

TASK 1 - Exploratory Data Analysis (EDA) with Twitter datasets given by Spotle.ai

Team Name - KeenNinjas

Team Members :-

Abhinav Sharma

Krishna Chaitanya Vadduri

Rishabh Purohit

Madhukar Paila

Importing the Libraries

We start off by importing the python libraries which will be used for the task.

List of Python Libraries Used (with their source code/documentation in braces):-

- Numpy (<https://numpy.org/> (<https://numpy.org/>))
- Pandas (<https://pandas.pydata.org/> (<https://pandas.pydata.org/>))
- Matplotlib (<https://matplotlib.org/> (<https://matplotlib.org/>))
- WordCloud (https://amueller.github.io/word_cloud/index.html (https://amueller.github.io/word_cloud/index.html))
- Plotly (<https://github.com/plotly/plotly.py> (<https://github.com/plotly/plotly.py>))
- Cufflinks (<https://github.com/santosjorge/cufflinks> (<https://github.com/santosjorge/cufflinks>))

```
In [1]: 1 import numpy as np
        2 import pandas as pd
        3 import warnings
        4 import re
        5 #Visualisation
        6 import matplotlib.pyplot as plt
        7 from wordcloud import WordCloud, STOPWORDS
        8 warnings.filterwarnings("ignore")
        9 %matplotlib inline
       10 from plotly.offline import iplot
       11 import plotly as py
       12 import plotly.tools as tls
       13 import cufflinks as cf
       14 py.offline.init_notebook_mode(connected = True)
```

```
In [2]: 1 def configure_plotly_browser_state():
2     import IPython
3     display(IPython.core.display.HTML('''
4         <script src="/static/components/requirejs/require.js"></script>
5         <script>
6             requirejs.config({
7                 paths: {
8                     base: '/static/base',
9                     plotly: 'https://cdn.plot.ly/plotly-latest.min.js?noext',
10                },
11            });
12        </script>
13        '''))
```

Here we imported the 'tweets_corona.txt' dataset which is the dataset provided by Spotle.ai to perform EDA . Also we edited the dataframe accordingly.

```
In [6]: 1 text = r"tweets_corona.txt"
2     splitLine = []
3     oFile = open(text, 'r', encoding="utf8")
4     line = oFile.readline()
5     while line:
6         splitLine.append(line.split('\n'))
7         line = oFile.readline()
8     oFile.close()
```

```
In [7]: 1 tweets = []
2     for sublist in splitLine:
3         for item in sublist:
4             tweets.append(item)
```

Sub Task 1 - A tag cloud depicting what topics / Word were being talked about on Twitter

[illegible]

Sub Task 2 - Which hashtag trended (Hashtags are words or phrases beginning with # eg #COVID)

localhost:8888/notebooks/Desktop/Data Analysis (final).ipynb#

```
In [9]: 1 raw = ' '.join(tweets)
2 tags = [re.sub(r"(\W+)$", "", j) for j in [i for i in raw.split() if i.startswith("#)]]
3 df_hash = pd.DataFrame({"hashtag": tags})
4 print(df_hash['hashtag'].value_counts().head(10))
```

```
#COVID19      49016
#Corona        30586
#lockdown      23420
#coronavirus    21504
#corona         19609
#covid19        18718
#Covid19         8530
#Covid_19        6145
#COVID-19        6056
#StayHome       4796
Name: hashtag, dtype: int64
```

```
In [10]: 1 df_hashtag = df_hash['hashtag'].value_counts().head(10)
```

We extracted the required hashtags from the database and made a dataframe of those twitter hastags and their value counts which were extracted from the dataset by the following code.

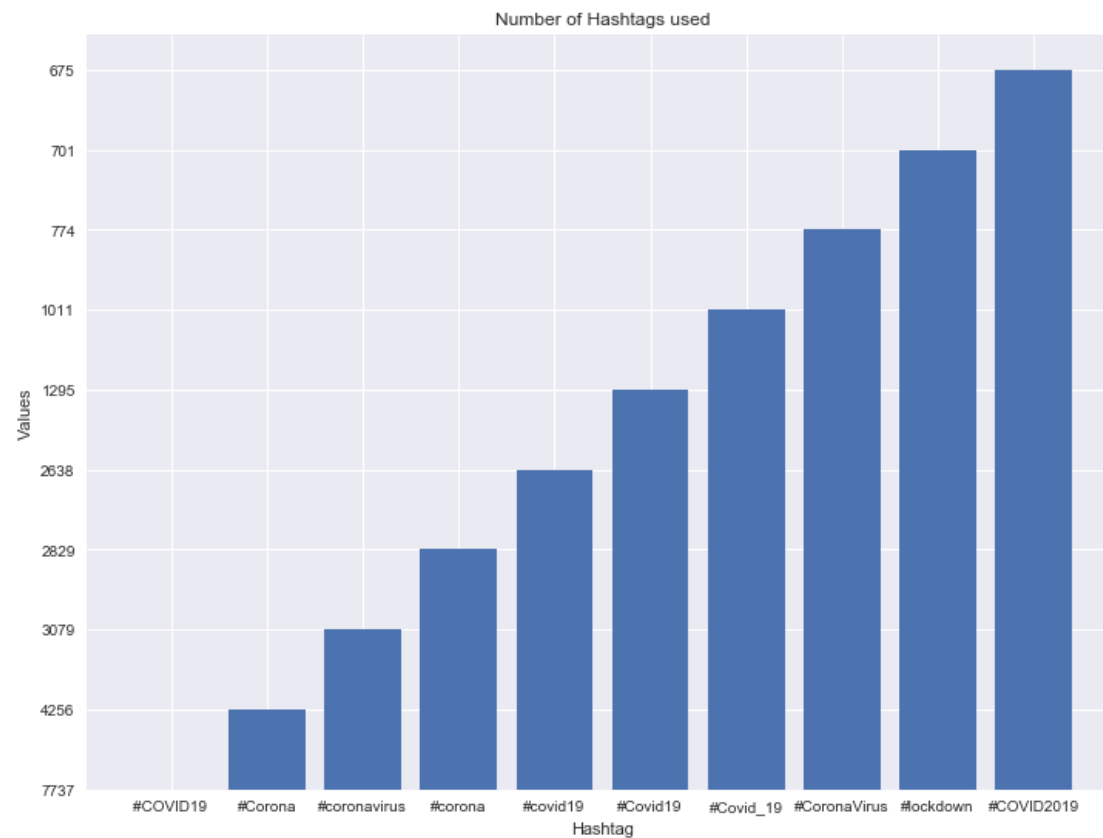
```
In [11]: 1 df_hashtag = pd.DataFrame({'hashtag':df_hashtag.index, 'val':df_hashtag.values})
2 df_hashtag
```

Out[11]:

	hashtag	val
0	#COVID19	49016
1	#Corona	30586
2	#lockdown	23420
3	#coronavirus	21504
4	#corona	19609
5	#covid19	18718
6	#Covid19	8530
7	#Covid_19	6145
8	#COVID-19	6056
9	#StayHome	4796

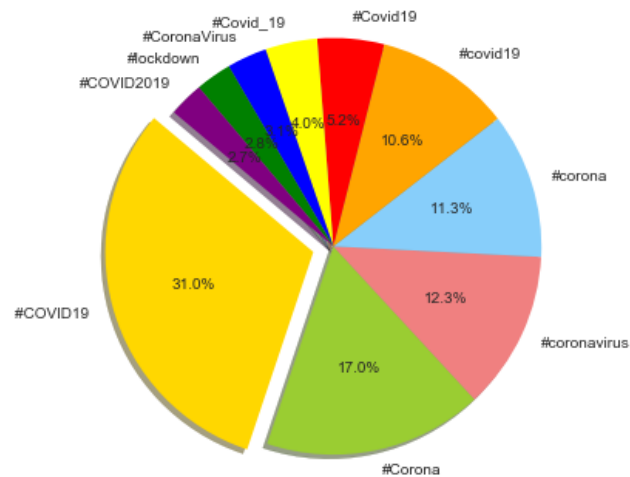
Following is the code to plot a bar chart for the most trending twitter hashtags:-

```
In [12]: 1 plt.figure(figsize=(12,9))
2 plt.style.use("seaborn")
3 objects = ('#COVID19 ', '#Corona', '#coronavirus', '#corona', '#covid19', '#Covid19', '#Covid_19', '#CoronaVirus', '#lockdown', '#COVID2019')
4 y_pos = np.arange(len(objects))
5 performance = ['7737', '4256', '3079', '2829', '2638', '1295', '1011', '774', '701', '675']
6
7 plt.bar(y_pos, performance, align='center', alpha=1.0, width=0.8)
8 plt.xticks(y_pos, objects)
9 plt.xlabel('Hashtag')
10 plt.ylabel('Values')
11 plt.title('Number of Hashtags used')
12
13 plt.show()
```



Following is the code to plot a pie chart for the most trending twitter hashtags:-

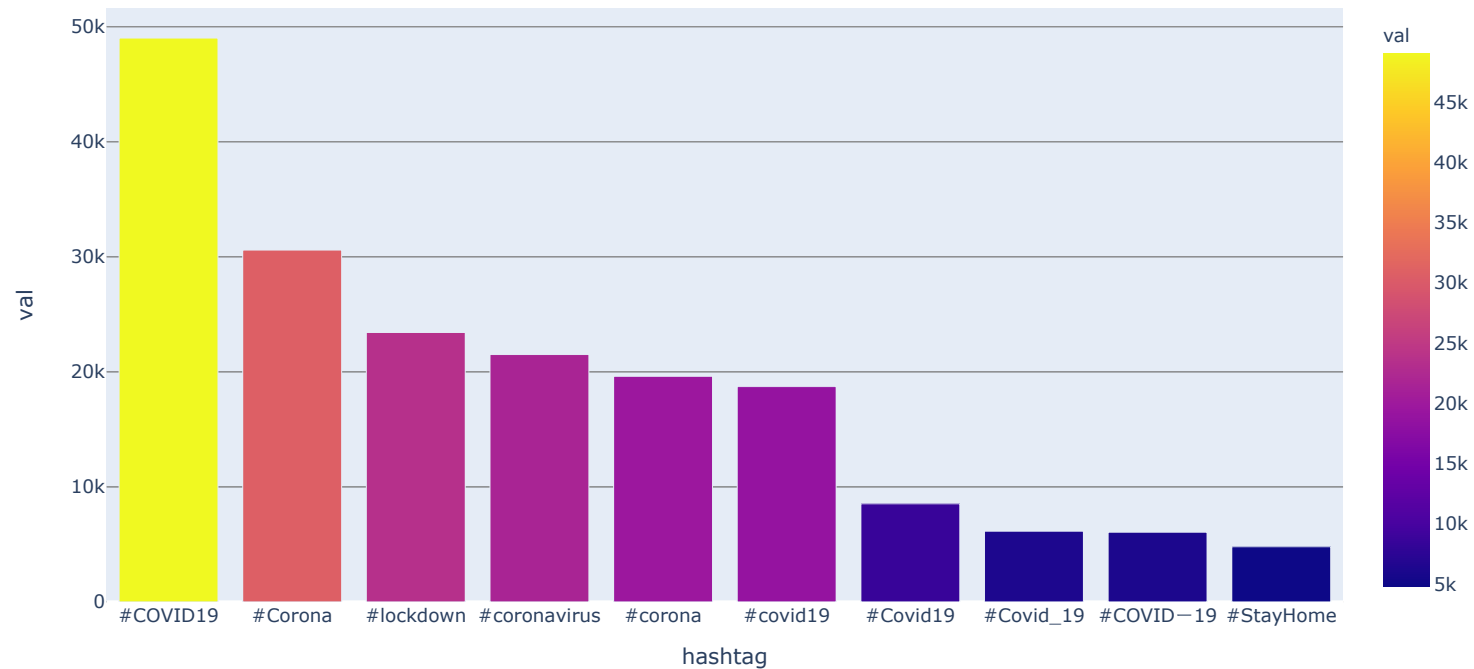
```
In [13]: 1 labels = ('#COVID19 ', '#Corona', '#coronavirus', '#corona', '#covid19', '#Covid19', '#Covid_19', '#CoronaVirus', '#lockdown', '#COVID2019')
2 sizes = ['7737', '4256', '3079', '2829', '2638', '1295', '1011', '774', '701', '675']
3 colors = ['gold', 'yellowgreen', 'lightcoral', 'lightskyblue', 'orange', 'red', 'yellow', 'blue', 'green', 'purple']
4 explode = (0.1, 0, 0, 0, 0, 0, 0, 0, 0, 0) # explode 1st slice
5
6 # Plot
7 plt.pie(sizes, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%', shadow=True, startangle=140)
8
9 plt.axis('equal')
10 plt.show()
```



We plotted a pie-chart for the trending Hashtags.

Following is the code to plot a bar-chart (using plotly library) for the the most trending twitter hashtags:-

```
In [14]: 1 configure_plotly_browser_state()
2 import plotly.express as px
3 fig = px.bar(df_hashtag, x='hashtag', y='val', hover_data=['val'], color='val')
4 fig.show()
```



We displayed a bar-chart in descending order for the trending Hashtags.

Sub Task 3 - Which Twitter Handler which dominated conversation on Twitter

Following is the code to get the most trending twitter handles by using the keyword '@' :-

```
In [15]: 1 raw = ' '.join(tweets)
2 tags = [re.sub(r"(\W+)$", "", j) for j in [i for i in raw.split() if i.startswith("@")]]
3 df_handler = pd.DataFrame({"handler": tags})
4 print(df_handler)
```

```

      handler
0    @HealthMedicalE1
1         @diprjk
2    @kansalrohit69
3    @DrSyedSehrish
4    @MoHFW_INDIA
...
59424    @DrRPNishank
59425    @CMOMaharashtra
59426    @HRDMinistry
59427    @narendramodi
59428    @cmnishank
```

[59429 rows x 1 columns]

```
In [16]: 1 df=df_handler['handler'].value_counts().head(10)
2 df
```

```
Out[16]: @narendramodi      1455
@PMOIndia      1295
@realDonaldTrump      837
@YouTube      725
@WHO      666
      649
@news_pandemic      559
@MoHFW_INDIA      446
@AmitShah      426
@Olacabs      348
Name: handler, dtype: int64
```

We made a dataframe of those twitter handles extracted from the dataset in the following code.

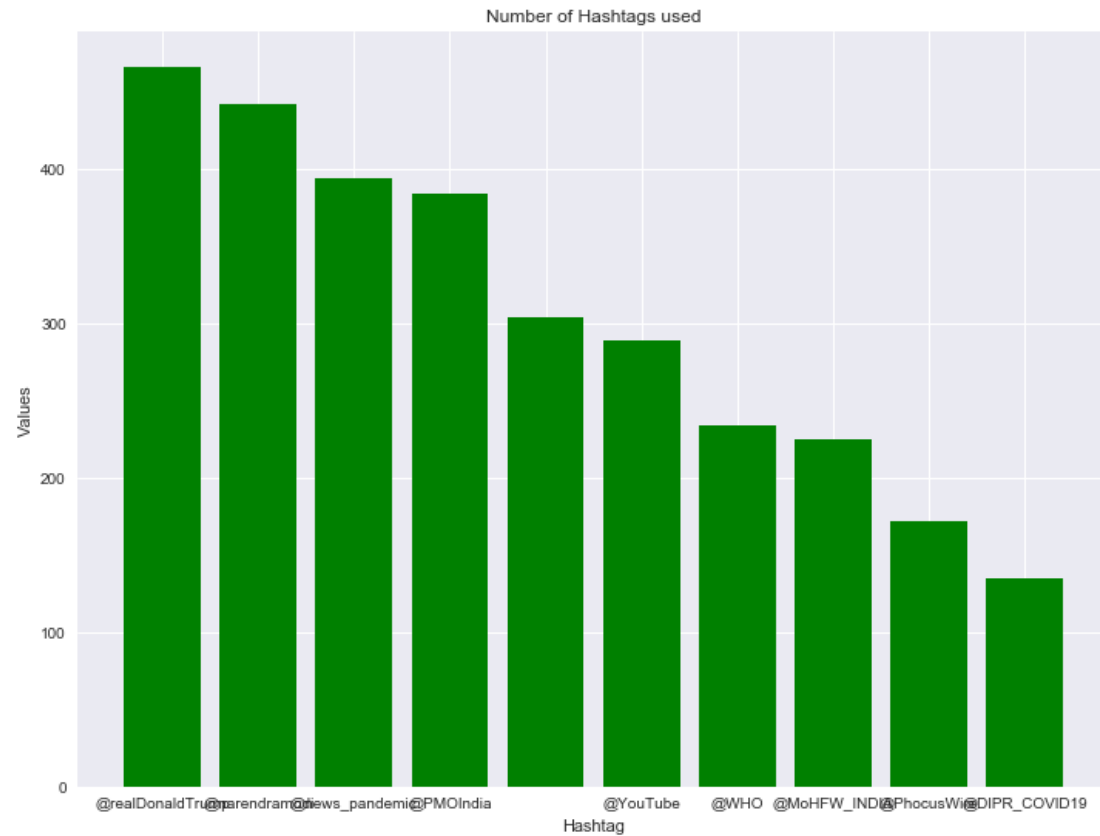

```
In [17]: 1 df_handlers = pd.DataFrame({'handler':df.index, 'val':df.values})  
        2 df_handlers
```

Out[17]:

	handler	val
0	@narendramodi	1455
1	@PMOIndia	1295
2	@realDonaldTrump	837
3	@YouTube	725
4	@WHO	666
5		649
6	@news_pandemic	559
7	@MoHFW_INDIA	446
8	@AmitShah	426
9	@Olacabs	348

Following is the code to plot a bar-chart for the the most dominating twitter handles :-

```
In [18]: 1 plt.figure(figsize=(12,9))
2 plt.style.use("seaborn")
3 objects = ('@realDonaldTrump', '@narendramodi', '@news_pandemic', '@PMOIndia', '@YouTube', '@WHO', '@MoHFW_INDIA', '@PhocusWire', '@DIPR_COVID19')
4 y_pos = np.arange(len(objects))
5 performance = [466,442,394,384,304,289,234,225,172,135]
6
7 plt.bar(y_pos, performance, align='center', alpha=1.0, width=0.8, color='green')
8 plt.xticks(y_pos, objects)
9 plt.xlabel('Hashtag')
10 plt.ylabel('Values')
11 plt.title('Number of Hashtags used')
12
13 plt.show()
```

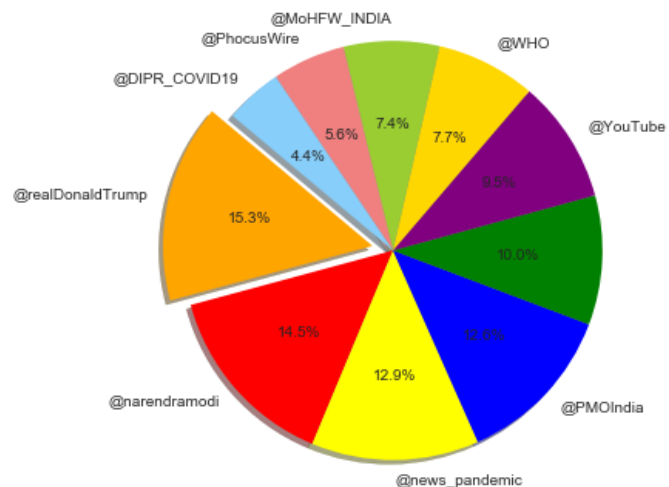


Following is the code to plot a pie chart for the the most dominating twitter handles :-

```

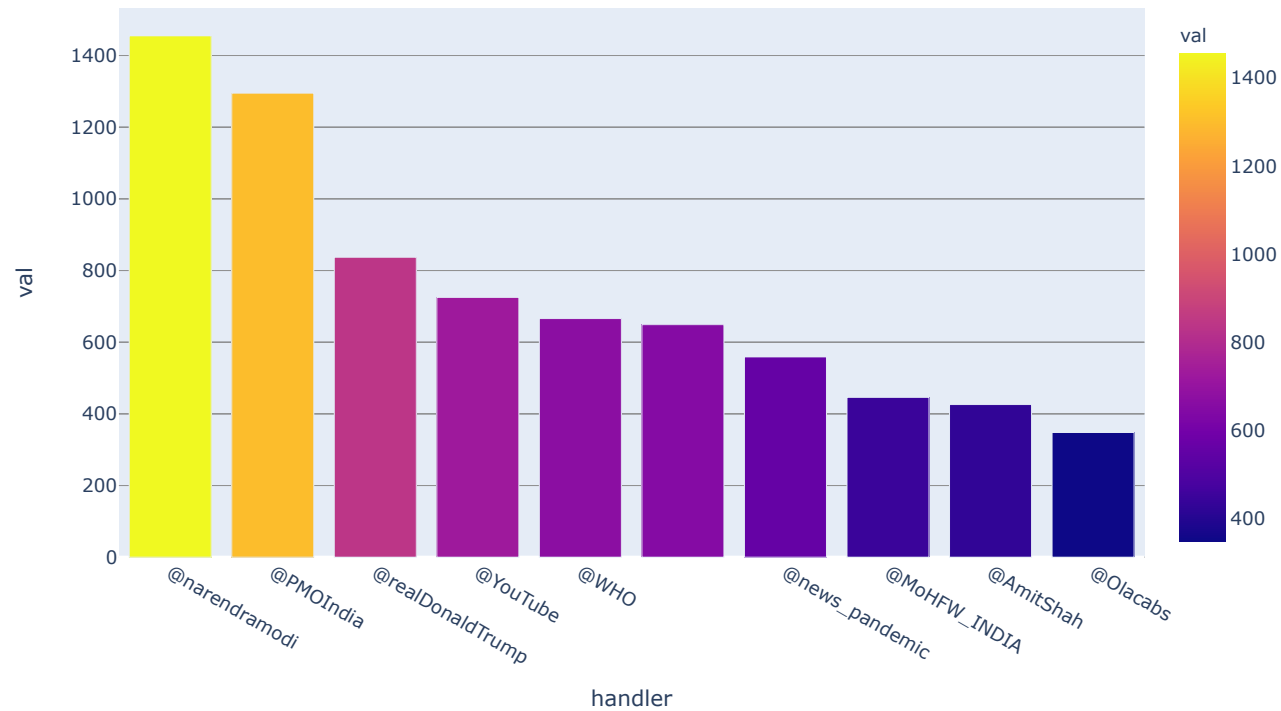
In [19]: 1 labels = ('@realDonaldTrump', '@narendramodi', '@news_pandemic', '@PMOIndia', ' ', '@YouTube', '@WHO', '@MoHFW_INDIA', '@PhocusWire', '@DIPR_COVID19')
2 sizes = [466, 442, 394, 384, 304, 289, 234, 225, 172, 135]
3 colors = ['orange', 'red', 'yellow', 'blue', 'green', 'purple', 'gold', 'yellowgreen', 'lightcoral', 'lightskyblue']
4 explode = (0.1, 0, 0, 0, 0, 0, 0, 0, 0, 0) # explode 1st slice
5
6 # Plot
7 plt.pie(sizes, explode=explode, labels=labels, colors=colors, autopct='%1.1f%%', shadow=True, startangle=140)
8
9 plt.axis('equal')
10 plt.show()

```



Following is the code to plot a bar-chart (using plotly library) for the the most dominating twitter handles :-

```
In [20]: 1 configure_plotly_browser_state()
2
3 fig = px.bar(df_handlers, x='handler', y='val' , hover_data=['val'], color='val')
4 fig.show()
```



You've reached the end of the notebook.

You've reached the end of the notebook

This IPython notebook contains all the code, required details, plots and all EDA performed on the dataset provided to us by Spotle.ai. Each step taken during the EDA has been explained here.

This concludes Assignment for Team - **KeenNinjas**

```
In [ ]: 1
```

```
In [ ]: 1
```

