# Problem Statement

In this Project we have given the data set of tweets done by the people during the coronavirus pandemic from various countries. We have to see in this data what type of comment that has been done during the pandemic.
We get to know from our data that there are five types of sentiment that are extremely positive, extremely negative , positive, negative and neutral.
From various classification algorithm we have to find out that whether the tweet that has been made belongs to which type of sentiment.
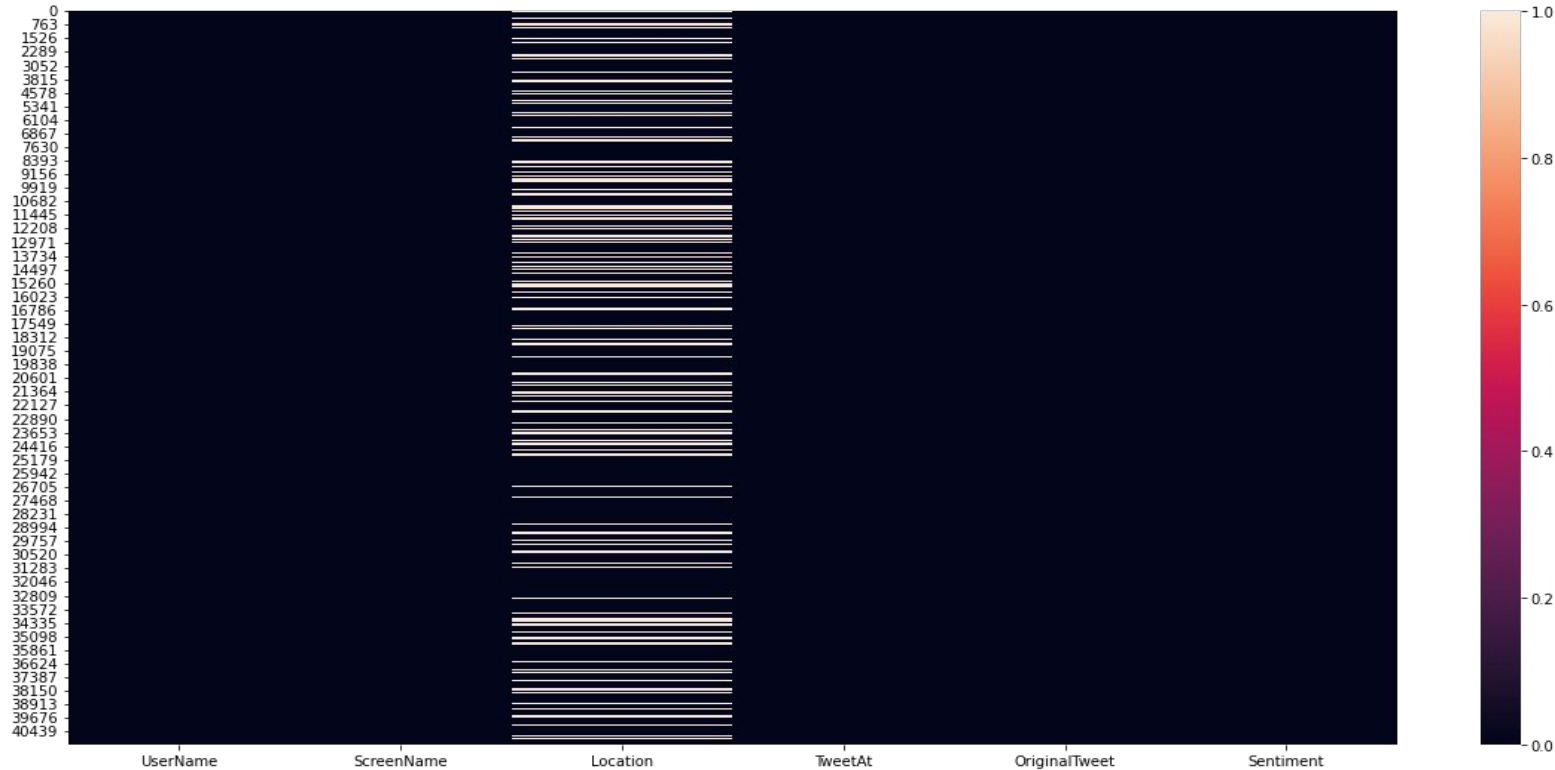
# Data Summary

● The original dataset has 6 columns and 41157 rows.

● In order to analyse various sentiments, We require just two columns named Original Tweet and Sentiment.

● There are four types of sentiments- Extremely Negative, Negative, Neutral, Positive and Extremely Positive.
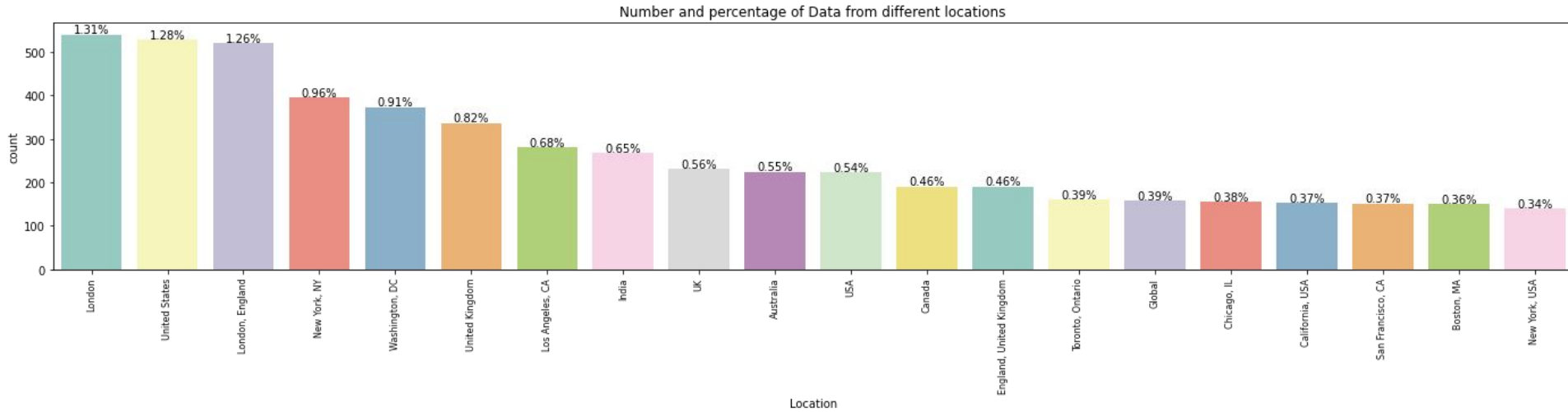
# Dataset Description

| Fields | Description |
|---|---|
| Username | Coded Username |
| ScreenName | Coded ScreenName |
| Location | Region of origin |
| TweetAt | Tweet Timing |
| OriginalTweet | First tweet in the thread |
| Sentiment-Target variable | Sentiment of the tweet |

# Heatmap to check the Nan values
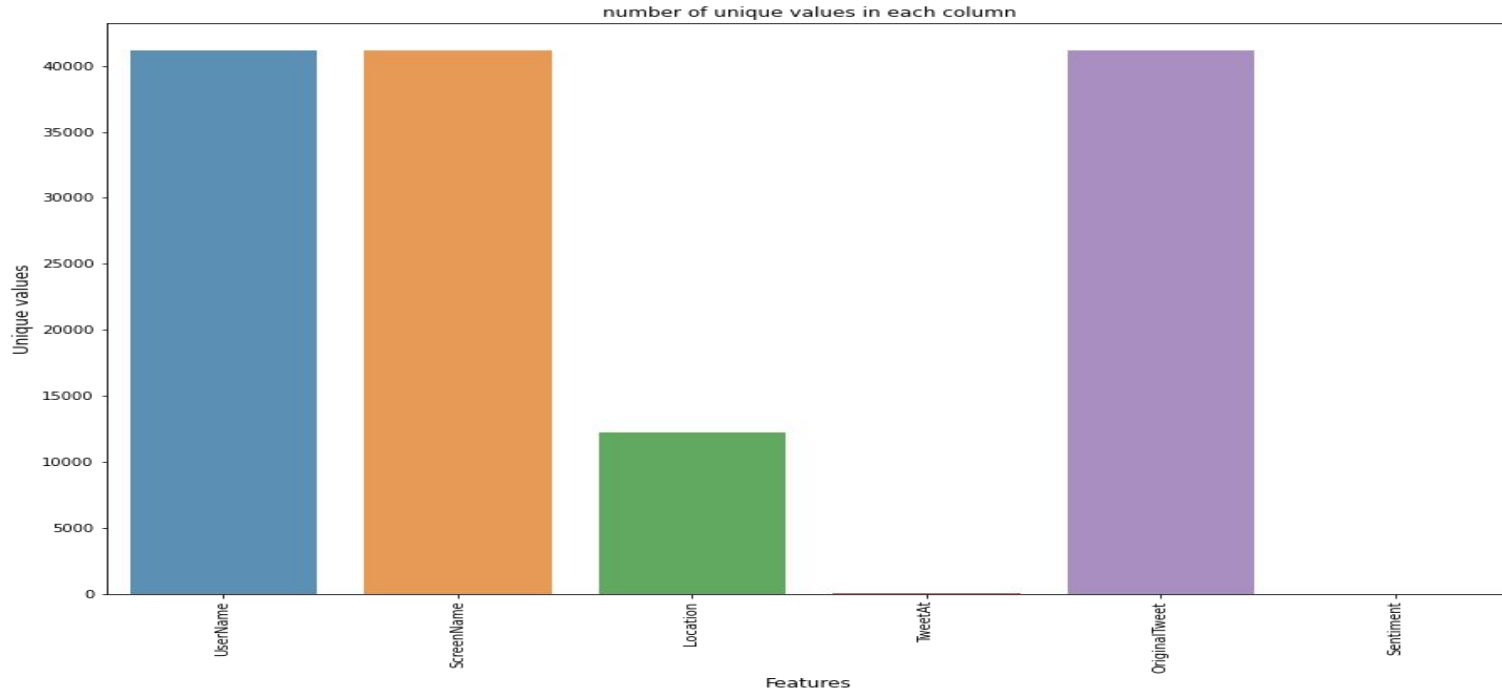
There are 20.87%(8567) null values in various places of location column.

# Number and percentage of data from different locations



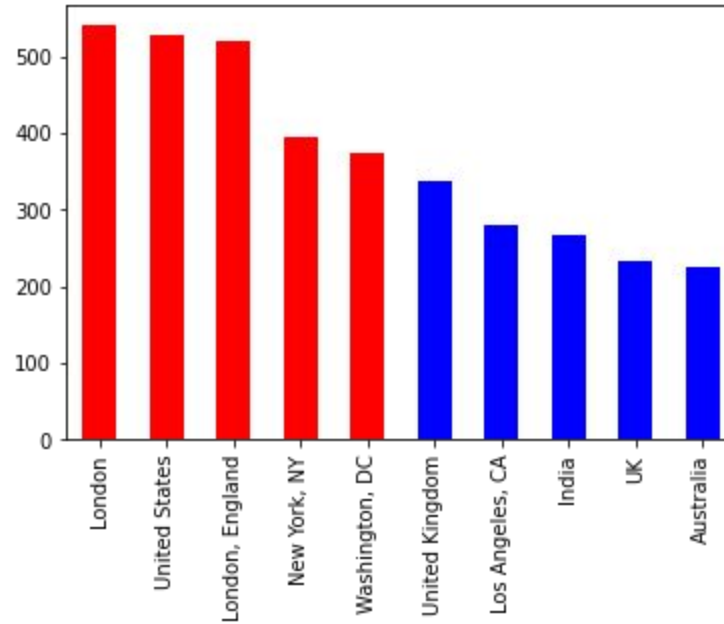Number and percentage of Data from different locations

The above observation shows that the number of data percentage wise from each country that of total data present. We can see that mainly United kingdom,United states of america , Indian and Canada are the countries from where tweet has been done.

# Number of unique values in each columns



The number of unique value in different feature are as followed Username , ScreenName , OriginalTweet has unique values of around 40000 then followed by Location that is between 10000-15000.
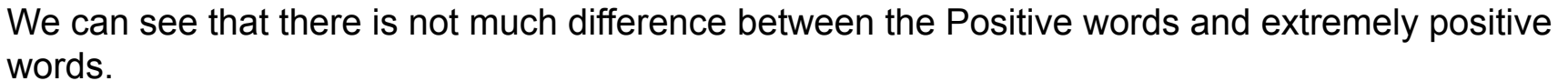
London has most number of unique tweet in developed cities and India has most number of tweet in list of developing nation.
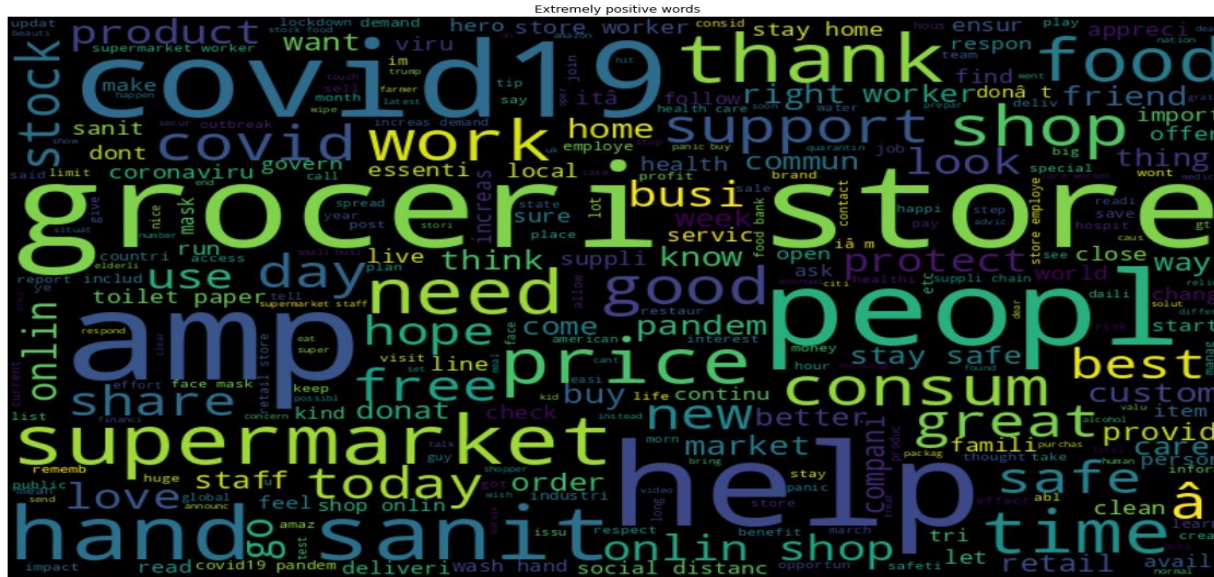
# Wordplay

## Wordcloud



In this analysis we get to know which type of word is most frequent and how many times that they have been used during the pandemic while tweeting.

**Positive words.**


positive words

We can see that there is not much difference between the Positive words and extremely positive words.

# Extremely Positive words
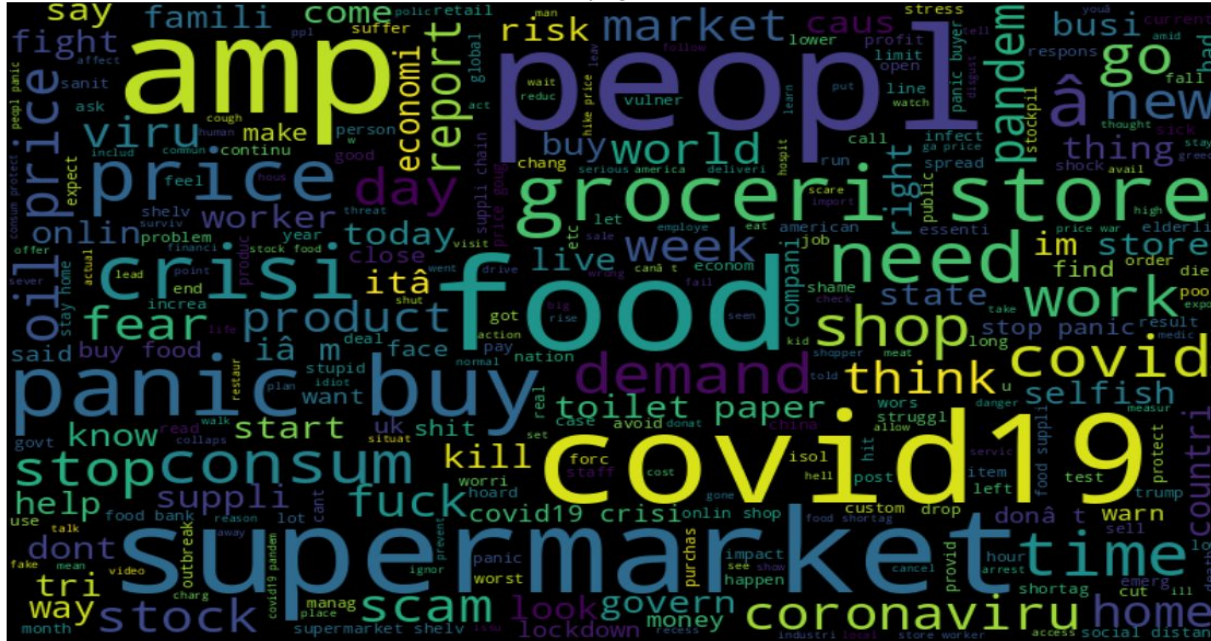


Here, like extremely positive and positive tweets an analogy has been developed the same analogy has been seen in negative types of tweets.
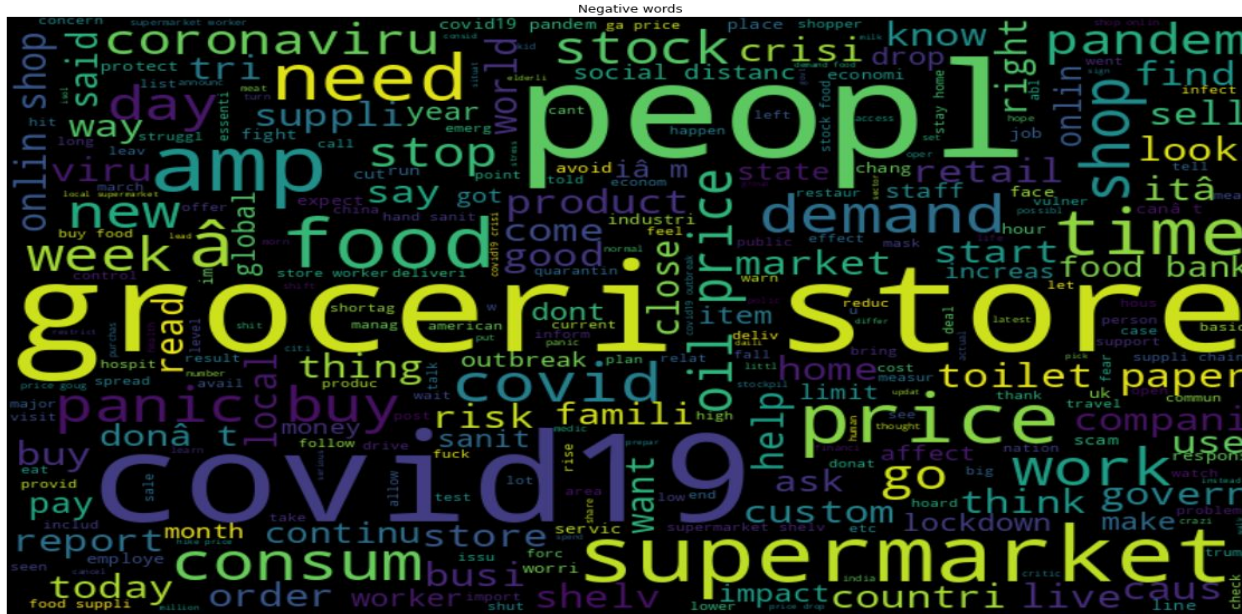
# Neutral words


Neutral words

Neutral words follows as positive type of word composition and there is no such difference.

# Extremely negative word



Extremely Negative words

From this analysis also we get to know that the same type of word has been used. There is no such difference that can be seen but the usage of negative word have increased like price, demand, crisis, oil prices and panic.

# Negative words



Here, like extremely positive and positive tweets an analogy has been developed the same analogy has been seen in negative types of tweets.

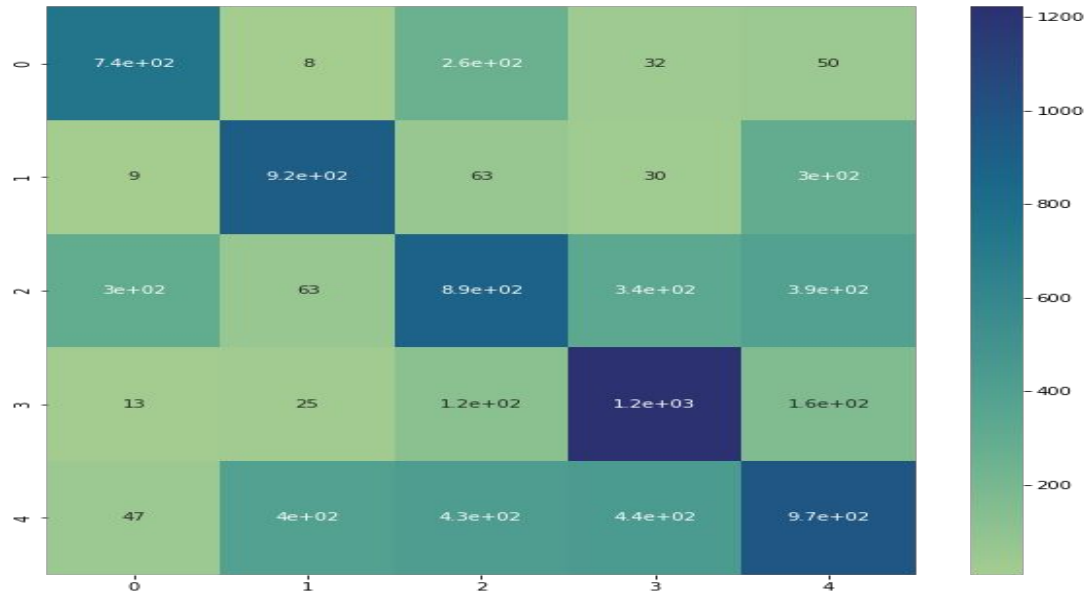# Classification

**Models Used:**

1. Naive Bayes
2. Logistic Regression
3. Random Forest
4. XGBoost
5. Support Vector Machines
6. Stochastic Gradient Descent
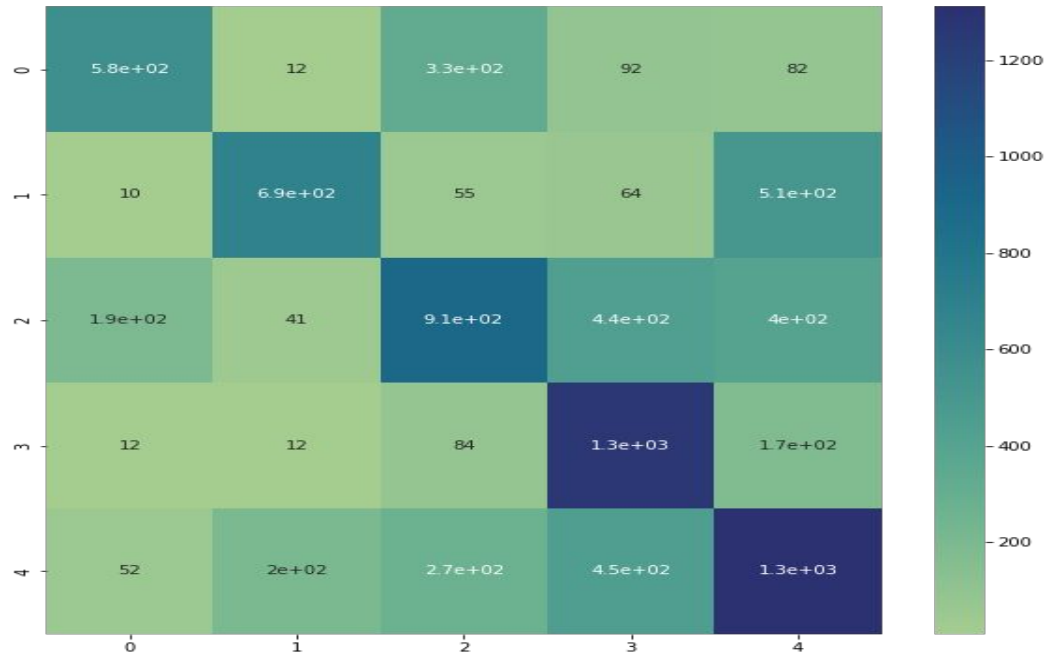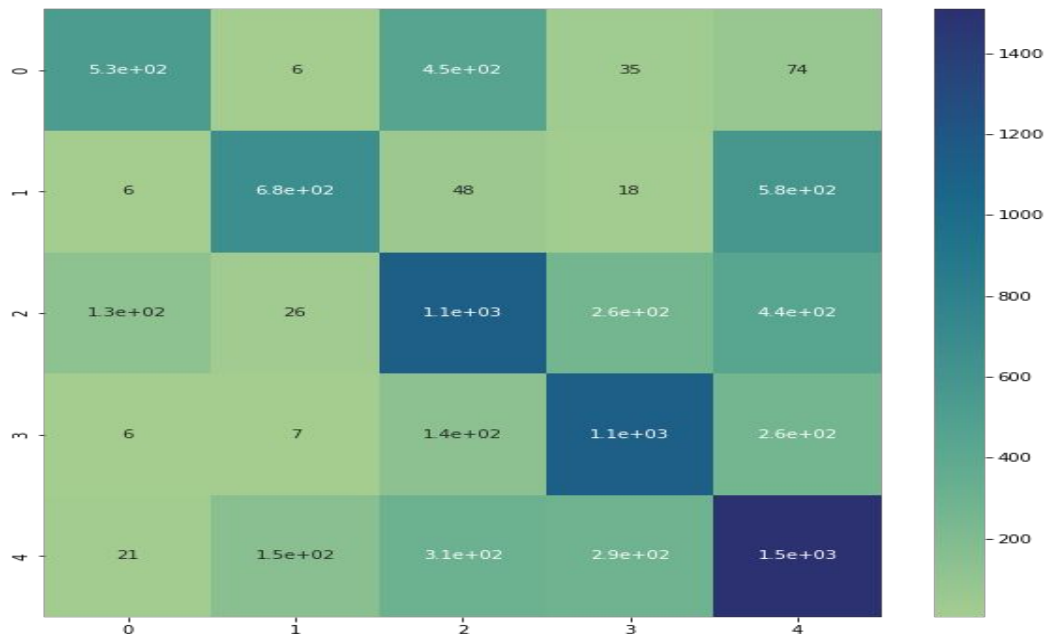7. K-Neighbors Classifier

# Naive Bayes

# K-Neighbors Classifier

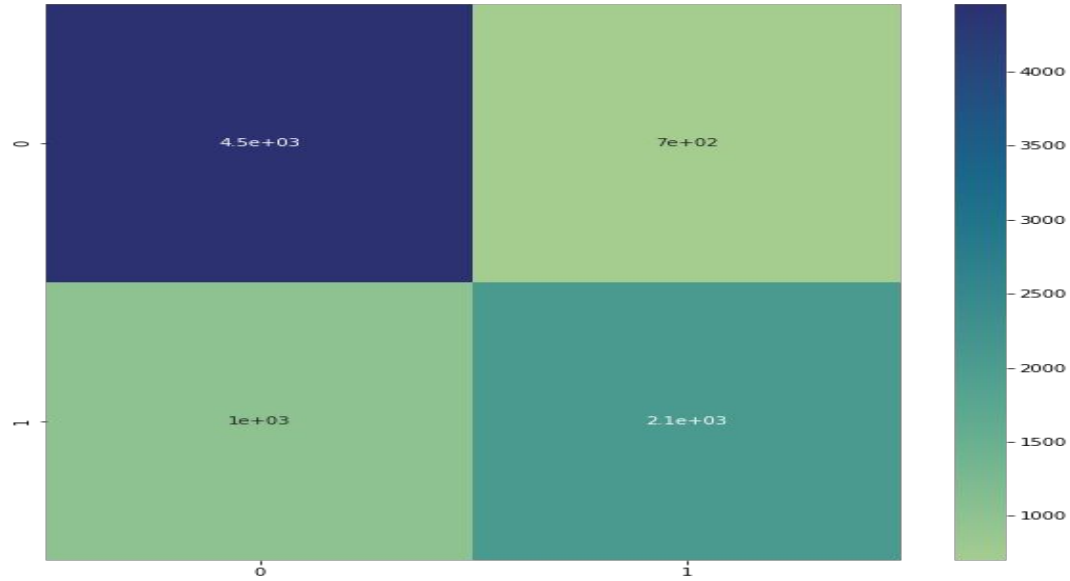# Stochastic Gradient Descent

# Extreme Gradient Boosting

# Logistic Regression

# Results Of Different Algorithm On Multi Classification

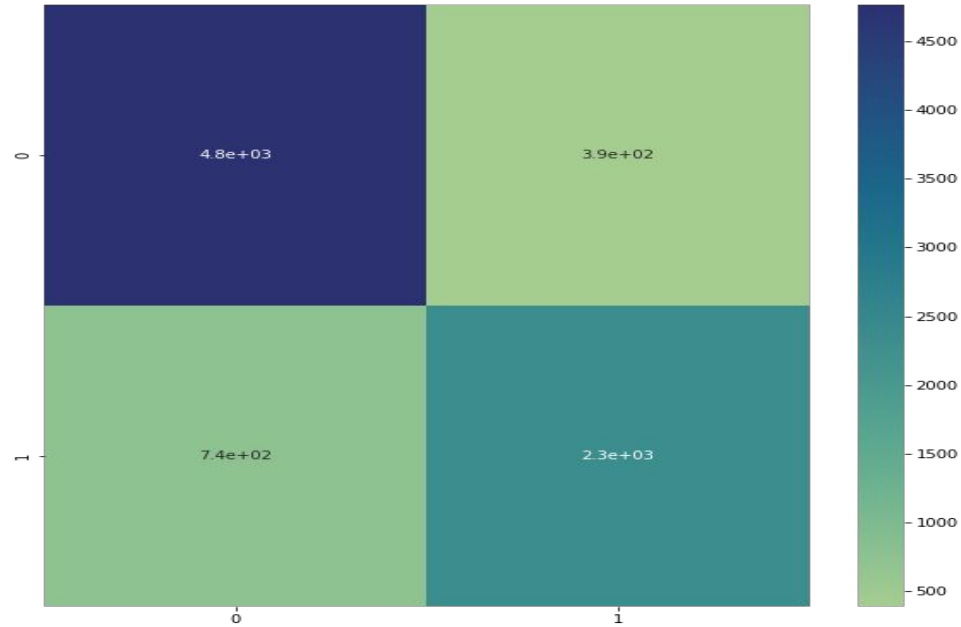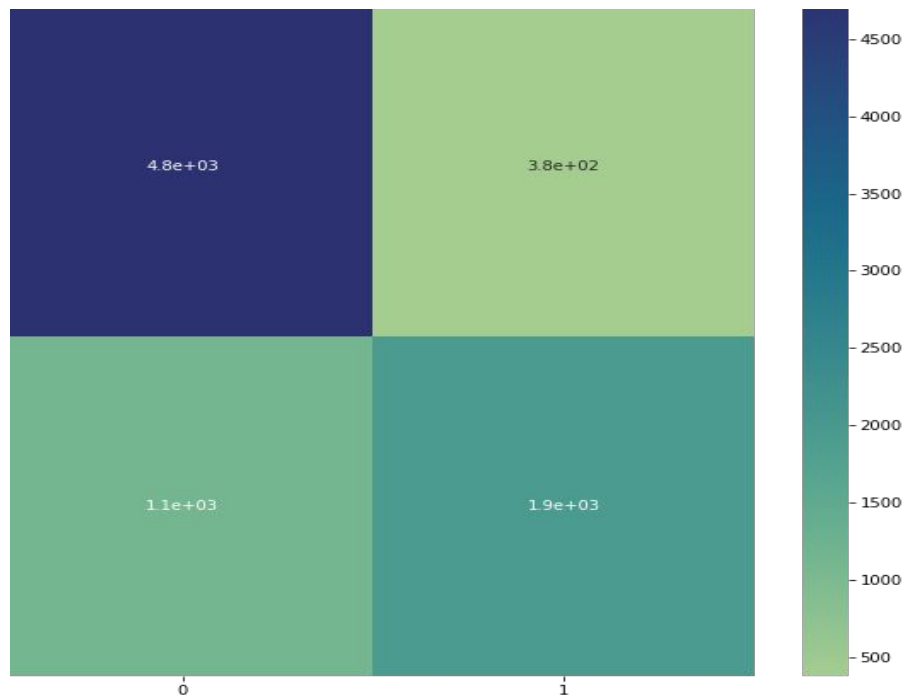| Model | Test accuracy |
|---|---|
| Logistic Regression | 0.618683 |
| Support Vector Machines | 0.602041 |
| Stochastic Gradient Decent | 0.577381 |
| xgboost | 0.576045 |
| Naive Bayes | 0.463800 |

# Naive Bayes

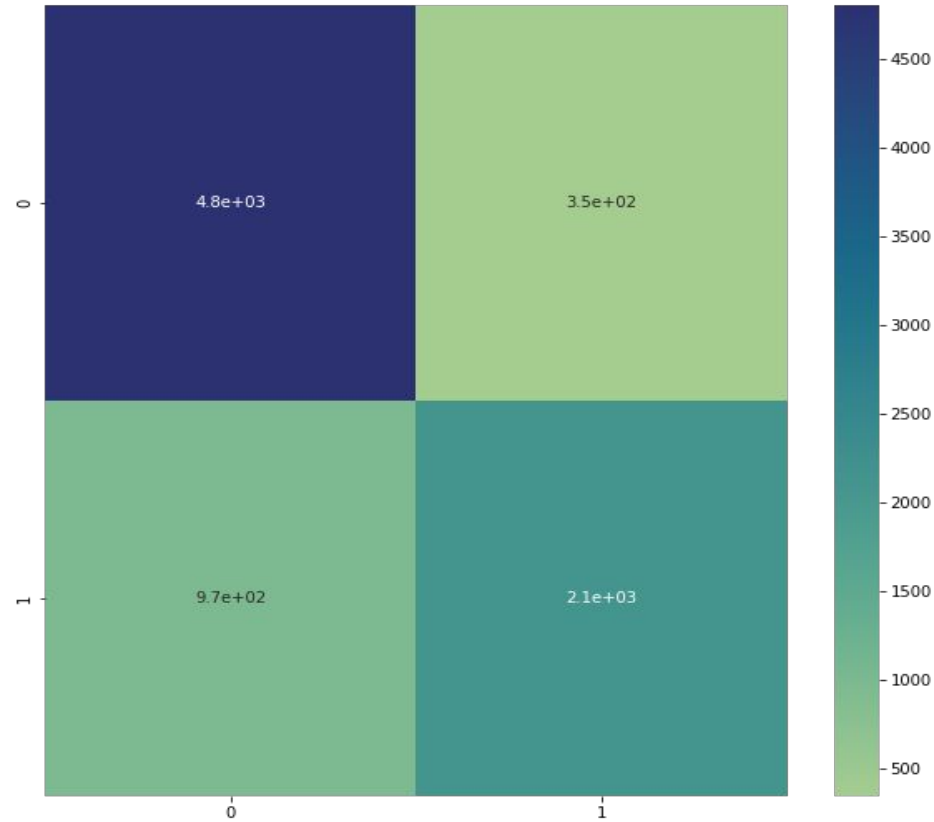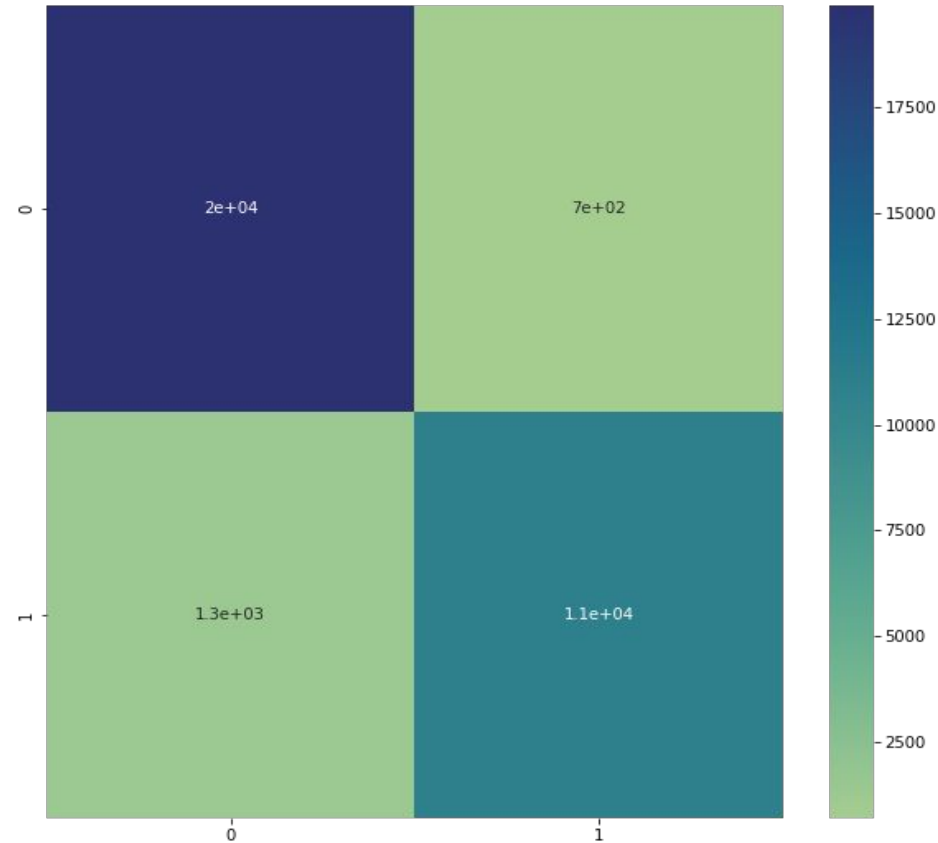# Random Forest Classifier
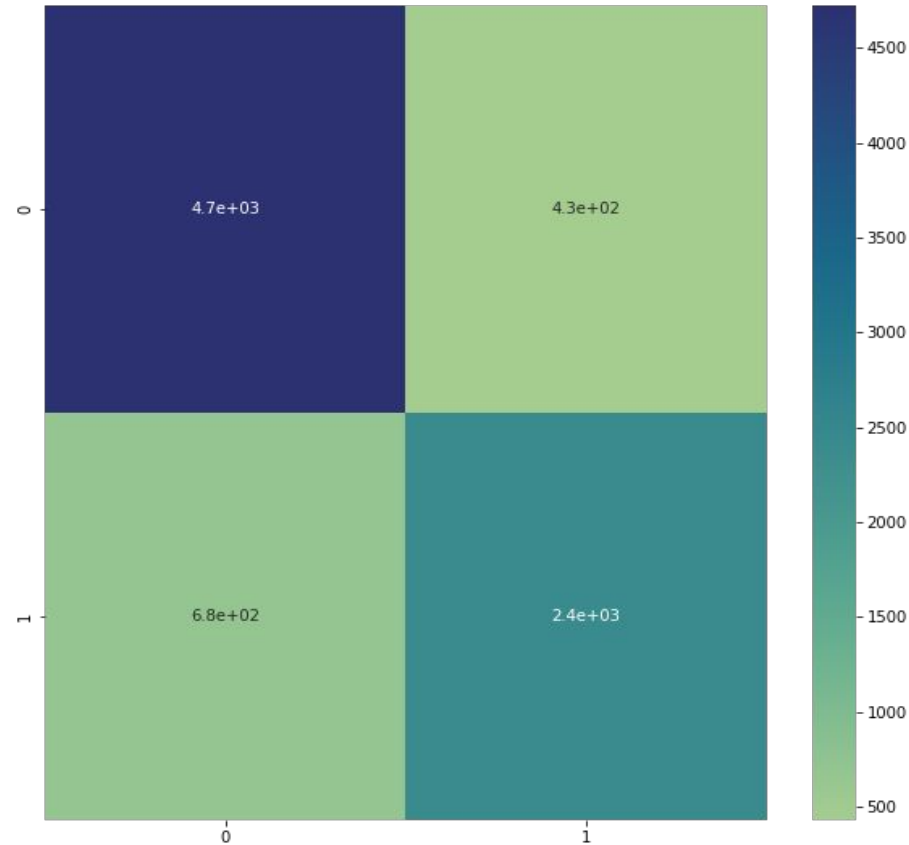
# Logistic Regression

# Extreme Gradient Boosting

# Support Vector Machines

# Stochastic Gradient Descent

# Stochastic Gradient Descent

# Results of binary classification of different algorithm

| Model | Test accuracy |
|---|---|
| Stochastic Gradient Decent | 0.865282 |
| Logistic Regression | 0.862366 |
| Support Vector Machines | 0.840136 |
| Random Forest | 0.828353 |
| XGBoost | 0.813776 |
| Naive Bayes | 0.790938 |

# ** Conclusion**

From above discussion we first evaluated the the model on the basis of multipleclassification where we have taken multiple sentiments into consideration which includes positive ,negative,neutral, extremely positive and extremely negative sentiment.

- After applying various classification algorithm we come to a conclusion that the best is logistic classification with a score of 61 percent.
- After multiple classification we divided our data into binary classification, in which extremely positive, positive and netral are taken as 0 and extremely negative and negative are taken as 1 for classification.
- The output after Deploying various Model into the system the best algorithm comes out to be Stochastic Gradient Descent.
- We have seen confusion matrix for individual cases prediction and seen that lighter part that is i has less frequency and darker part has high frequency.