

Name: Souvik Banerjee

Data Science March Major Project

Problem statement:

Create a classification model to predict the sentiment either (Positive or Negative) based on Covid Tweets

Context:

The tweets have been pulled from Twitter and manual tagging has been done then. The names and usernames have been given codes to avoid any privacy concerns.

Imports

```
In [2]: import numpy as np
import pandas as pd
import seaborn as sns
import nltk
from nltk.stem.porter import PorterStemmer
import sklearn.model_selection
import sklearn.svm, sklearn.neighbors, sklearn.metrics
import matplotlib.pyplot as plt
import nltk.corpus
nltk.download('wordnet')
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('punkt')
nltk.download('punkt')
nltk.download('punkt')
```

Dataset

```
In [3]: df=pd.read_csv('dataset.csv',encoding='latin-1')
```

```
In [4]: df.head()
```

```
Out[4]:
```

	UserName	ScreenName	Location	TweetAt	OriginalTweet	Sentiment
0	3799	48751	London	16-03-2020	@Mehyrbt @Phil_Gahan @Chrisv https://t.co...	Neutral
1	3800	48752	UK	16-03-2020	advice Talk to your neighbours family to excha...	Positive
2	3801	48753	Vagabonds	16-03-2020	Coronavirus Australia: Woolworths to give elde...	Positive
3	3802	48754	U	16-03-2020	My food stock is not the only one which is emp...	Positive
4	3803	48755	Nah	16-03-2020	Me, ready to go at supermarket during the #COVID...	Extremely Negative

```
In [5]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 41157 entries, 0 to 41156
Data columns (total 6 columns):
 #   Column          Non-Null Count  Dtype
---  --
 0   UserName        41157 non-null   int64
 1   ScreenName      41157 non-null   int64
 2   Location        41157 non-null   object
 3   TweetAt        41157 non-null   object
 4   OriginalTweet   41157 non-null   object
 5   Sentiment       41157 non-null   object
dtypes: int64(2), object(4)
memory usage: 1.9+ MB
```

```
In [6]: df.describe()
```

```
Out[6]:
```

	UserName	ScreenName
count	41157.000000	41157.000000
mean	24377.000000	69329.000000
std	11481.166667	11081.166667
min	3799.000000	48751.000000
25%	14008.000000	69040.000000
50%	24377.000000	69329.000000
75%	34666.000000	79616.000000
max	44955.000000	89987.000000

```
In [7]: df.shape
(41157, 6)
```

```
In [8]: df.columns
```

```
Index(['UserName', 'ScreenName', 'Location', 'TweetAt', 'OriginalTweet',
       'Sentiment'],
      dtype='object')
```

Remove null value (if any)

```
In [9]: df.isnull().sum()
```

```
Out[9]:
```

	UserName	ScreenName	Location	TweetAt	OriginalTweet	Sentiment
0	3799	48751	London	16-03-2020	@Mehyrbt @Phil_Gahan @Chrisv https://t.co...	Neutral
1	3800	48752	UK	16-03-2020	advice Talk to your neighbours family to excha...	Positive
2	3801	48753	Vagabonds	16-03-2020	Coronavirus Australia: Woolworths to give elde...	Positive
3	3802	48754	U	16-03-2020	My food stock is not the only one which is emp...	Positive
4	3803	48755	Nah	16-03-2020	Me, ready to go at supermarket during the #COVID...	Extremely Negative

```
In [10]: df['Location'].fillna('U')
```

```
Out[10]:
```

	UserName	ScreenName	Location	TweetAt	OriginalTweet	Sentiment
0	3799	48751	London	16-03-2020	@Mehyrbt @Phil_Gahan @Chrisv https://t.co...	Neutral
1	3800	48752	UK	16-03-2020	advice Talk to your neighbours family to excha...	Positive
2	3801	48753	Vagabonds	16-03-2020	Coronavirus Australia: Woolworths to give elde...	Positive
3	3802	48754	U	16-03-2020	My food stock is not the only one which is emp...	Positive
4	3803	48755	Nah	16-03-2020	Me, ready to go at supermarket during the #COVID...	Extremely Negative

```
In [11]: df['Location'].value_counts()[0:10]
```

```
Out[11]:
```

	Location	count
0	London	8582
1	United States	548
2	London, England	528
3	New York, NY	395
4	Staffordshire Moorlands	1
5	Kilchobert, OH	1
6	Tulsa, OK	1
7	Wentford, South Yorkshire, Rotherham	1
8	I love you so much he/she	1
9	Name: Location, dtype: int64	

```
In [12]: df['Location'].value_counts().sort_values(ascending=False)
```

```
Out[12]:
```

	Location	count
0	London	8582
1	United States	548
2	London, England	528
3	New York, NY	395
4	Staffordshire Moorlands	1
5	Kilchobert, OH	1
6	Tulsa, OK	1
7	Wentford, South Yorkshire, Rotherham	1
8	I love you so much he/she	1
9	Name: Location, dtype: int64	

```
In [13]: df['Sentiment'].value_counts()
```

```
Out[13]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [14]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[14]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [15]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[15]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [16]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[16]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [17]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[17]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [18]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[18]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [19]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[19]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [20]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[20]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [21]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[21]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [22]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[22]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [23]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[23]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [24]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[24]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [25]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[25]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [26]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[26]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [27]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[27]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [28]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[28]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [29]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[29]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [30]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[30]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [31]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[31]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [32]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[32]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [33]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[33]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [34]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[34]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [35]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[35]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [36]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[36]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [37]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[37]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [38]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[38]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [39]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[39]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [40]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[40]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [41]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[41]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [42]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[42]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [43]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[43]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	7713
3	Extremely Positive	6024
4	Extremely Negative	5481
5	Name: Sentiment, dtype: int64	

```
In [44]: df['Sentiment'].value_counts().sort_values(ascending=False)
```

```
Out[44]:
```

	Sentiment	count
0	Positive	14422
1	Negative	9917
2	Neutral	771