

Understanding the Correlation between COVID-19 Related Tweets and National Policies through Sentiment Analysis

Yue Liu¹ and Senqi Zhang²

Abstract—During the COVID-19 pandemics, billions of people are impacted during this global crisis. However, in the United States, the unavailability of public data in early 2020 makes it difficult for smaller groups to study and analyze the spread of the virus as well as the mental health of people during isolation. This paper introduces a new way to monitor the pandemics as well as a new approach for the government to understand the public using social media.

I. INTRODUCTION

The COVID-19 pandemic has rapidly impacted over 200 countries, areas, and territories. On December 8, 2019, several cases of COVID-19 were first reported in Wuhan, China [1]. Just one and a half months later, on January 23, Wuhan, a city of more than 11 million, was locked down to prevent the spread [2]. In March, while China starts to see its turning point, COVID-19 is just starting to become a global threat. On March 13, President of the United States, Donald Trump declared a national emergency, and people are encouraged to practice social distancing [2]. As of May 4, in the United States, there are 1.21 Million reported cases and nearly 70 thousand deaths [3].

In a difficult time like this, we would like to know COVID-19 has impacted the sentiment of the public and how sentiment is correlated with National Policies. In light of the deteriorating situation in the United States, the discussions related to “COVID19” and “Donald Trump” on Twitter has increased drastically. Therefore, we choose to analyze the streaming data on Twitter between March 14, 2020, and April 14, 2020, for this study.

In this study, we calculate the average sentiment score for each day by implementing Google NLP API and study the fluctuation of the average sentiment score. Then we related interesting changes to important national policies that are announced on that day. Then to further assure the correlation between the sentiment and the national policies, we use LDA and WorldCloud to testify whether the popular topic on Twitter is indeed related to the national policies.

II. RELATED WORK

Since the beginning of the COVID-19 pandemic, researchers have worked diligently in studying various aspects of fighting COVID-19. Data mining research focuses on studying the impacts and implication, including protection against the virus[5], new treatment method[6][7], real-time tracking of the situation[8], the effectiveness of social distancing[9][10], and spreading of misinformation[12][13].

Research on people’s sentiment has been conducted on the Chinese largest social media Weibo[10]. In our study, we want to focus on how public policies affect people’s sentiment over time and provide insights into new ways of understanding and tracking pandemics.

III. METHODOLOGY

A. Data Collection

All of our data are collected using Python Package Tweepy [14]. From March 14 2020 to April 14 2020, the data is collected every day from 10 a.m to 12 p.m EST. We choose this specific time period because statistics shows Twitter have highest volume of user activity from 10 a.m. to 1p.m[17]. About 200 thousands to 300 thousands tweets in English are collected each day and a total of 7 million tweets are collected during this period. The keyword we use to filter the streaming data are “Coronavirus”, “COVID19”, and “Trump” and all retweets are filtered. Noted that “Trump” is chosen for keyword as President Trump’s usage of Twitter has attracted worldwide attention ever since the beginning of his presidency. It has become somewhat a part of the way Trump Administration connect with publics. We believe adding this keyword will help in analyzing the public sentiments towards the national policies.

B. Data Preprocessing

Streaming tweets contain a lot of information. The first thing we do is disable collecting retweets because retweets are a reoccurring sample, we want to treat every tweet we collect with the same weight in our sentiment analysis. Then we only extract the “text” field for our sentiment analysis. At first, we plan to use other fields such as “geo” and “location” to do a location-based sentiment analysis with more dimensions, but the reality is most people either do not have the location service enabled for Twitter or they can simply customize the location to wherever they like, thus disqualifying the validity of the data. After that, we have a clean JSON file with one dictionary each line that looks like “Text:”, “COVID-19”. However, some sentences still contain extra apostrophe and quotation marks inside the dictionary structure. This causes our algorithm to report bugs and having difficulty working with the data structure. Therefore, we removed any additional apostrophe and quotation marks inside the dictionary. At last, we eliminate the weblinks and Unicode in the text.

¹University of Rochester, Email: yliu165@u.rochester.edu

²University of Rochester, Email: szhang71@u.rochester.edu

time have a noticeable lower sentiment score across the board which indicates the public's dissatisfaction with the Trump Administration, or at least, Trump himself.

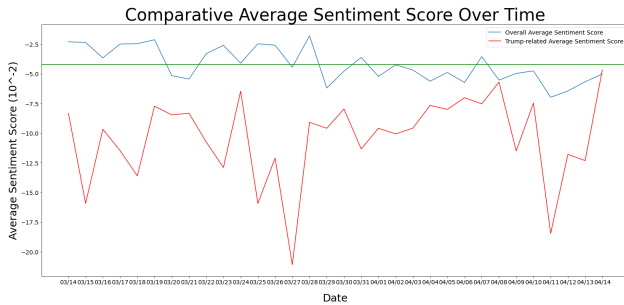


Fig. 11: Comparative Sentiment Score Graph

V. CONCLUSION AND FUTURE WORK

We have presented a study focusing on discovering a correlation between COVID-19 related national policies and individual's sentiment behavior on social media during the ongoing COVID-19 pandemic. Utilizing Word Cloud and LDA analysis, we discovered an overlap between heated topic on twitter and national policies announced that day. By analyzing the fluctuation of sentiment scores over time, we further reassured the correlation between public sentiment and national policies. In addition, by comparing the sentiment score of tweets containing "Trump" and "Coronavirus" to all the tweets, we discovered a distinct dichotomy between the two curves. Results from these analyses show analyzing sentiment of users on social media on a large scale can be used as a metric to measure the popularity of a national policy among people. We believe we have introduced a new way of monitoring pandemics as well as a new approach for the government to connect and understand the public.

In the future, we want to try out different sentiment analysis tools such as VADER and see if we can have more accurate and credible results. We want to implement the location of the tweets to our algorithm so we can visualize the sentiment based on location using social media. In addition, we hope to develop a predictive algorithm based on data from social media that can be used for prediction of the spread of pandemics. Last but not the least, we will continue working on improving our visualization of LDA because it is not as intuitive as we expected in its current state.

REFERENCES

- [1] Timeline, the central events regarding Coronavirus in Wuhan., <https://www.nytimes.com/article/coronavirus-timeline.html>
- [2] How the Coronavirus Pandemic Unfolded: A Timeline. , https://www.sohu.com/a/373647917_14988
- [3] Coronavirus Disease 2019 (COVID-19), Cases, Data, Surveillance., <https://www.cdc.gov/coronavirus/2019-ncov/cases-updates/>
- [4] Twitter by the Numbers: Stats, Demographics Fun Facts., <https://www.omnicoreagency.com/twitter-statistics/>.
- [5] Zheng, Y., Ma, Y., Zhang, J. et al, *COVID-19 and the cardiovascular system*, Nat Rev Cardiol 17, 259–260 (2020).

- [6] Dong, L., Hu, S., Gao, J, *Discovering drugs to treat coronavirus disease 2019 (COVID-19)*, Drug Discoveries Therapeutics 14,58-60 (2020), 259–260 (2020).
- [7] Mehta, P., McAuley, D., Brown, M. et al, *COVID-19: consider cytokine storm syndromes and immunosuppression.* ,The Lancet 395,1033-1034 (2020).
- [8] Dong, E., Du, H., Gardner, L, *An interactive web-based dashboard to track COVID-19 in real time*,The Lancet Infectious Disease 20,533-534 (2020)
- [9] Anderson, M., Heesterbeek, H,*How will country-based mitigation measures influence the course of the COVID-19 epidemic*,The Lancet 395,931-934 (2020)
- [10] Hellewell, J., Abbott S, Gimma,*Feasibility of controlling COVID-19 outbreaks by isolation of cases and contacts*,The Lancet Global Health 8, 488-496(2020).
- [11] Zhao, Y., Xu, H,*Chinese Public Attention to COVID-19 Epidemic: Based on Social Media.* medRxiv
- [12] Kouzy, R., Jaoude, J., Kraitem, A. et al,*Coronavirus Goes Viral: Quantifying the COVID-19 Misinformation Epidemic on Twitter*,Cureus. 2020 Mar; 12(3): e7255
- [13] Singh, L., Bansal, S., Bode L. et al,*A first look at COVID-19 information and misinformation sharing on Twitter*,arXiv:2003.13907 [cs.SI]
- [14] Tweepy Documentation, <http://docs.tweepy.org/en/latest/>
- [15] WordCloud Documentation, <https://amueller.github.io/>
- [16] White House Briefing Statements, <https://www.whitehouse.gov/briefings-statements/>
- [17] The Biggest Social Media Science Study: What 4.8 Million Tweets Say About the Best Time to Tweet, <https://buffer.com/resources/best-time-to-tweet-research>