

Original Article

# An Unsupervised Clustering Approach for Twitter Sentimental Analysis: A Case Study for George Floyd Incident

Balaji Karumanchi

Senior Software Engineer & Natsoft Corporation, 6804 Prompton Bnd, Irving, Texas, USA

Received Date: 02 May 2020

Revised Date: 12 June 2020

Accepted Date: 14 June 2020

**Abstract** - Performing sentiment analysis is vital, which can be used to find out the public review about a product or ongoing events in the world. The public can easily and efficiently express their perspectives and ideas on a wide variety of topics like events, services and brands via social networking websites. Social networks, especially Twitter, is continuously updated with public views, expressions and opinions. In this, we have performed a sentimental twitter analysis to review public opinion about the George Floyd incident using Twitter data. Text mining and sentimental analysis are used. Text mining and sentiment analysis to analyze unstructured tweet text to extract positive and negative polarity about this incident. Moreover, tweet frequency analysis has been done to view trends in public opinion across 9 days' tweet text data. We found out that the majority of the people have an attitude towards this incident by using 3 hashtags and overall data.

**Keywords** - George Floyd, sentimental twitter analysis, K-Means clustering, data mining

## I. INTRODUCTION

Catastrophic national events have a ground effect on the lives of national citizens or even on people around the globe. The consequences of such a scenario are experienced in diverse and multiple ways. The outcome and reaction of the general public can be felt through media, news reports and social media. One such footprint has been left since the murder of George Floyd on Jun 25, 2020. Since his death, protests have flared across the country demanding the arrest of the officers involved in the killing and for systemic change to put an end to police brutality. There has also been an overwhelming response on social media.

One of the best mechanisms to capture human views and emotions is to analyze the content they post on social media. People's opinions are always an important piece of information for businesses and other people to make them aware of the current trends. Since the introduction of the World Wide Web, the world has become a global village. According to [1], approximately 4.57 billion people have access to the internet, which makes up more than 57% of the world's population. And out of these 4.57 billion

people, 3.81 billion use the internet for social media. Hence public opinions on social media are the best source to know about a person's mood and opinion about a product or an event [2].

This paper is an effort to analyze the current trends and reactions of people towards this incident using statistical analysis through Twitter data. Here we have proposed two methods to achieve our goal; one is through tweet frequency analysis and second sentiment analysis using unsupervised learning to learn about the polarity of tweets.

In the first proposed method, we are going to see the trends in the top 3 hashtags selected from downloaded tweets. Based on these hashtags, we will give a strong idea about the reaction of people. This method is also vital in pre-processing data before passing it for sentimental analysis. Sentiment analysis is a technique that can be used to gauge the polarity of writing. It can help analyze the attitudes in a text related to a particular event or subject. A mathematical model is created to know about people's opinions and expressions. In [3], researchers have effectively used for political predictions, marketing strategy, e-commerce, and brand reputation management. Since we have used unlabelled data, we will be using unsupervised sentimental analysis techniques to derive results and opinions.

The remaining paper is structured as follows: In section 2, we discuss the related work on sentimental analysis. Then in section 3, we discuss methodology, which includes discussions on dataset curation, pre-processing and a brief introduction of the model. In section 4, we discuss the results. And final section concludes the paper.

## II. RECENT WORK

With the introduction of the World Wide Web and Web2.0 technology, we can see a sudden surge in consumer voices and public opinion over social media. One of the resulting emerging fields is sentiment analysis [4]. Natural Language Processing (NLP) is being widely used in opinion mining. A review of existing methods on opinion mining and sentimental analysis is done by Pang and Lee [5].



The power of social media marketing is influencing the consumer and companies as well by spreading useful information and exchange of positive or negative values. Companies are learning the customer views and discussions to support their own mission and performance goals. In [6], Rathod et al. have employed a weightage classification model based on a self-learning model to study public opinion about smartphone products. Based on words from the tweet, they categorize the tweet to be positive, negative or neutral.

In [7], the authors used hashtags, URLs and emoticons to create new tweet specific features. Becker et al. [8] proposed an online clustering framework to identify different types of real-world events and their associated social media documents. This technique can categorize similar events and non-events. Bhuvan et al. [9] proposed a sentimental model using a naïve-based algorithm to classify the polarity of the trained dataset and to validate the model to get the percentage for three categories like positive, negative, or neutral for the automotive industry.

Several methods have been deployed for sentimental analysis, including work based on Support Vector Machines [10], Naïve Bayes [11] and K-Means clustering [12]. In [13], the authors have reviewed several works for Twitter data analysis using machine learning techniques and Naïve Bayes classifier for public opinion extraction. In a more recent work [14], the authors have used graphs using the Clauset-Newman-Moore algorithm to create clusters and groupings.

The latest advent of deep learning has also been leveraged in sentiment analysis by [15]. They provide a comparison between traditional machine learning methods, polarity based methods and deep learning methods. They employed various datasets to train deep learning models hence covering a wide variety of tweets. Their results show that deep learning methods can provide up to 97% accuracy as compared to the machine learning model, which achieved a maximum of 84 % accuracy. The problem with deep learning methods (LSTM, CNN) is that they require a large dataset, very high computation power and take time to train. Hence for the specific type of sentimental analysis, traditional deep learning methods are recommended to be used.

### III. METHODOLOGY

The proposed methodology for sentimental analysis is shown in Fig. 1. Each step is explained below in detail.

#### A. Dataset

The dataset is curated using a scrapper built with Python. We retrieved 8,86,579 tweets by using 5 hashtags (blacklivesmatter, georgefloyd, icantbreathe, riots) over a span of 11 days from May 25, 2020, to Jun 04, 2020. The columns in the dataset include username, tweet text, date of tweet and link to the tweet. The dataset is further divided into four sub-datasets which are created using the top 2 hashtags (blacklivesmatter and georgefloyd). Whereas to give insight about public reaction to protests,

tweets related to protest are separated for tweet frequency analysis.

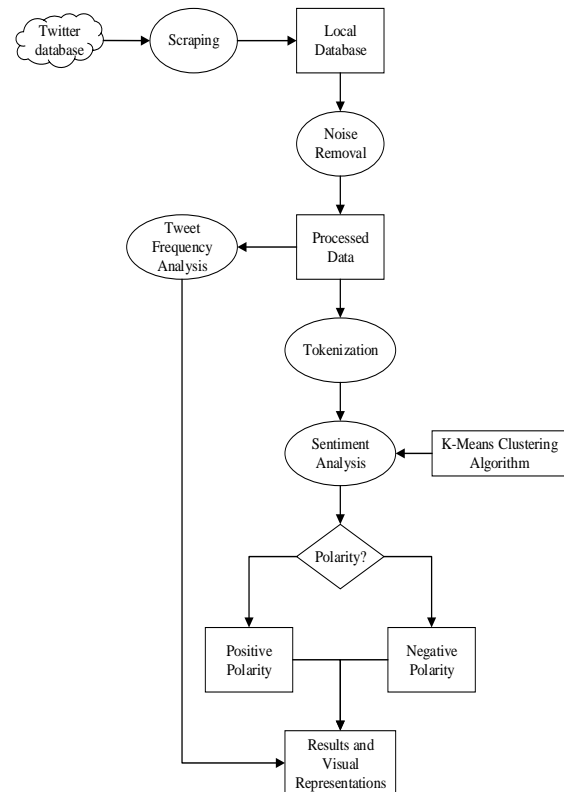


Fig. 1 Proposed methodology

#### B. Pre-Processing

The text data from tweets is raw in its original form. The text is filled with noise data like punctuation marks, emoji text, removal of duplicate data, hashtags, mentions and links. So it is necessary to clean the data to make it feasible for the next stage. Python's regular expression library (re) and Natural Language Toolkit (nltk) have been used for this purpose. Moreover, subsets are also created for tweet frequency analysis based on 3 hashtags. During this stage, the tweet feature vector using tokenization to create unigrams is also created, which is passed to the clustering algorithm to perform sentimental analysis.

#### C. Tweet Frequency Analysis

In this stage, we analyze the trend of public reaction using all data and three hashtags. The tweeter feature vector and dates are used to create groups by date. The number of dates in each group is created, and bar graphs are created. This gives insight into the public trend towards this incident.

#### D. Sentimental Analysis

We employed an unsupervised machine learning algorithm called K-Means clustering [16] for the clustering of the polarity of each tweet for the sentiment analysis. This seemed most suitable for the given problem since the data is unlabelled, and it is not possible to manually annotate 800k+ tweets. It works by taking an input number N of necessary clusters and outputs coordinates of

calculated central points of discovered clusters. Basically, the algorithm aims at minimizing the objective function known as squared error function given by:

$$J(V) = \sum_{i=1}^c \sum_{j=1}^{c_i} (||x_i - v_j||)^2$$

Where  $V$  corresponds to the current cluster,  $x_i$  is the centre of the current cluster, and  $v_j$  is the position of data points. It is an iterative algorithm, in which in the first step,  $N$  random data points are chosen as coordinates for the centre of clusters. In each step, all using Euclidean distance points are assigned to their closest centroid. Then new coordinates of centroids are calculated by taking the mean of coordinates of all the data points in that cluster. These steps are iterated until a minimum value of mean squared error between points assigned to centroids is achieved.

Since the output of the K-Means algorithm is the distance of data points from centroids, we have to convert it to a polarity score. For this purpose, the distance is multiplied by the inverse of the closeness score. After this step, we will have a dictionary containing a word and a weighted sentimental score. Later then, we calculated the term frequency-inverse document frequency score (tfidf). This is a numerical statistic that points out how much important a word is to a sentence/document and is used as a weighting vector for information retrieval or text mining. To achieve this, we used the sklearn library. After this step, we have 2 vectors for each sentence; one vector containing weighted sentimental score and the other one tfidf score. Finally, these 2 vectors are multiplied to achieve the final polarity of the sentence being positive or negative.

#### IV. RESULTS AND DISCUSSIONS

##### A. Tweet Frequency Analysis

The analysis for the whole data is given in Fig. 2. It shows the number of tweets per day spanning over a period of 9 days. It can be seen from the figure that the first tweet arose on early May 26, 2020. And it starts rising till May 28, 2020. Later it remains almost constant and then starts decreasing from Jun 02, 2020.

The tweet frequency analysis for hashtag George Floyd is shown in Fig. 3. A similar trend as seen in the whole data plot can be seen here. The plot first increases to May 25, 2020, then remains almost constant till Jun 04, 2020, and then starts decreasing. The tweet frequency analysis for the hashtag

The tweet frequency analysis for hashtag protest is shown in Fig. 5. We can see that during the days when protests were at their peak, we can see a similar trend in tweet data also. The protests were at their peak during these days. Then we can see a drop in data as there were fewer protests during that time.

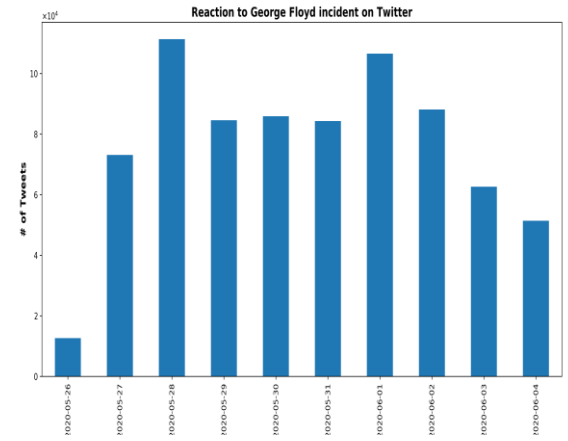


Fig. 1 Tweet Frequency Analysis (GeorgeFloyd)

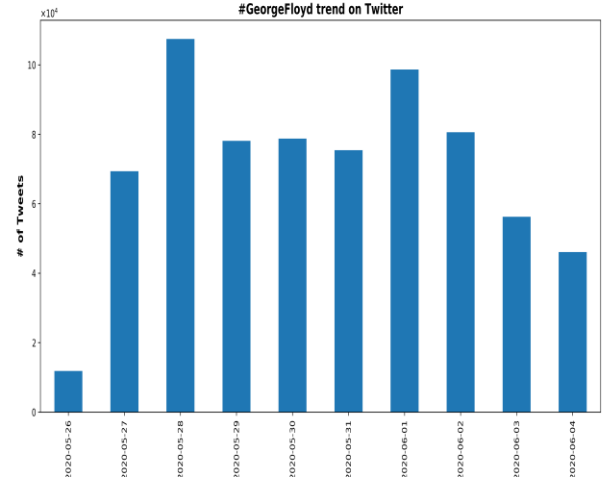


Fig. 2 Tweet Frequency Analysis

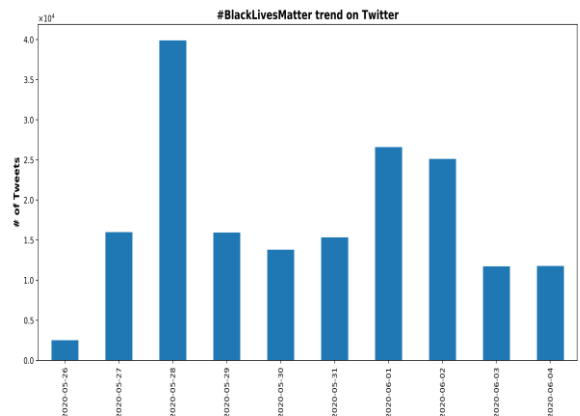


Fig. 3 Tweet Frequency Analysis (BlackLivesMatter)

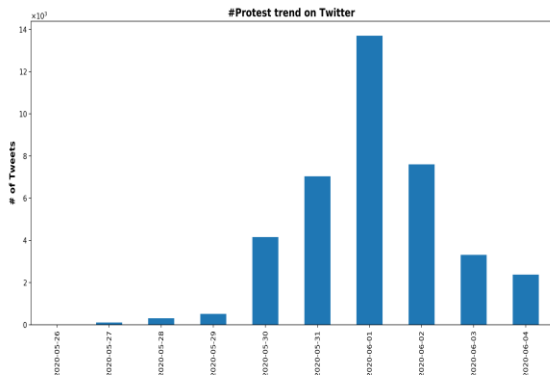


Fig. 4 Tweet Frequency Analysis (#protest)

**B. Sentimental Analysis**

After the implementation of sentiment classification, we got the values of sentiment distribution of negative and positive polarities. Each tweet is analyzed to be a positive or negative tweet based on a query term and polarity clustering. Some tweet samples of positive and negative polarities are given in Table I.

Looking at the table, we can conclude that tweets having harsh words or too much anger are marked as negative, whereas tweets with good words or very less harsh words are marked as positive. Moreover, it is obvious from Fig. 6, Fig. 7 and Fig. 8 that the majority of the people have a positive attitude towards this matter which means that majority of the tweets do not show hatred but still want justice for George Floyd and have a positive attitude towards this incident.

Table 1. Tweet Samples and their Polarity

Tweet	Polarity
can you imagine feeling so empowered that you allow yourself to be recorded while you commit murder in broad daylight in front of several witnesses??? #blacklivesmatter	Positive
rest in peace # icantbreathe #georgefloyd	Positive
it is heartbreaking & terrifying living in a country where I wouldn't call the police if i needed help, in fear that someone in my family could be wrongfully killed. #minneapolis #philandocastille #centralparkkaren #minneapolispolice #blacklivesmatter #icantbreathe	Negative
im so at a loss for words. this shit is just so hard to watch. #georgefloyd #rip! I hope all them cops a painfully slow death fr	Negative

It can be seen that for all data, we have 72.2 % positive and 27.3 % negative tweets. For #GeorgeFloyd, we can see that 67.9% percent of people have positive and 32.1 % of people have negative polarity. Finally, for

#BlackLivesMatter, it can be seen that 70% of people have positive and 30 % of people have a negative attitude.

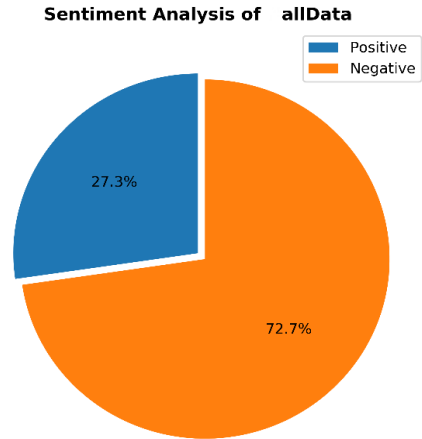


Fig. 2 Sentiment distribution of all tweets

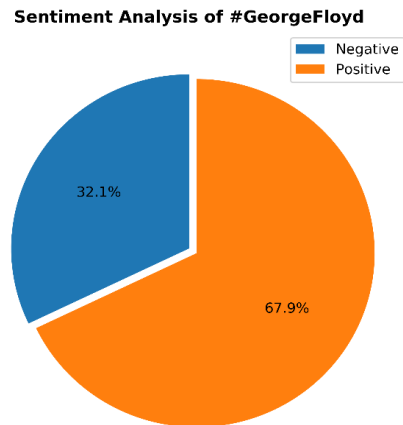


Fig. 3 Sentiment analysis of tweets (GeorgeFloyd)

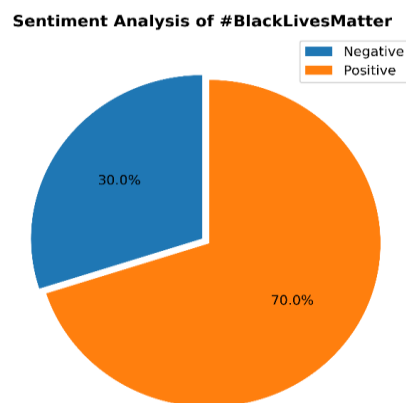


Fig. 4 Sentiment distribution of tweets (BlackLivesMatter)

## V.CONCLUSION

Sentiment analysis is one of the most attractive and shining fields of text and data mining with vast applications in numerous sectors. In this paper, we have shown how a dataset can be curated in its pre-processing stages. Then we have done Tweet Frequency Analysis to monitor the trend towards this incident. Then we did the sentimental analysis using the K-Means clustering algorithm to know about the polarity of tweets, i.e. positive or negative. The results showed that the Twitter response after 2 days of the incident was at its peak, then it declined slowly. Moreover, we have shown that, on average, 70.2 % of people have a positive attitude towards this incident. The polarity results were consistent in all three samples of the dataset.

In future work, these results can be further improved by labelling the dataset. Then using this labelled dataset, a comparison between different machine learning algorithms can be made. Moreover, deep learning algorithms like CNN, LSTM or RNN can also be employed for improvement in polarity clustering

## REFERENCES

- [1] J. Clement, Global Digital Population as of Statista.Com, (2020) [Online]. Available: <https://www.statista.com/statistics/617136/digital-population-worldwide/>, Accessed On (2020).
- [2] E. Cambria, B. Schuller, Y. Xia, and C. Havasi, New Avenues in Opinion Mining and Sentiment Analysis, *Ieee Intelligent Systems*, 28(2) (2013) 15-21.
- [3] B. Liu, Sentiment Analysis and Opinion Mining, *Synthesis Lectures On Human Language Technologies*, 5(1) (2012) 1-167.
- [4] N. Majumder Et Al., Improving Aspect-Level Sentiment Analysis with Aspect Extraction, *Arxiv Preprint Arxiv:2005.06607*, (2020).
- [5] B. Pang And L. Lee, *Opinion Mining and Sentiment Analysis. Foundations and Trends in Information Retrieval*, (2008).
- [6] J. A. Rathod, S. Vignesh, and A. J. Shetty, Sentiment Analysis of Smartphone Product Reviews Using Weightage Calculation, in *Advances in Computing and Intelligent Systems: Springer*, (2020) 427-437
- [7] A. Ritter, S. Clark, And O. Etzioni, Named Entity Recognition in Tweets: An Experimental Study, in *Proceedings of the Conference on Empirical Methods in Natural Language Processing, Association For Computational Linguistics*. (2011) 1524-1534.
- [8] H. Becker, M. Naaman, And L. Gravano, Beyond Trending Topics: Real-World Event Identification on Twitter, In *Fifth International Aaai Conference on Weblogs And Social Media*, (2011).
- [9] Malladihalli S Bhuvan, Vinay D Rao, Siddharth Jain, T S Ashwin, And R. M. R. Guddeti, Semantic Sentiment Analysis Using Context-Specific Grammar, Presented at the *International Conference on Computing, Communication and Automation*.(2015).
- [10] C. C. Chen And Y.-D. Tseng, Quality Evaluation of Product Reviews Using An Information Quality Framework, *Decision Support Systems*, 50(4) (2011) 755-768.
- [11] M. Gayathri, S. S. Nisha, And M. M. Sathik, Twitter Sentiment Analysis Using Naive Bayes Classification, *Studies In Indian Place Names*, 40(71) (2020) 1473-1478.
- [12] H. Suresh, An Unsupervised Fuzzy Clustering Method for Twitter Sentiment Analysis *International Conference on Computation System And Information Technology For Sustainable Solutions (Csitss)*, (2016) 80-85.
- [13] T. Vaseeharan And A. Aponso, Review on Sentiment Analysis of Twitter Posts About News Headlines Using Machine Learning Approaches And Naïve Bayes Classifier, in *Proceedings of the 12th International Conference on Computer and Automation Engineering*, (2020) 33-37.
- [14] W. Ahmed, J. Vidal-Alaball, J. Downing, And F. L. Seguí, Covid-19 And The 5g Conspiracy Theory: Social Network Analysis of Twitter Data, *Journal of Medical Internet Research*, 22(5) (2020).
- [15] Y. Chandra And A. Jana, Sentiment Analysis Using Machine Learning and Deep Learning, in *2020 7th International Conference On Computing for Sustainable Global Development (Indiacom)*, (2020) 1-4.
- [16] J. Macqueen, Some Methods for Classification and Analysis of Multivariate Observations, in *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics And Probability*,. Oakland, Ca, Usa, 1(14) (1967) 281-297.
- [17] B.Srinivasa Rao, S.Vellusamy Raddy, A Hard K-Means Clustering Techniques for Information Retrieval from Search Engine Ssrg *International Journal of Computer Science and Engineering* 4(2) (2017).