# SENTIMENT ANALYSIS ON STREAMING TWEETS

*Laura Soto*

University of Miami
lms338@miami.edu

## ABSTRACT

This project analyses the sentiment classification on tweets that share the same keyword. The implementation consists of live streaming tweets that contain the given keyword into Google Cloud Storage, to then perform data preprocessing and sentiment analysis. Finally, after comparing and applying the algorithms, the obtained results were stored and analyzed. Five different keywords were analyzed, but overall a pattern between the obtained results was identified. The datasets for each keyword were obtained on three different timeframes and had different sizes. The performance on raw data and on saved results was analyzed, as well as the computational time on different datasets sizes.

## 1. INTRODUCTION

Twitter is one of the most known, and most used social media platforms. Indeed, according to Twitter Usage Statistics, on average, every minute 350,000 tweets are produced, which equivalents to 500 million tweets a day. The main usage of this social media platform is to share ideas and opinions. Therefore, analyzing the sentiments behind these tweets can provide numerous knowledge. In fact, it allows exploring how a topic is overall being observed.
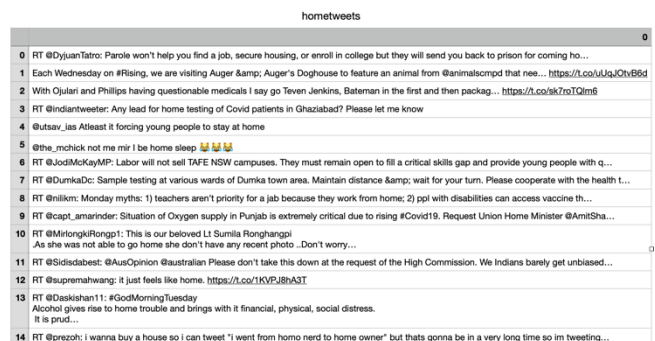
This project analyses the sentiment behind five different keywords. The initial thesis was that these keywords were clearly going to follow the predicted sentiment. In fact, the keywords used with the initial belief to have a positive sentiment behind them are 'family', 'day', and 'home'. On the other hand, the words that were predicted to have most likely a negative sentiment behind them are 'trump' and 'lockdown'.

In order to prove or contradict this initial thesis, sentiment analysis on collected datasets needed to be performed. In order to do so, accessing the Twitter API was needed, as well as stabling a Dataflow between Google Cloud Storage, and Twitter.

## 2. IMPORTING THE TWEETS

Twitter provides developer access, which allows people like me to perform algorithms and analysis on their data. That is why the first step to start this project was to set up the Twitter Developer account and to ensure access to the Twitter API. After this account was created, some personal access keys were provided. By using these personal keys, a simple python code was created to download a given amount of tweets that contained the given keyword. Even though a Tweet has different components like creation date, username, and others, for the objective of this project only the text of the tweet was imported and stored as a CSV file.

Initially, in order to import these datasets into the Google Cloud Storage, a manual command was inputted. By using this first approach, I was able to both create the tweets dataset and import them into a Bucket. On the other hand, this initial approach presented some challenges, since there is a restriction of around 2500 tweets that can be downloaded at a time, and the process of manually importing these tweets can be tedious. That is why overcoming these challenges was the main goal of my second and final approach.



**Fig. 1.** Screenshot of how the tweets were stored in the local machine before automation.

For my second and final approach, a Dataflow Job in Google Cloud was applied. This Dataflow job automates the inputting the tweets into the Google Cloud Storage. In order to do so, the initial retrieving of the tweet's python code was modified in order to push each tweet into a Pub Topic. The

Dataflow Job imports these tweets at the same time they are being exported from Twitter into the specified Bucket. The advantage of this approach is that the process of streaming tweets is automated, as well as the limit of tweets that can be downloaded.
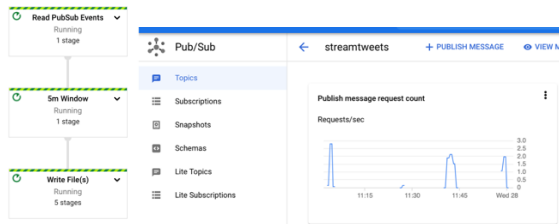


**Fig. 2.** Screenshot of Dataflow job on Pub/Sub, and the Pub/Sub

Finally, after automating and optimizing the streaming of the tweets, the result is different text files containing the imported tweets, into the Google Cloud Storage that are ready to be clean and analyzed.
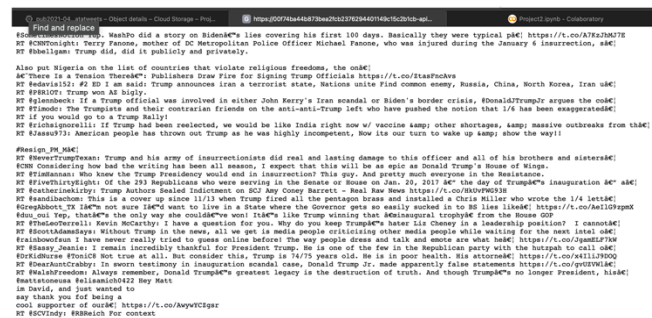


**Fig. 3.** Screenshot of the text file containing the tweets after automation

## 3. DATA PREPROCESSING

As mentioned before, the data files imported into the Google Cloud Bucket only contain the text of a tweet which is not only composed of English words. In fact, in many cases, it contains numbers, emojis, symbols, links, and others. Therefore, in order to prepare the data for a sentiment analysis algorithm all the non-alphabet and non-key words needed to be cleaned.

### 3.1 Data Cleaning

Because of the magnitude of many of these datasets, the most effective time and space complexity approach to clean this data is to use PySpark. A python code was created to perform this task. This code initially converts the dataset text file into an RDD, and then applies three different functions. The first function removes all the non-alphabet things on the text, therefore after applying this initial function the remaining RDD is composed only by string

text. The second function removes all non-English words as well as empty spaces. Finally, after cleaning the datasets, the function distinct was applied to make sure we remove all the duplicates. As a result of these functions, the remaining RDD was composed of cleaned and ready-to-be analyzed data.

## 4. SENTIMENT ANALYSIS

The main objective of this project was to be able to identify the sentiment behind a tweet. In order to achieve so, the python library Textblob was used. This library contains a well-known algorithm to identify the polarity and the subjectivity of a given text.

Since most of our datasets have a large size, applying PySpark was also observed to be the best approach. The code to implement this sentiment analysis algorithm consists of initially using the cleaned RDD to perform the analysis. Moreover, Textblob was used as a function, which was mapped on each element of the RDD. As a result of this mapping, each RDD element had a polarity value. This polarity value consisted of a float number between negative one and positive one. A zero value for this polarity represents neutral so that no sentiment could be identified by the input. On the other hand, a negative value represents a negative sentiment, and a positive value represents a positive sentiment. Furthermore, a second function was applied to identify each RDD element sentiment based on the above criteria. After applying this function a reduced map function was applied to the RDD, which allowed to have as an output the count of tweets per sentiment. This output was stored as a text file in the Bucket.
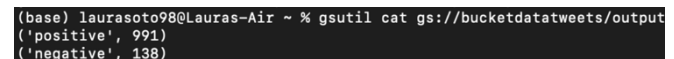


**Fig. 4.** Screenshot of the output stored after performing the sentiment analysis

Applying the described above algorithm was performed through a Dataproc Job which was executed using the following command.

```
gcloud dataproc jobs submit pyspark sentimentanalysis.py --cluster="cluster-twitter" --region="us-central1" --project="projecttwitter-311814" --
gs://bucketdatatweets/input/textfilename gs://bucketdatatweets/output/name
```

**Fig. 5.** Command executed to perform a Dataproc Job

## 5. RESULTS

For each keyword, three different sizes datasets were stored on different dates. Overall, the percentage of positive and negative sentiments through these different timeframes and sizes followed a pattern. On the other hand, the results showed that the initial prediction of the most prominent sentiment was wrong for some keywords. These results analyses were performed by importing the output of the

sentiment analysis into Google Colab and using the Matplot library and are the following.

## 5.1 Number of Tweets Vs Sentiments

The number of tweets versus the number of sentiments varies depending on each keyword. In general, contradicting the initial hypothesis, all the keywords had more positive sentients than negative. In fact, the average proportion was 70% positive 30% negative.

### 5.1.1 Lockdown

The initial belief for this keyword was that the prominent sentiment was going to be negative. On the other hand, the results show that there are more positive tweets that involve this word than negative. The three initial datasets, before the data cleaning and dropping neutral values, had a size of 500, 5000, and 2000 instances.
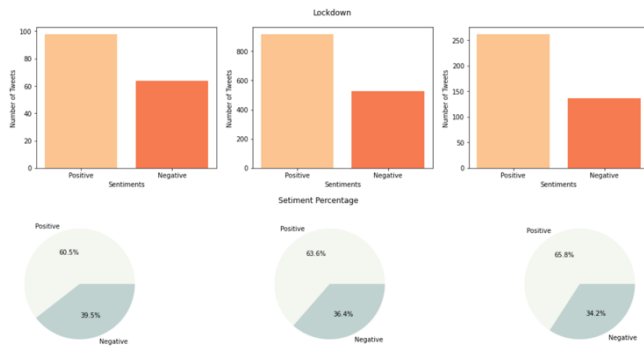


**Fig. 6.** Plots obtain from keyword 'Lockdown'

### 5.1.2 Home

This keyword had a proportion of 70% positive 30% negative, which proved the initial hypothesis. On the other hand, the still high percentage of negative sentiment when using this keyword may be caused by covid impacts during the past months. The initial datasets had a size of 5000, 3000, and 2000 tweets.
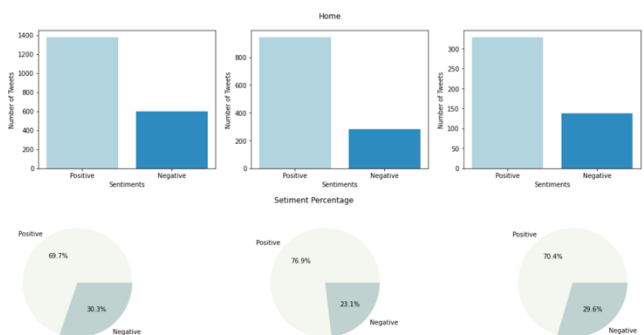


**Fig. 7.** Plots obtain from keyword 'Home'

### 5.1.3 Trump

Contrary to the belief this keyword was involved in more positive tweets than negative. On the other hand, out of the used keywords, this was the one with the highest negative percentage, with a distribution of 60% positive 40% negative. The initial dataset sizes are 3000, 5000, and 1000 respectively.
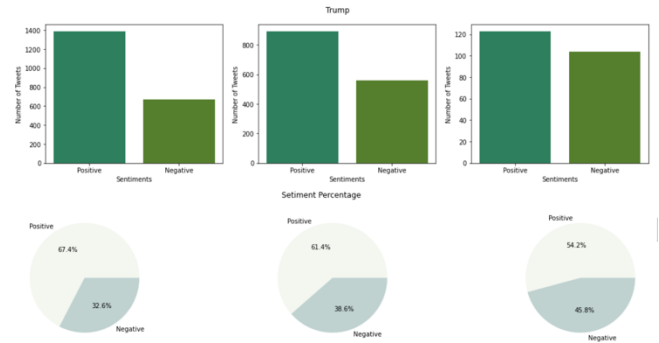


**Fig. 8.** Plots obtain from keyword 'Trump'

### 5.1.4 Family

As predicted, this keyword is most likely used on a positive tweet. Its original datasets contain 10000, 5000, and 300 tweets.
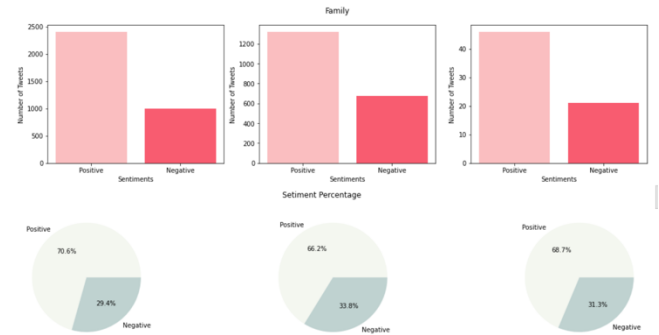


**Fig. 9.** Plots obtain from keyword 'Family'

### 5.1.5 Day

Finally, the last keyword did not only prove to be most likely used on positive tweets, but it was also the one with the highest positive rate. The initial datasets for this keyword contained 5000, 5000, and 2000 tweets.
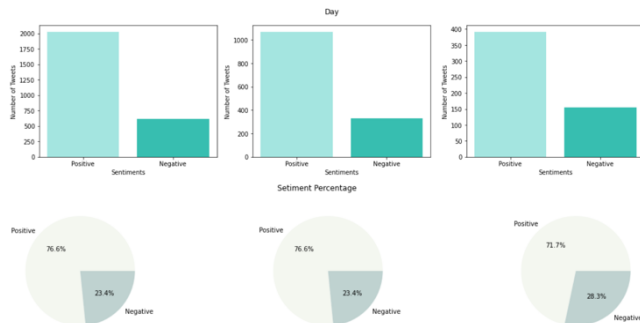
**Fig. 10.** Plots obtain from keyword 'Day'

## 5.2 Time vs Amount of Tweets

When measuring the computational time when locally downloading different sizes of datasets, we obtained the following results:
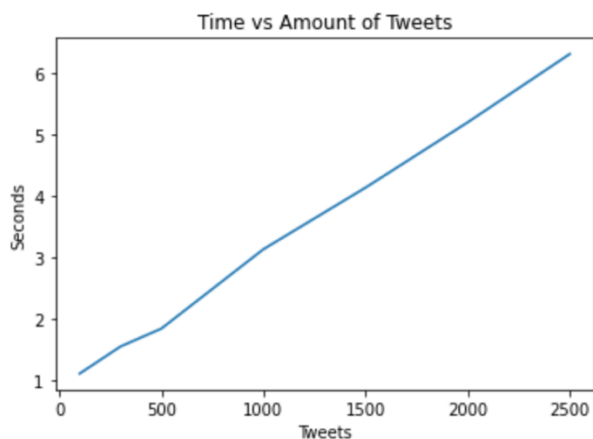


**Fig. 11.** Plot obtain from imputing the time vs the amount of tweets

This time measure was obtained using the command 'time' before compiling the code. What we can observe in this figure is that the more tweets we want to download, the longer it takes. Overall, it follows almost a straight line.
On the other hand, since Twitter API has a request limit, I wasn't able to download more than 2500 tweets per request locally.

## 5.3 Performance on Raw Data vs Saved Results

When measuring the computational time on raw data and on saved results, I used the time it took to complete the Dataproc job, which is displayed on Google Cloud Dataproc Jobs. Overall, I did not notice a big change between the performance on raw data or on saved results. In fact, the average time that it takes to completes the job is 36s no matter if it was performed on the 10.000 or the 300 datasets.

## 6. CONCLUSION

In conclusion, this project was able to perform sentiment analysis on stream tweets. Moreover, using Dataflow and Pub to live stream tweets into a Google Cloud Storage Bucket to perform data cleaning and sentiment analysis using a PySpark code applied as a Dataproc Job. The results of these executions showed in some cases different results of what it was initially predicted. Proving that for these keyword cases the sentiment percentage proportion was on average 70%-30%, prevailing the positive tweets. On the other hand, when comparing the computational time for the different dataset sizes, and for new or already used datasets, on average they all took 36s to compute. Furthermore, as expected the more tweets are downloaded the more time it takes. Finally, this project is an example of how to optimize the performance of applying a sentiment analysis algorithm to big tweets dataset by using the Google Cloud tools, and PySpark.

## 7. REFERENCES

[1
    https://github.com/Laurasoto98/SentimentAnalysisonStreamingTweets

[2]
    https://colab.research.google.com/drive/1fo2xMNmkasUP6lxae3IDOlq0GjTQAaLv?authuser=1#scrollTo=MuuA_E6E7NZg