## **Summary of EDA**

## **Tweet Dataset:**

Dataset	Multivariate,	Number of	50000	Area:	Text
Characteristics:	Text	Instances:			
Attribute	Real	Number of	11	Date Created:	1 <sup>st</sup> of March,
Characteristics:		Attributes:			2020
Associated	Classification,	Missing	None	Type of	Unsupervised
Task:	Clustering	Values?		Learning:	

**Dataset Information:** The dataset was collected using the Twitter API called Tweepy. Tweepy API provides access to the entire RESTful API methods. Methods accept various parameters and returns responses. To collect the records, the tweets were streamed and stored in a json format first and then converted to the csv format. It has 11 attributes and more than 50000 records collected over three days' period.

**Attribute Information:** The attributes collected here are the information about the tweets like, from where it was tweeted (location), who tweeted it (id), when this was tweeted (created\_at), how many followers does the person doing tweet have (followers\_count) etc. These attributes are considered to be parameters that will help in determining the sentiments of the tweets.

**Time Period Covered:** As the Tweepy API helps in extracting the tweets of the time when it is running, we have collected the recent tweets over the three days' period, that was since 28<sup>th</sup> Feb, 2020.

Brief summary of any data cleaning steps you have performed. For example, are there any particular observations / time periods / groups / etc. you have excluded?

## **Data Cleaning:**

There were two different stages where data cleaning was performed:

- 1) During Streaming: This was the time when with the use of Tweepy API we listened to the tweets and started collecting the tweet's information. Every tweet had over 25 attributes and not all of them were needed for the data analysis. We captured those features and attributes that were necessary to evaluate the sentiments of the tweets and discarded rest of them. This is one of the dimensionality reduction technique where unnecessary attributes were removed. To tackle redundancy, we removed the retweets from consideration and so set the retweeted status to False.
- 2) After Streaming: The first thing to be taken care was the missing values. Sometimes there remains some missing value that can result in bad analysis, so removal of those missing values was a good step. We could have replaced the value with the mean or median of the rest of the tuples but it doesn't seem right because it is hard to judge something like sentiments by just some random numbers. The next step of cleaning was to remove all the stop words from the tweets. Stopwords are useless in such analysis and will only add to the computation and complexity of the task. This will help us further in just focusing on the words that actually add meaning to our analysis. This way we could focus more on the adjectives that was used which is considered to be an important factor of analyzing the sentiments of any tweet.

Description of outcome with an appropriate visualization technique.

• The visualization technique we will use is the scatter plot for different features to understand the relation between each features with every other. The following will be the graph of all the plots that we have found:

Description of key predictors with appropriate visualization techniques that compare predictors to the response. You should investigate all predictors in your data as part of your project. For the purpose of this assignment, pick the one or two predictors that you think are going to be most important in explaining the outcome. Your selection of predictors can either be guided by your domain knowledge or be the result of your EDA on all predictors.