

UNDERSTANDING THE PERCEPTION OF COVID-19 POLICIES BY MINING A MULTILANGUAGE TWITTER DATASET

Christian E. Lopez^{1,2*}, Malolan Vasu¹ and Caleb Gallemore³

¹ Computer Science Department, Lafayette College, Easton, PA 18042

² Mechanical Engineering Department, Lafayette College, Easton, PA 18042

³ International Affairs Program, Lafayette College, Easton, PA 18042

Abstract

The objective of this work is to explore popular discourse about the COVID-19 pandemic and policies implemented to manage it. Using Natural Language Processing, Text Mining, and Network Analysis to analyze corpus of tweets that relate to the COVID-19 pandemic, we identify common responses to the pandemic and how these responses differ across time. Moreover, insights as to how information and misinformation were transmitted via Twitter, starting at the early stages of this pandemic, are presented. Finally, this work introduces a dataset of tweets collected from all over the world, in multiple languages, dating back to January 22nd, when the total cases of reported COVID19 were below 600 worldwide. The insights presented in this work could help inform decision makers in the face of future pandemics, and the dataset introduced can be used to acquire valuable knowledge to help mitigate the COVID-19 pandemic.

Link for dataset: https://github.com/lopezbec/COVID19_Tweets_Dataset

1 Introduction

The Coronavirus Disease of 2019, known as COVID-19, is a rapidly spreading disease caused by the Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV2). The COVID-19 is now considered a pandemic that has affected countries in all inhabited continents. Since the first cases of COVID-19 reported in Wuhan, China, in December 2019, the number of fatalities worldwide has increased rapidly. Due to its high infection and death rate, governments have implemented a wide range of policies aimed at mitigating the spread of this virus and its impact. Such actions began with the Chinese government order to quarantine Wuhan on January 23rd, 2020, to, most recently, multiple countries declaring state of emergency and implementing strict quarantine and social distancing measurements (e.g., US, Italy, Argentina, Spain).

Most government leaders have implemented measures to incentivize, and in some cases enforce, “social distancing” to reduce the spread of COVID-19. These measures have resulted in the cancelled entertainment events, closures of schools and colleges, and businesses reducing hours of operation, implement telecommuting, or close altogether. There is no doubt the pandemic and the measures set in place to mitigate it have and will continue to drastically impact the lives of millions. As this pandemic and the responses to it are unprecedented, however, we are likely to be surprised by how people respond.

Since the early stages of the disease, people have expressed their opinion and shared information, as well as misinformation, about it via social media platforms, such as Twitter. As COVID-19 spreads to other countries and governments try to mitigate its impact by implementing counter measures, people have also used social media platforms to express their opinion about the measures themselves, the leaders implementing them, and the ways their lives are changing. The use of social media, such as Twitter, as platforms to express opinions and share information about COVID-19, will only continue to grow, precisely because of the “social distancing” measures set in place to mitigate it.

Policymakers could mine this social media data to explore popular discourse about the pandemic and the measures set in place to mitigate it. We plan to analyze a corpus of tweets that relate to COVID-19 with the objective of identifying common responses to the pandemic and how these responses differ across time, countries, and policies. Moreover, insights as to how information and misinformation about this pandemic and the policies are transmitted are presented. Finally, we introduce and share with the research community a dataset of tweets collected from all over the world, in multiple languages, dating back to January 22nd when the total cases of reported COVID-19 were below 600 worldwide. Here, we describe and present descriptive statistics of this dataset, and explain our data collection methods and intended analyses.

2 Dataset Description

The dataset presented is being continuously collected using the Twitter API. The dataset presented here (v1.3) covers Jan 22, 2020 to Apr 05, 2020 and contains 40,067,960 tweets. The keywords used for search tweets are: *virus* and *coronavirus* since January 22, 2020, *ncov19* and *ncov2019* since

*Corresponding Author. 569 Rockwell Integrated Science Center, Lafayette College, Easton, PA 18042, lopezbec@lafayette.edu

February 26, 2020, and *covid* since March 7, 2020.

The average daily number of tweets collected on dataset v1.3 was 534,239.47 (SD=764,426.85, Mdn=159,813.0). In the first few months, the number of tweets collected increased steadily from 724,877 in January and 3,084,728 in February, 27,414,279 in March, to 8,844,076 in April. Table 1 shows the summary statistics for the daily number of tweets collected each month and on the first 05 day(s) of April.

Table 1: Statistics on number of daily tweets per month

Month	Mean	SD	Median
Jan	72,487.7	22,964.94	82,153.0
Feb	106,369.93	40,156.08	86,126.0
Mar	884,331.58	884,625.83	296,543.0
Apr	1,768,815.2	461,889.93	1,879,773.0

While the dataset contains tweets from 63 languages, only English-language tweets were collected from 22 January to 30 January 2020. English-language tweets remain the most prominent in the dataset accounting for 58.74% of the total. Figure 1 presents the distribution of collected tweets by language.

Information about retweets and likes was also collected. Figure 2 presents the distribution of collected tweets' retweets over the observed period. While the overall level of retweeting appears to have declined in February, retweets rose abruptly as the crisis ramped up in Europe in late February and early March.

In addition to language, a relatively small percentage of

Table 2: Distribution of Tweets per Language

Language	Number of Tweets	Percentage
English	23,315,199	58.74%
Spanish; Castilian	6,662,566	16.79%
Portuguese	2,648,300	6.67%
Bahasa	1,938,992	4.89%
French	1,554,369	3.92%
Thai	680,075	1.71%
Italian	591,873	1.49%
Japanese	350,149	0.88%
Turkish	291,978	0.74%
Tagalog	288,697	0.73%
Hindi	279,104	0.70%
Catalan; Valencian	238,538	0.60%
German	235,563	0.59%
Dutch; Flemish	112,906	0.28%
Other	502,355	1.27%

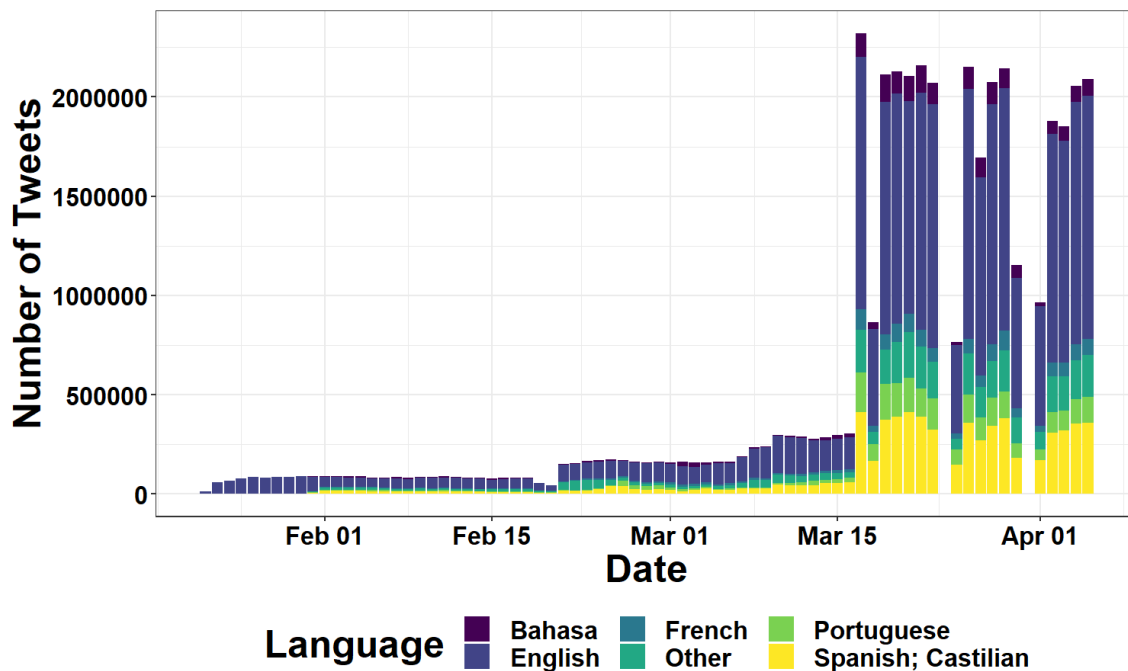


Figure 1: Collected tweets, by language, as of Apr 05, 2020.

Table 3: Summary statistics for collected tweets

Number of Geolocated Tweets	14,116
Maximum Observed Retweets	1,064,693
Median Observed Retweets	73
25th Percentile of Observed Retweets	1
Maximum Observed Likes	8732
Median Observed Likes	0
25th Percentile of Observed Likes	0

3 Dataset Accessibility

Figure 2: Distribution of observed retweets, on a log10 scale, across the observed period.

the collected tweets contain geolocations. Figure 3 presents the locations observed as of Apr 05, 2020. Table 3 provides numerical summaries of the observed tweets number of retweets, likes and geolocation information.

The dataset v1.3 was realized on Apr 05, 2020. The dataset is available on Github at: https://github.com/lopezbec/COVID19_Tweets_Dataset. The dataset is released in compliance with the Twitter’s Terms & Conditions. Hence, only the tweets-IDs are made available to researchers. However, using the Twitter API the tweets can be “rehydrated” and the data of tweets that have not been deleted can be accessed (more details on the GitHub page). This dataset is still being continuously collected and routinely updated. If you have technical questions about the data collection, please contact the corresponding author.

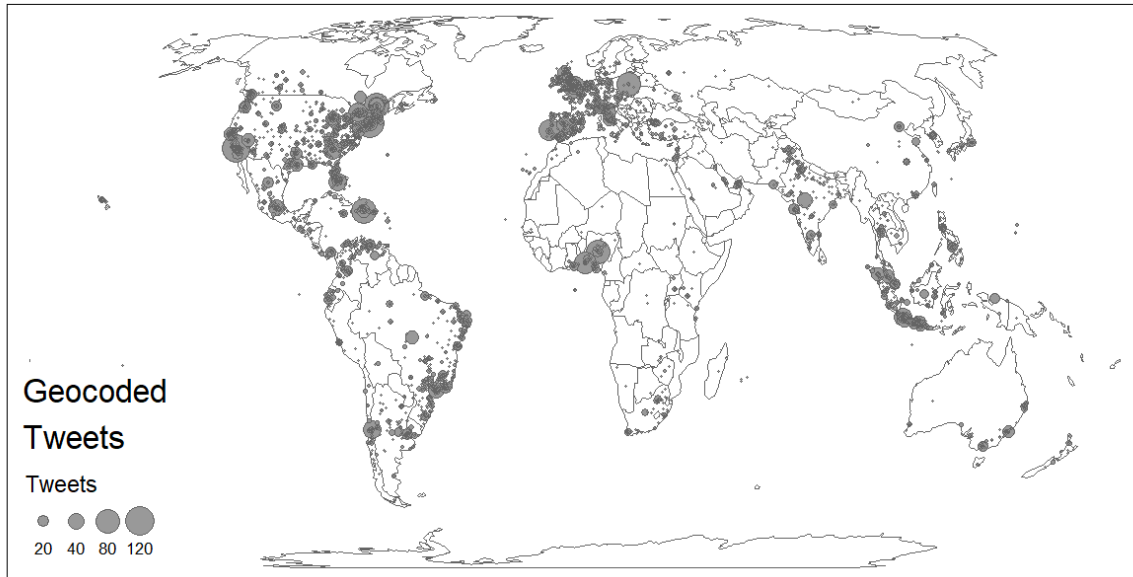


Figure 3: Geolocated tweets as of Apr 05, 2020.

4 Future Work

This research note's main objective was to introduce and share with the research community a dataset of tweets related to the COVID-19. We are continuously collecting and routinely updating the dataset. Similarly, we will be using Natural Language Processing and Text Mining, and Network Analysis to analyze the corpus of tweets to identify common popular responses to the pandemic and how these responses differ across time. Moreover, with this dataset we will explore to how information and misinformation about COVID-19 is transmitted via Twitter.