1.503516e+18 | 2022-03-17 04:06:25 UTC | 2022-03-17 | 04:06:25 | 0.0 | 1.486028e+18 |
| **10010** | 1.504308e+18 | 1.504308e+18 | 2022-03-17 04:06:24 UTC | 2022-03-17 | 04:06:24 | 0.0 | 1.504306e+18 |
| **10011** | 1.504308e+18 | 1.486862e+18 | 2022-03-17 04:06:24 UTC | 2022-03-17 | 04:06:24 | 0.0 | 1.470945e+18 |
| **10012** | 1.504308e+18 | 1.504289e+18 | 2022-03-17 04:06:24 UTC | 2022-03-17 | 04:06:24 | 0.0 | 1.239372e+18 |
| **10013** | 1.504308e+18 | 1.504111e+18 | 2022-03-17 04:06:23 UTC | 2022-03-17 | 04:06:23 | 0.0 | 1.464508e+18 |

5 rows × 36 columns

In [113…

```
#the below line calls first 5 rows from the excel
data.head()
```

Out[113…

| | id | conversation_id | created_at | date | time | timezone | user_id | |
|---|---|---|---|---|---|---|---|---|
| **0** | 1.504326e+18 | 1.504083e+18 | 2022-03-17 05:15:51 UTC | 2022-03-17 | 05:15:51 | 0.0 | 1.016938e+09 | bowti |
| **1** | 1.504326e+18 | 1.504323e+18 | 2022-03-17 05:15:51 UTC | 2022-03-17 | 05:15:51 | 0.0 | 1.420232e+18 | the |
| **2** | 1.504326e+18 | 1.504326e+18 | 2022-03-17 05:15:51 UTC | 2022-03-17 | 05:15:51 | 0.0 | 1.387731e+18 | rosaor |
| **3** | 1.504326e+18 | 1.504326e+18 | 2022-03-17 05:15:50 UTC | 2022-03-17 | 05:15:50 | 0.0 | 5.421008e+07 | |
| **4** | 1.504326e+18 | 1.504325e+18 | 2022-03-17 05:15:50 UTC | 2022-03-17 | 05:15:50 | 0.0 | 6.432839e+07 | ar |

5 rows × 36 columns

In [114…

```
#understanding rows and columns present in the excel
data.shape
```

Out[114…

```
(10014, 36)
```

In [115…

```
#retreives basic info about the excel data
data.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10014 entries, 0 to 10013
Data columns (total 36 columns):
 #   Column            Non-Null Count  Dtype
---  ------            --------------  -----
 0   id                10014 non-null  float64
 1   conversation_id   10014 non-null  float64
 2   created_at        10014 non-null  object
 3   date              10014 non-null  datetime64[ns]
 4   time              10014 non-null  object
 5   timezone          10014 non-null  float64
 6   user_id           10014 non-null  float64
 7   username          10014 non-null  object
 8   name              10014 non-null  object
 9   place             1 non-null      object
 10  tweet             10014 non-null  object
 11  language          10014 non-null  object
 12  mentions          10014 non-null  object
 13  urls              10014 non-null  object
 14  photos            10014 non-null  object
 15  replies_count     10014 non-null  float64
 16  retweets_count    10014 non-null  float64
 17  likes_count       10014 non-null  float64
 18  hashtags          10014 non-null  object
 19  cashtags          10014 non-null  object
 20  link              10014 non-null  object
 21  retweet           10014 non-null  bool
 22  quote_url         876 non-null    object
 23  video             10014 non-null  float64
 24  thumbnail         936 non-null    object
 25  near              0 non-null      float64
 26  geo               0 non-null      float64
 27  source            0 non-null      float64
 28  user_rt_id        0 non-null      float64
 29  user_rt           0 non-null      float64
 30  retweet_id        0 non-null      float64
 31  reply_to          10014 non-null  object
 32  retweet_date      0 non-null      float64
 33  translate         0 non-null      float64
 34  trans_src         0 non-null      float64
 35  trans_dest        0 non-null      float64
dtypes: bool(1), datetime64[ns](1), float64(18), object(16)
memory usage: 2.7+ MB
```

In [116…

```python
#a brief description about the data
data.describe()
```

Out[116…

| | id | conversation_id | timezone | user_id | replies_count | retweets_count |
|---|---|---|---|---|---|---|
| **count** | 1.001400e+04 | 1.001400e+04 | 10014.0 | 1.001400e+04 | 10014.000000 | 10014.000000 |
| **mean** | 1.504317e+18 | 1.502877e+18 | 0.0 | 6.984499e+17 | 0.313661 | 0.552227 |
| **std** | 5.075717e+12 | 2.728863e+16 | 0.0 | 6.443610e+17 | 2.549457 | 10.848945 |
| **min** | 1.504308e+18 | 4.371802e+17 | 0.0 | 7.421430e+05 | 0.000000 | 0.000000 |
| **25%** | 1.504312e+18 | 1.504181e+18 | 0.0 | 4.921743e+08 | 0.000000 | 0.000000 |
| **50%** | 1.504317e+18 | 1.504309e+18 | 0.0 | 8.388104e+17 | 0.000000 | 0.000000 |
| **75%** | 1.504321e+18 | 1.504316e+18 | 0.0 | 1.354872e+18 | 0.000000 | 0.000000 |
| **max** | 1.504326e+18 | 1.504326e+18 | 0.0 | 1.504322e+18 | 142.000000 | 666.000000 |

In [117…

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

```
Out[117…    id                    0
            conversation_id       0
            created_at            0
            date                  0
            time                  0
            timezone              0
            user_id               0
            username              0
            name                  0
            place             10013
            tweet                 0
            language              0
            mentions              0
            urls                  0
            photos                0
            replies_count         0
            retweets_count        0
            likes_count           0
            hashtags              0
            cashtags              0
            link                  0
            retweet               0
            quote_url          9138
            video                 0
            thumbnail          9078
            near              10014
            geo               10014
            source            10014
            user_rt_id        10014
            user_rt           10014
            retweet_id        10014
            reply_to              0
            retweet_date      10014
            translate         10014
            trans_src         10014
            trans_dest        10014
            dtype: int64
```

In [118…

```python
#retreives all the columns
data.columns
```

Out[118…
```
Index(['id', 'conversation_id', 'created_at', 'date', 'time', 'timezone',
       'user_id', 'username', 'name', 'place', 'tweet', 'language', 'mentions'
,
       'urls', 'photos', 'replies_count', 'retweets_count', 'likes_count',
       'hashtags', 'cashtags', 'link', 'retweet', 'quote_url', 'video',
       'thumbnail', 'near', 'geo', 'source', 'user_rt_id', 'user_rt',
       'retweet_id', 'reply_to', 'retweet_date', 'translate', 'trans_src',
       'trans_dest'],
      dtype='object')
```

In [119…
```python
#lists first 5 data(tweets) listed under the column "tweet"
data[["tweet"]].head()
```

Out[119…

| | tweet |
|---|---|
| 0 | @PeterSchiff @PadaPrabu @SteveKrohn1 If it wer... |
| 1 | @meatballsubzero Are you pro russia or pro Ukr... |
| 2 | @SUBWAY Please stop doing business in Russia.... |
| 3 | Is Russia prepared for an economic crisis? Dev... |
| 4 | @BW Putin is Fake News 🔴 The Ruble is trash... |

In [120…
```python
#lists first 5 data(username) listed under the column "username"
data[["username"]].head()
```

Out[120…

| | username |
|---|---|
| 0 | bowtiedbeyonce |
| 1 | theshydoomer |
| 2 | rosaort91373426 |
| 3 | woodsallan |
| 4 | artemistweet |

In [121…
```python
#lists first 5 data(langauge) listed under the column "language"
data[["language"]].head()
```

Out[121…

| | language |
|---|---|
| 0 | en |
| 1 | en |
| 2 | en |
| 3 | en |
| 4 | en |

In [122…
```python
#displays the tweets posted in corresponding languages
data["language"].value_counts()
```

```
Out[122...   en    9018
             pt     211
             und    158
             it     118
             hi      80
             in      79
             ru      69
             ja      54
             es      22
             pl      19
             tl      18
             nl      15
             de      14
             ar      13
             fr      13
             zh      11
             th      10
             ca       9
             ta       8
             ro       6
             et       6
             bn       5
             fi       5
             mr       5
             ne       5
             or       5
             uk       4
             kn       4
             cs       4
             ml       4
             te       3
             el       3
             ur       3
             no       3
             gu       3
             tr       2
             iw       2
             sl       1
             am       1
             fa       1
             Name: language, dtype: int64
```

In [123...

```python
#lists first 5 data(URL's) listed under the column "link"
data[["link"]].head()
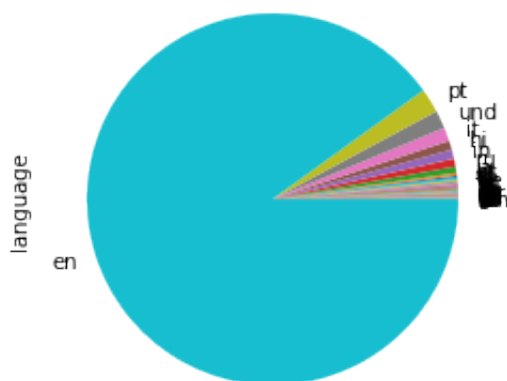```

Out[123...                                                  **link**

0        https://twitter.com/bowtiedbeyonce/status/1504...

1        https://twitter.com/TheShyDoomer/status/150432...

2        https://twitter.com/RosaOrt91373426/status/150...

3        https://twitter.com/WoodsAllan/status/15043256...

4        https://twitter.com/ArtemisTweet/status/150432...

In [124...
```python
#sorting the languages
pi = data.language.value_counts().sort_values()
```

In [125...
```python
#displaying the sorted lanuages in a pie chart
displ = pi.plot(kind = 'pie')
```
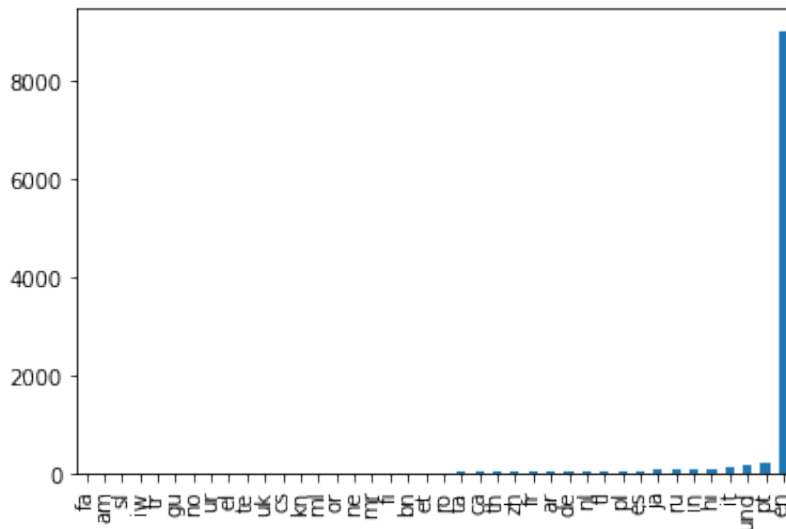


In [126...
```python
#displaying the sorted lanuages in a bar chart
displ1 = pi.plot(kind = 'bar')
```

In [127…

```
#displays the 369th tweet
data["tweet"][369]
```

Out[127…
```
'@MarshaHairbrush @RusEmbJakarta @mfa_russia @natomission_ru @NATO @Kemlu_RI @
RusEmbUSA @RusEmbIndia @EmbassyofRussia  https://t.co/dSosWvqgMo'
```

In [128…

```
# defining function for twitter hashtag extraction to classify sentiment anal
def hashtag_extract(text_list):
    hashtags = []
    for text in text_list:
        ht = re.findall(r"#(\w+)", text)
        hashtags.append(ht)
    return hashtags
```

In [133…

```
#importing seaborn library
import seaborn as sns
```

In [134…

```
# defining function for generating frequent hashtag used
def generate_hashtag_freqdist(hashtags):
    a = nltk.FreqDist(hashtags)
    b = pd.DataFrame({'Hashtag': list(a.keys()),'Count': list(a.values())})
    # selecting top 15 most frequent hashtags
    b = b.nlargest(columns="Count", n = 25)
    plt.figure(figsize=(16,7))
    ax = sns.barplot(data=b, x= "Hashtag", y = "Count")
    plt.xticks(rotation=80)
    ax.set(ylabel = 'Count')
    plt.show()
```

In [135...
```python
hashtags = hashtag_extract(data["tweet"])
hashtags = sum(hashtags, [])
```

In [136...
```python
#frequently used hastags are displayed using seaborn library
generate_hashtag_freqdist(hashtags)
```



In [188...
```python
data['total_length_characters'] = data['tweet'].str.len()
print(data['total_length_characters'])
total_length_characters = data['total_length_characters'].sum()
print(total_length_characters)
count = 0
for y in data["tweet"]:
    count = count + 1
print(count)
average_length = total_length_characters / count
print (average_length)
```

```
0          130
1          162
2          167
3          220
4           87
         ...
10009      255
10010       84
10011      176
10012      249
10013      216
Name: total_length_characters, Length: 10014, dtype: int64
1831052
10014
182.84921110445376
```

In [189...

```python
data['total_count_words'] = data['tweet'].str.split().str.len()
print(data['total_count_words'])
total_words = data['total_count_words'].sum()
print(total_words)
count = 0
for y in data["tweet"]:
    count = count + 1
print(count)
average_words = total_words / count
print (average_words)
```

```
0           22
1           28
2           26
3           32
4           15
          ..
10009       44
10010       11
10011       32
10012       39
10013       32
Name: total_count_words, Length: 10014, dtype: int64
271703
10014
27.13231475933693
```

In [190…
```python
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('[%s]' % re.escape(string.punctuation), '', text)
    text = re.sub('\n', '', text)
    text = re.sub('\w*\d\w*', '', 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)
```

In [191…
```python
data['total_length_characters'] = data['tweet'].str.len()
print(data['total_length_characters'])
total_length_characters = data['total_length_characters'].sum()
print("Total_length :",total_length_characters)
count = 0
for y in data["tweet"]:
    count = count + 1
print("Total_rows :",count)
average_length = total_length_characters / count
print ("Averge length :",average_length)
```

```
0           64
1           98
2          121
3          134
4           59
         ...
10009      126
10010       74
10011      115
10012      130
10013      137
Name: total_length_characters, Length: 10014, dtype: int64
Total_length : 1142035
Total_rows : 10014
Averge length : 114.04383862592371
```

In [192…
```python
data['total_count_words'] = data['tweet'].str.split().str.len()
print(data['total_count_words'])
total_words = data['total_count_words'].sum()
print(total_words)
count = 0
for y in data["tweet"]:
    count = count + 1
print(count)
average_words = total_words / count
print (average_words)
```

```
0            9
1           16
2           19
3           19
4           11
            ..
10009       20
10010       10
10011       19
10012       20
10013       18
Name: total_count_words, Length: 10014, dtype: int64
163674
10014
16.344517675254643
```

In [68]:
```python
!pip3 install textblob
```

```
Defaulting to user installation because normal site-packages is not writeable
Collecting textblob
  Downloading textblob-0.17.1-py2.py3-none-any.whl (636 kB)
     |████████████████████████████████| 636 kB 1.7 MB/s
Requirement already satisfied: nltk>=3.1 in ./Library/Python/3.9/lib/python/si
te-packages (from textblob) (3.7)
Requirement already satisfied: joblib in /Library/Frameworks/Python.framework/
Versions/3.9/lib/python3.9/site-packages (from nltk>=3.1->textblob) (1.0.1)
Requirement already satisfied: click in ./Library/Python/3.9/lib/python/site-p
ackages (from nltk>=3.1->textblob) (8.0.3)
Requirement already satisfied: tqdm in /Library/Frameworks/Python.framework/Ve
rsions/3.9/lib/python3.9/site-packages (from nltk>=3.1->textblob) (4.62.3)
Requirement already satisfied: regex>=2021.8.3 in ./Library/Python/3.9/lib/pyt
hon/site-packages (from nltk>=3.1->textblob) (2022.3.15)
Installing collected packages: textblob
Successfully installed textblob-0.17.1
WARNING: You are using pip version 21.3.1; however, version 22.0.4 is availabl
e.
You should consider upgrading via the '/Library/Frameworks/Python.framework/Ve
rsions/3.9/bin/python3.9 -m pip install --upgrade pip' command.
```

In [193...
```python
from textblob import TextBlob
```

In [194...
```python
def analyze_sentiment(tweet):
    analysis = TextBlob(clean(tweet))
    if analysis.sentiment.polarity > 0:
        return 1
    elif analysis.sentiment.polarity == 0:
        return 0
    else:
        return -1
```

In [195...
```python
data['Sentiment'] = data['tweet'].apply(lambda x:analyze_sentiment(x))
data['Source'] = 'random_user'
data['Length'] = data['tweet'].apply(len)
data['Word_counts'] = data['tweet'].apply(lambda x:len(str(x).split()))
```

In [196...
```python
data1=data[['tweet','retweets_count', 'Sentiment', 'Source',
'Length','Word_counts']]
data1.head()
```

Out[196...

|   | tweet | retweets_count | Sentiment | Source | Length | Word_counts |
|---|-------|----------------|-----------|--------|--------|-------------|
| 0 | peterschiff padaprabu would shit pant chang n... | 0.0 | -1 | random_user | 64 | 9 |
| 1 | meatballsubzero pro russia pro ukrain cannot ... | 0.0 | 0 | random_user | 98 | 16 |
| 2 | subway pleas stop busi russia everi dollar sp... | 0.0 | 1 | random_user | 121 | 19 |
| 3 | russia prepar econom crisi develop expert nata... | 0.0 | 0 | random_user | 134 | 19 |
| 4 | bw putin fake news ðÿ"° rubl trash ðÿ— russia... | 0.0 | -1 | random_user | 59 | 11 |

In [197...
```python
data1['Clean tweet'] = data1['tweet'].apply(lambda x:clean(x))
```

In [198...
```python
data1[["Clean tweet","Sentiment"]].iloc[369]
```

Out[198...
```
Clean tweet     marshahairbrush rusembjakarta mfarussia natomi...
Sentiment                                                        0
Name: 369, dtype: object
```

In [200... 
```python
#displaying total number of neutral, positive and negative sentiments
sentiment = data1['Sentiment'].value_counts()
sentiment
```

Out[200...
```
 0    5094
 1    2788
-1    2132
Name: Sentiment, dtype: int64
```

In [201...
```python
#plotting the sentiments using seaborn library
plt.figure(figsize = (10,8))
sns.countplot(data = data1, x = 'Sentiment')
plt.show()
```



In [202...
```python
#defining values for neutral, positive and negative sentiments as 0, 1 and -1
neutral = data1[data1['Sentiment'] == 0]
positive = data1[data1['Sentiment'] == 1]
negative = data1[data1['Sentiment'] == -1]
```

In [203...
```python
#retrieving details about 2001th negative tweet
negative.iloc[2001]
```

Out[203...
```
tweet              russia would like get game
retweets_count                          0.0
Sentiment                                -1
Source                          random_user
Length                                   26
Word_counts                               5
Clean tweet        russia would like get game
Name: 9385, dtype: object
```

In [204...
```python
#retrieving details about 400th postive tweet
positive.iloc[400]
```

Out[204...
```
tweet              one thing love russia there much steak japan
retweets_count                                          0.0
Sentiment                                                 1
Source                                          random_user
Length                                                   44
Word_counts                                               8
Clean tweet            one thing love russia much steak japan
Name: 1428, dtype: object
```

In [205...
```python
#retrieving details about 4300th neutral tweet
neutral.iloc[4300]
```

Out[205...
```
tweet              kyivindepend russia lost  lost
retweets_count                              0.0
Sentiment                                     0
Source                              random_user
Length                                       31
Word_counts                                   4
Clean tweet        kyivindepend russia lost  lost
Name: 8454, dtype: object
```

In [206...
```python
print ("**********************************************************
#neutral_tweet
print("Example of a neutral tweet :",neutral['tweet'].values[3])
print ("**********************************************************
#positive tweet
print("Example of a positive tweet :",positive['tweet'].values[6])
print ("**********************************************************
#negative_text
print("Example of a negative tweet :",negative['tweet'].values[9])
print ("**********************************************************
```

```
********************************************************************************
********************
Example of a neutral tweet : prayerfeath russia putin war crimin putin held ac
count russia choke sanctionsaggress invas  total unaccept russia etern dame
********************************************************************************
********************
Example of a positive tweet : itâ€™ show  america amp nato arenâ€™t tri win gr
ound war russia give scrap ukrain compar could give  zelenski isnâ€™t tri stri
ke deal  long go drag
********************************************************************************
********************
Example of a negative tweet : mani compani step help put pressur russia list h
avent long includ reebok eddiebau ninewest subway halliburton dunkindonut gene
ralmil hiltonhotel marriott hyatt mani
********************************************************************************
********************
```

In [207…
```python
from wordcloud import WordCloud
```

In [208…
```python
txt = ' '.join(text for text in data1['Clean tweet'])
wordcloud = WordCloud(
            background_color = 'white',
            max_font_size = 100,
            max_words = 100,
            width = 800,
            height = 500
            ).generate(txt)
plt.imshow(wordcloud,interpolation = 'bilinear')
plt.axis('off')
plt.show()
```

In [209...

```python
#displaying the positive words using wordcloud
positive_words =' '.join([text for text in data1['Clean tweet'][data1['Sentim
#wordcloud = WordCloud(width=800, height=500, random_state=21, max_font_size=
wordcloud1 = WordCloud(
            random_state=21,
            max_font_size = 110,
            max_words = 100,
            width = 800,
            height = 500
            ).generate(positive_words)
plt.figure(figsize=(10, 7))
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.show()
```

In [210…
```python
#displaying the negative words using wordcloud
negative_words =' '.join([text for text in data1['Clean tweet'][data1['Sentime
#wordcloud = WordCloud(width=800, height=500, random_state=21, max_font_size=
wordcloud1 = WordCloud(
            random_state=21,
            max_font_size = 110,
            max_words = 100,
            width = 800,
            height = 500
            ).generate(negative_words)
plt.figure(figsize=(10, 7))
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.show()
```



In [102…
```python
#displaying the neutral words using wordcloud
neutral_words =' '.join([text for text in data1['Clean tweet'][data1['Sentime
wordcloud = WordCloud(width=800, height=500, random_state=21, max_font_size=1
plt.figure(figsize=(10, 7))
plt.imshow(wordcloud, interpolation="bilinear")
plt.axis('off')
plt.show()
```

# Thank you for time.

Dataset may be shared upon request.

In [ ]: