

Twitter Sentiment Analysis

In [1]:

```
1 # import 'Numpy'
2 import numpy as np
3
4 # import 'Pandas'
5 import pandas as pd
6
7 # import subpackage of Matplotlib
8 import matplotlib.pyplot as plt
9
10 # import color package from matplotlib
11 from matplotlib.colors import ListedColormap
12
13 # import 'Seaborn'
14 import seaborn as sns
15
16 # import datetime
17 import datetime
18
19 # to suppress warnings
20 from warnings import filterwarnings
21 filterwarnings('ignore')
22
23 # import regular expression
24 import re
25
26 # import ast
27 import ast
28
29 # display all columns of the dataframe
30 pd.options.display.max_columns = None
31
32 # import train-test split
33 from sklearn.model_selection import train_test_split
34
35 # import Tfidf
36 from sklearn.feature_extraction.text import TfidfVectorizer
37
38 # import various functions from sklearn
39 from sklearn.metrics import accuracy_score, roc_curve, roc_auc_score, classification_re
40 from sklearn.metrics import precision_score, recall_score
41
42 # import Logistic Regression
43 from sklearn.linear_model import LogisticRegression
44
45 # import DecisionTree Classifier
46 from sklearn.tree import DecisionTreeClassifier
47
48 # import RandomForest Classifier
49 from sklearn.ensemble import RandomForestClassifier
50
51 #import Naivebayes
52 from sklearn.naive_bayes import MultinomialNB
53
54 # import XtremeGradientBoost Classifier
55 from xgboost import XGBClassifier
56
57 # import NeuralNetwork
58 from sklearn.neural_network import MLPClassifier
59
```

```

60 # import svm
61 from sklearn.svm import SVC
62
63 #import xgboost
64 from xgboost import XGBClassifier
65
66 # import wordlcoud, stopwords
67 from wordcloud import WordCloud,STOPWORDS
68
69 #import nltk
70 import nltk
71 from nltk.corpus import stopwords

```

In [2]:

```

1 # Reading the dataset and viewing the first 5 rows of it
2
3 train = pd.read_csv('twitter_training.csv' , header = None)
4 test = pd.read_csv('twitter_validation.csv', header = None)
5
6 df = pd.concat([train , test] , ignore_index=True )
7
8 df.head()

```

Out[2]:

	0	1	2	3
0	2401	Borderlands	Positive	im getting on borderlands and i will murder yo...
1	2401	Borderlands	Positive	I am coming to the borders and I will kill you...
2	2401	Borderlands	Positive	im getting on borderlands and i will kill you ...
3	2401	Borderlands	Positive	im coming on borderlands and i will murder you...
4	2401	Borderlands	Positive	im getting on borderlands 2 and i will murder ...

In [3]:

```

1 # Renaming the columns
2
3 df.rename(columns = {0:'Id',1:'Platform',2:'Sentiment',3:'Review'} , inplace = True)

```

In [4]:

```

1 df.head()

```

Out[4]:

	Id	Platform	Sentiment	Review
0	2401	Borderlands	Positive	im getting on borderlands and i will murder yo...
1	2401	Borderlands	Positive	I am coming to the borders and I will kill you...
2	2401	Borderlands	Positive	im getting on borderlands and i will kill you ...
3	2401	Borderlands	Positive	im coming on borderlands and i will murder you...
4	2401	Borderlands	Positive	im getting on borderlands 2 and i will murder ...

In [5]:

```
1 # Dropping of id and platform
2
3 df.drop(columns = ['Id','Platform'] , inplace = True)
```

In [6]:

```
1 # Checking for missing values
2
3 df.isnull().sum()
```

Out[6]:

```
Sentiment      0
Review         686
dtype: int64
```

In [7]:

```
1 # Dropping of null values
2
3 df.dropna(subset = 'Review' , inplace = True)
```

In [8]:

```
1 df.isnull().sum()
```

Out[8]:

```
Sentiment      0
Review         0
dtype: int64
```

In [9]:

```
1 # Converting string into lower case characters
2
3 df['Review'] = df['Review'].apply(lambda x : str(x).lower())
```

In [10]:

```
1 # Removing special characters (Numbers , Punctuations , Characters)
2
3 df['Review'] = df['Review'].apply(lambda x : re.sub('[^a-z]+',' ',x))
```

In [11]:

```
1 # Doing Lemmatization
2
3 from nltk.stem import WordNetLemmatizer
4
5 wnl = WordNetLemmatizer()
6
7 df['Review'] = df['Review'].apply(lambda x : ' '.join([ wnl.lemmatize(i) for i in x.
```

In [12]:

```
1 # Creating a list of stopword
2
3 stopword = list(set(list(STOPWORDS) + list(stopwords.words('english'))))
4
5 len(stopword)
```

Out[12]:

227

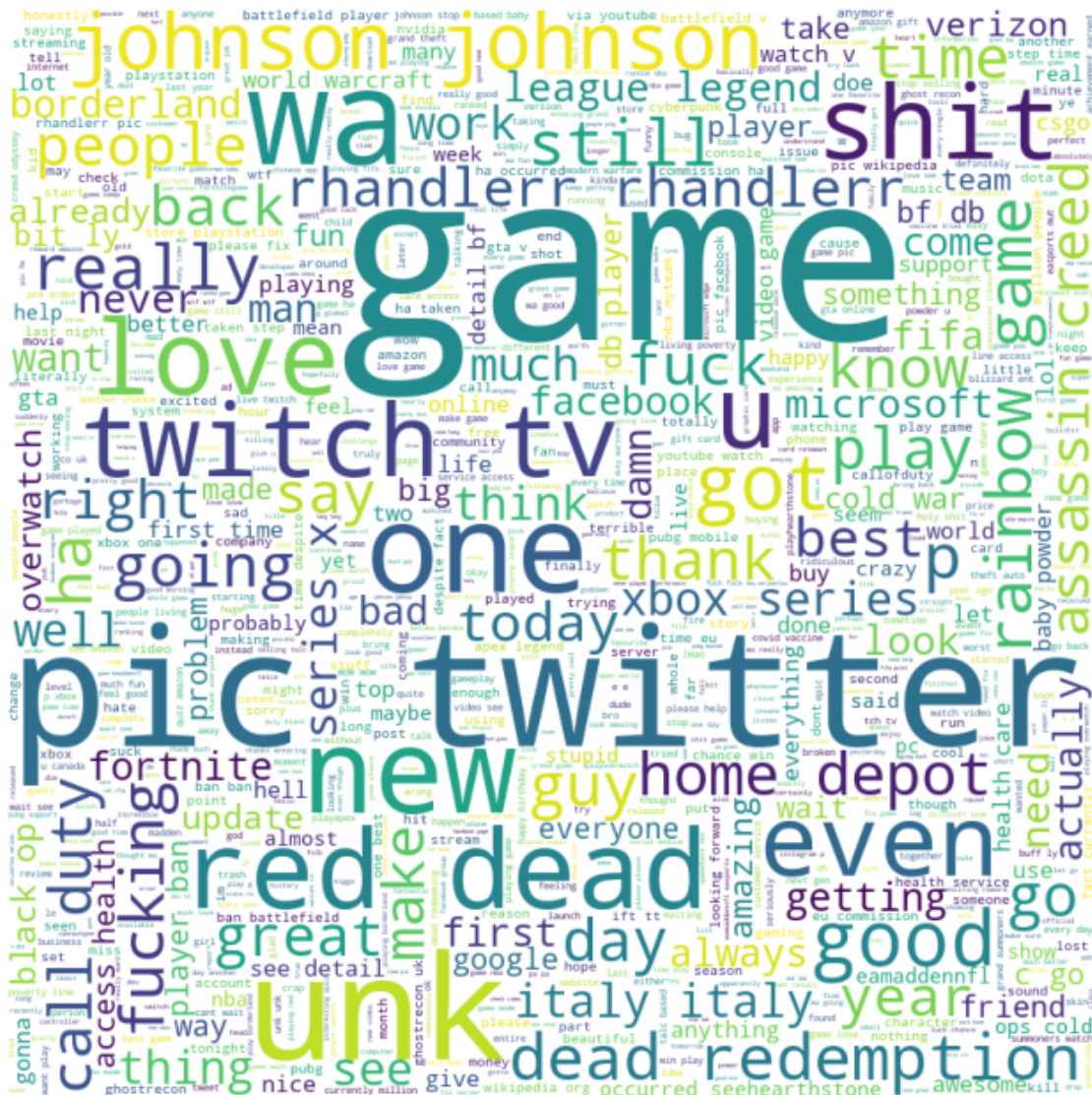
In [13]:

```
1 # Removing stopwords
2
3 df['Review'] = df['Review'].apply(lambda x : ' '.join([ i for i in x.split() if i not in stopword ]))
```

In [14]:

```
1 # Plotting of WordCloud
2
3 words = ''
4
5 for i in df['Review']:
6     words = words + i + ' '
```

```
1 wordcloud = WordCloud(width = 800 , height = 800 , background_color = 'white', max_w
2
3 plt.figure(figsize = (15,8))
4 plt.imshow(wordcloud)
5 plt.axis('off')
6 plt.show()
```



In [16]:

```

1 # Extracting Top 1000 words
2
3 word = pd.Series(words.split()).value_counts()
4 word

```

Out[16]:

```

game          11927
wa             5518
pic           4441
twitter       4390
one           3826
...
youthful      1
sloggettin    1
ascent        1
brooch        1
battleroyale  1
Length: 27145, dtype: int64

```

In [17]:

```
1 frequent_words = list(word[0:10000].index)
```

In [18]:

```
1 df['Review'] = df['Review'].apply(lambda x : ' '.join([i for i in x.split() if i in
```

In [19]:

```

1 # Encoding Target variable
2
3 df['Sentiment'].replace({'Positive':1 , 'Neutral':0 , 'Irrelevant':0 , 'Negative':-1

```

In [20]:

```

1 # Splitting the dataset randomly into train and test dataset using ratio of 70:30
2
3 x = df['Review']
4 y = df['Sentiment']
5
6 xtrain , xtest , ytrain , ytest = train_test_split(x,y,test_size = 0.30 , random_sta

```

In [21]:

```

1 # Intialising Tfidf and fitting into train dataset
2
3 vectorizer = TfidfVectorizer()
4
5 vectorizer.fit(xtrain)

```

Out[21]:

TfidfVectorizer()

In [22]:

```
1 len(vectorizer.get_feature_names_out())
```

Out[22]:

9984

In [23]:

```
1 # Transforming train and test dataset
2
3 tf_train = vectorizer.transform(xtrain)
4 tf_test = vectorizer.transform(xtest)
```

In [24]:

```
1 # Creating a dataframe for both train and test dataset
2
3 train_df = pd.DataFrame(tf_train.toarray() , columns = vectorizer.get_feature_names_out())
4 test_df = pd.DataFrame(tf_test.toarray() , columns = vectorizer.get_feature_names_out())
```


Building different models and evaluating using appropriate technique

In [25]:

```

1  # Creating a user defined function to store values of accuracy , precision , recall
2
3  performance_df = pd.DataFrame(columns = ['Model Name','Train Accuracy','Train Precision',
4                                           'Test Accuracy','Test Precision','Test Recall'])
5
6  def model_performance(model,name,xtrain = xtrain,xtest = xtest):
7      global performance_df
8
9      # Predicting train and test data
10
11     pred_train = model.predict(xtrain)
12     pred_test = model.predict(xtest)
13
14     # Calculating metrics for train and test data
15
16     train_acc = round(accuracy_score(ytrain,pred_train),2)*100
17     test_acc = round(accuracy_score(ytest,pred_test),2)*100
18     train_precision = round(precision_score(ytrain,pred_train, average = 'macro'),2)
19     test_precision = round(precision_score(ytest,pred_test, average = 'macro'),2)
20     train_recall = round(recall_score(ytrain,pred_train, average = 'macro'),2)
21     test_recall = round(recall_score(ytest,pred_test, average = 'macro'),2)
22     train_f1 = round(f1_score(ytrain,pred_train, average = 'macro'),2)
23     test_f1 = round(f1_score(ytest,pred_test, average = 'macro'),2)
24
25
26     # Adding train and test scores in performance_df dataframe
27
28     performance_df = performance_df.append({'Model Name': name, 'Train Accuracy':train_acc,
29                                             'Train Precision': train_precision,
30                                             'Train Recall': train_recall, 'Train f1_score': train_f1,
31                                             'Test Accuracy':test_acc,'Test Precision': test_precision,
32                                             'Test Recall':test_recall, 'Test f1_score': test_f1})
33
34     print('Train Report \n',classification_report(ytrain,pred_train),'\n')
35     print('Test Report \n',classification_report(ytest,pred_test))

```

In [26]:

```

1 # Building Model using Logistic Regression
2
3 model_lr = LogisticRegression(multi_class='multinomial').fit(train_df,ytrain)
4
5 model_performance(model_lr, name = 'Logistic Regression' , xtrain = train_df , xtest

```

Train Report

	precision	recall	f1-score	support
-1	0.85	0.82	0.83	15848
0	0.80	0.87	0.84	22067
1	0.83	0.76	0.79	14582
accuracy			0.82	52497
macro avg	0.83	0.82	0.82	52497
weighted avg	0.82	0.82	0.82	52497

Test Report

	precision	recall	f1-score	support
-1	0.80	0.76	0.78	6776
0	0.75	0.82	0.78	9373
1	0.77	0.70	0.74	6350
accuracy			0.77	22499
macro avg	0.77	0.76	0.77	22499
weighted avg	0.77	0.77	0.77	22499

In [27]:

```

1 # Building Model using Naive bayes
2
3 model_nb = MultinomialNB().fit(train_df,ytrain)
4
5 model_performance(model_nb , name = 'Naive Bayes' , xtrain = train_df , xtest = test

```

Train Report

	precision	recall	f1-score	support
-1	0.79	0.78	0.79	15848
0	0.76	0.84	0.80	22067
1	0.80	0.69	0.74	14582
accuracy			0.78	52497
macro avg	0.79	0.77	0.78	52497
weighted avg	0.78	0.78	0.78	52497

Test Report

	precision	recall	f1-score	support
-1	0.76	0.73	0.74	6776
0	0.71	0.80	0.75	9373
1	0.76	0.64	0.70	6350
accuracy			0.74	22499
macro avg	0.74	0.73	0.73	22499
weighted avg	0.74	0.74	0.73	22499

In [28]:

```

1 # Building Model using DecisionTree
2
3 model_dt = DecisionTreeClassifier(max_depth=750,criterion='gini').fit(train_df,ytrain)
4
5 model_performance(model_dt, name = 'DecisionTree' , xtrain = train_df , xtest = test_df)

```

Train Report

	precision	recall	f1-score	support
-1	0.99	0.95	0.97	15848
0	0.93	0.99	0.96	22067
1	0.99	0.94	0.96	14582
accuracy			0.97	52497
macro avg	0.97	0.96	0.97	52497
weighted avg	0.97	0.97	0.97	52497

Test Report

	precision	recall	f1-score	support
-1	0.84	0.80	0.82	6776
0	0.79	0.86	0.83	9373
1	0.83	0.77	0.80	6350
accuracy			0.82	22499
macro avg	0.82	0.81	0.81	22499
weighted avg	0.82	0.82	0.82	22499

In [29]:

```

1 # Building Model using RandomForest
2
3 model_rf = RandomForestClassifier(max_depth=750,criterion='gini',n_estimators = 150,
4
5 model_performance(model_rf , name = 'Random Forest' , xtrain = train_df , xtest = te

```

Train Report

	precision	recall	f1-score	support
-1	0.99	0.96	0.98	15848
0	0.94	0.99	0.97	22067
1	0.99	0.95	0.97	14582
accuracy			0.97	52497
macro avg	0.98	0.97	0.97	52497
weighted avg	0.97	0.97	0.97	52497

Test Report

	precision	recall	f1-score	support
-1	0.95	0.90	0.92	6776
0	0.87	0.96	0.91	9373
1	0.94	0.86	0.90	6350
accuracy			0.91	22499
macro avg	0.92	0.91	0.91	22499
weighted avg	0.92	0.91	0.91	22499

In [31]:

```

1 performance_df

```

Out[31]:

	Model Name	Train Accuracy	Train Precision	Train Recall	Train f1_score	Test Accuracy	Test Precision	Test Recall	Test f1_score
0	Logistic Regression	82.0	0.83	0.82	0.82	77.0	0.77	0.76	0.77
1	Naive Bayes	78.0	0.79	0.77	0.78	74.0	0.74	0.73	0.73
2	DecisionTree	97.0	0.97	0.96	0.97	82.0	0.82	0.81	0.81
3	Random Forest	97.0	0.98	0.97	0.97	91.0	0.92	0.91	0.91