Bryant University

# Bryant Digital Repository

4-22-2022

# Predicting COVID-19 Fake News

Nilanjana Nambiar
*Bryant University*, nambiar@bryant.edu

# Bryant University

## HONORS THESIS

# Predicting COVID-19 Fake News

BY Nilanjana Nambiar

ADVISOR • Suhong Li
EDITORIAL REVIEWER • Tingting Zhao

_____

## Table of Contents

# ABSTRACT

The aim of this project is to use a machine learning model to identify COVID-19 fake news on Twitter and perform additional analysis on the fake news tweets to distinguish any common trends. As misinformation is very common in the information found online, the purpose of the study is to see how machine learning can be used to discern what information can be classified as true versus what is false. Prior research regarding fake news detection, modeling, and analysis was conducted to familiarize on the current studies provided in predicting and analyzing COVID-19 fake news on Twitter. In this study, over one billion tweets were collected and analyzed between March 2020 and October 2021 via the Twitter API using #covid and COVID-19 as keywords. Specific keywords were set to further filter the data. The training data set was retrieved from the University of Pennsylvania's COVID-19 healthcare misinformation dataset (CoAID), which included 137,799 tweets that were manually classified as 'Real' or 'Fake'. A binary logistic regression model was used to classify the COVID-19 tweets, resulting in approximately 92% of the tweets containing accurate information regarding COVID-19, while the remaining 8% containing false information about the coronavirus. Of the English-based tweets, Twitter posts with false information on COVID-19 were commonly related to politics, vaccinations, and social distancing, regardless of location.

## INTRODUCTION

In December 2019, the WHO was first notified of an increasing number of respiratory and pneumonia cases in Wuhan, China. Shortly after, the emergence of the coronavirus (SARS-CoV-2) caused millions of cases and deaths worldwide, the World Health Organization (WHO) declaring it a global health emergency. Due to the initial lack of knowledge, effective treatment, or vaccine, the COVID-19 outbreak led to a large amount of speculation and false information spread about the virus, especially on social media and online news sources.

The WHO classified this overabundance of misinformation as an 'infodemic,' which relates to the spread of false information about COVID-19 on social media. The world is driven by technology, and social media platforms currently have an undeniable amount of power in the contributions to the types of information and media that people consume daily. On the other hand, social media platforms such as Twitter, also provide the opportunity for people to voice their thoughts and opinions. While the internet can be a resource for providing accurate and factual information, rumors and false claims are also bound to spread equally as fast. Especially when it comes to something as crucial as the COVID-19 pandemic, misinformation can make it difficult for users to find reliable information online, potentially distorting peoples' perceptions. Social media companies today have made great strides to ensure that intentional and malicious false claims targeted at specific groups of people are removed, while still allowing people to exercise their right of free speech on the internet.

In this paper, the following literature review will cover related studies that have been performed in relation to analyzing COVID-19 fake news, as well as cover machine learning modeling techniques. The methodology section will discuss the process of obtaining the tweets, finding a usable training dataset, and running the model. The results section will provide the results of the model once it was applied to the current Twitter dataset and additional analysis that was performed to specifically analyze the tweets that were classified as fake. Finally, the conclusion section will summarize the outcome, relevance, limitations, and future implications for the study.

## LITERATURE REVIEW

Misinformation and the Infodemic Crisis

Research has shown that there was an initial fear of the severe consequences regarding the coronavirus pandemic with the surging number of cases and deaths continuously on the rise (Li et al., 2020). Just between the months of March and April 2020 alone, the number of COVID-19 cases grew from under 100,000 to over 3 million worldwide (Cuomo et al., 2020). This eventually led to an increased speculation online through user preferences and attitudes regarding COVID-19 (Cinelli et al., 2020).

However, as this topic tends to be subjective, there is no agreement amongst researchers over the definition of 'fake news' (Patwa et al., 2021). As defined by Shahi et al. (2020), misinformation refers to the circulation of false information. It can also be defined as "a claim of fact that is currently false due to lack of scientific evidence" and can be propagated "without constraints, does not entail any curation or peer-review, and does not require any professional verifications" (Kouzy et al., 2020). This is ideal for misinformation to be spread and become amplified on social media through personally tailored content. In another definition, "deceptive news which includes news fabrications, satire, hoaxes, etc., are considered as fake news" (Rubin et al., 2016). Existing research has also referenced the term 'infodemic,' which represents the increasing amount of misinformation spread uninhibited over traditional and social media at a rapid pace (Kouzy et al., 2020). Patwa et al. (2021) also mentions that "[d]espite the existence of several works dedicated for fake news, accurate automatic fake news detection is an extremely challenging task. The lack of a common acceptable benchmark dataset for this task is one of the key problems."

User, Content, and Keyword Analysis of Misinformation

An exploratory study conducted by Shahi et al. (2021) focused on the propagation, authors, and content of misinformation on Twitter around the topic of COVID-19 to gain insights. The tweets collected were fact-checked by over 92 professional fact-checking organizations between January and mid-July 2020, resulting in 1500 tweets relating to 1274 false and 226 partially false claims, respectively. The exploratory analysis that conducted on Twitter users

revealed that verified users are also involved in either creating (new tweets) or spreading (retweet) the misinformation. False claims propagate faster than partially false claims, and tweets with misinformation are often more concerned with discrediting other information on social media. The study also had various limitations. Potential selection bias was known in the data collection stage as Shahi et al. (2021) only considered rumors with specific tweet IDs that were eventually investigated by a fact-checking organization; thus, the data most likely excluded less viral rumors. Secondly, the interpretation of both hashtag and emoji usage by Twitter authors of misinformation were both culturally and contextually bound. The analysis was also influenced by facts such as age and gender; however, there was an overall lack of information on how the users intended them. Finally, the tweets utilized in this study were solely written in English; therefore, the range of topics discussed may have differed between English and non-English tweets. Overall, the work of Shahi et al. (2021) had considerable possibilities and possibly extend to topics such as embedding the category of fake news within news articles, as well as detecting the class of fake news on the web or social media, potentially with a time series machine learning model.

Kouzy et al. (2020) found that most tweets circulating around COVID-19 were issued by informal individuals and groups, while only a small percentage of the collected sample tweets belonged to verified Twitter accounts. In this study, a Twitter Archiver add-on was used to search Twitter for tweets containing one or more 11 common hashtags and three common key terms pertaining to the COVID-19 epidemic that were identified by the Symplur analytical tool. Samples were selected based on computer-generated random sequence, and a set of predetermined variables were collected for every individual tweet. User accounts were classified based on content into the following categories: informal individual/group, business/NGO/government, news outlet/journalist, and healthcare/public health/medical. Tweets were also categorized based on the content tone in the following categories: serious, humorous, and opinions. Tweets that contained genuine information regarding the COVID-19 epidemic were identified by cross-matching information presented by the World Health Organization (WHO), the Center for Disease Control and Prevention (CDC), peer-reviewed scientific journals, and prominent news outlets. Tweets that included information that could

be clearly refuted using one of the above-mentioned references were considered under misinformation, and tweets that could not be proven correct or incorrect by the references were designated as unverifiable information. The results concluded that majority of the tweets pertained to serious and genuine information on the COVID-19 epidemic. Of the collected tweets dataset, 24.8% of the tweets included misinformation, and 17.4% of the tweets included unverifiable information. The rate of misinformation spreading was higher among individual and group accounts, and the number of likes and retweets per tweet were not associated with a difference in either false or unverifiable content. It was also indicated that medical information and unverifiable content that pertained to the global pandemic propagated at an alarming rate on social media. Overall, this study demonstrated the idea that social media continues to have an immense impact on the spread of false information on social media platforms, especially on the current COVID-19 pandemic situation.

Research done by Li et al. (2020) showed that there was increased stigma against various societal groups due to the COVID-19 pandemic. Topics such as flu-like symptoms, personal protective equipment (PPE), Asian origin, and certain careers were used to mark individuals who may have spread the virus. Regarding the stigma against Asian populations in particular, numerous tweets implied that people who had different eating habits should be held responsible for the COVID-19 virus outbreak. This had created a cultural divide in the responsibility for the COVID-19 pandemic due to food preferences. Of the data collected, about one in five tweets highlighted the peril of the coronavirus outbreak in various aspects of peoples' lives, and the fear and anxiety associated with the threats of COVID-19 potentially encouraged individuals to blame others for the coronavirus situation, overall increasing stigma within the population. Compared to tweets without misinformation, online messages that included misinformation and conspiracy theories about COVID-19 were less likely to mention the actual threats and peril of the coronavirus situation and focused more on increasing stigma against certain groups of people. The authors emphasized that public health messages should attempt to minimize the unintentional stigmatization of infectious diseases such as COVID-19 by emphasizing the effectiveness of prevention measures, rather than associating discussions on the coronavirus to certain groups of people as disease-related

stigma and misinformation spread online directly correlated to adverse mental health outcomes.

Previous research has also shown that the communication and language used by world leaders and the media has an influence on the opinions and the behavior of the public, as well as have the power to address crises such as the current COVID-19 pandemic (Rufai & Bunce, 2020). The potential roles that Twitter have played in a crisis include infectious disease surveillance, predicting the dissemination of disease, and the spreading of public health information while assessing public views toward public health outbreaks. In a study done by Evanega et al. (2020), the most prominent topics of COVID-related misinformation that emerged in traditional media were analyzed between January 1, 2020, and May 26, 2020, with a total sample of over 38 million articles published in English-language media worldwide. The results showed that the media mentions of Donald Trump within the context of COVID-19 misinformation made up by far the largest share of the infodemic, and only 16.4% of the misinformation conversation was "fact-checking" in nature, thus suggesting that the majority COVID-related misinformation was conveyed by the media without any questions or corrections. Rufai and Bunce (2020) explored the role of Twitter as used by the Group of Seven (G7) world leaders in response to the coronavirus pandemic and used content analysis to categorize tweets into appropriate themes. The Group of Seven (G7) world leaders include: Justin Trudeau, Emmanuel Macron, Angela Merkel, Giuseppe Conte, Shinzo Abe, Boris Johnson, Donald Trump, Charles Michel, and Ursula von der Leyen of the EU council. The results of this study determined that majority of the tweets collected fell within the 'Informative' category, followed by 'Moral-boosting' and 'Political' categories, respectively. All the G7 leaders had a large following on Twitter; however, a disproportionally high number of followers were attributed to the Twitter account of former US President Donald Trump. Further, this study demonstrated that the number of followers did not necessarily correlate into more viral tweets as there was no requirement to follow a Twitter account to view, like, retweet, or comment on a tweet sent from a particular account (Rufai & Bunce, 2020).

The current literature commonly mentions the usage of searching specific keywords related to COVID-19 understand various features and trends in their analysis of misinformation online. In the work done by Cinelli et al. (2020), researchers collected more than eight million comments and posts over the time span of 45 days to analyze user engagement and interest about COVID-19 by performing a comparative analysis on five social media platforms, including Twitter, Instagram, YouTube, Reddit, and Gab. The topics related to the COVID-19 content were extracted and analyzed with Natural Language Processing (NLP) techniques. The data collected filtered contents according to a selected few of Google Trends' COVID-19 related queries, such as coronavirus, coronavirusoutbreak, imnotavirus, ncov, ncov-19, pandemic, wuhan, nCoV, IamNotAVirus, coronavirus_update, coronavirus_transmission, coronavirusnews, and coronavirusoutbreak. Specifically with Twitter, the researchers collected tweets related to the topic of the coronavirus by using both the search and stream endpoint of the Twitter API. The data derived from the search API represented a random sample of the tweets containing the selected keywords up to a maximum rate limit of 18000 tweets every ten minutes. With the data collected, the researchers were able to classify and divide their data into three labels: Conspiracy-Pseudoscience, Pro-Science or Questionable to check whether a data source was questionable or reliable. Natural Language Processing techniques were also used for text analysis for each social media platform, and the words were clustered by running the Partitioning Around Medoids (PAM) algorithm on their vector representations. Through this study, Cinelli et al. (2020) found that users in mainstream platforms were less susceptible to the spreading of information from questionable sources. There was also no significant difference in the way information derived from news outlets marked as either reliable or questionable spread; however, the findings did suggest that the interaction patterns between each social media platform and the online audience played a pivotal role in information and misinformation spreading.

Defining the Research Objective

The goal of this project is to machine learning modeling to identify fake news. From there, perform further analysis on specifically the fake tweets to look for any common trends.

# RESEARCH METHODOLOGY

Data Collection

Over two billion tweets were collected since March 2020 with the Twitter Application Programming Interface (API) using 'COVID-19' and '#covid' as keywords. For this project, the Twitter dataset used was a subset of tweets collected from March 2020 to October 2021, approximating to 1,314,705,279 collected tweets. The Tweets were originally uploaded to Amazon S3 storage system and were analyzed in the Databricks PySpark platform. PySpark is the Python API for Apache Spark, an open source, distributed computing framework that provides a set of libraries for real-time, large scale data processing.

Additional keywords were provided to filter the data, including 'masks', 'vaccines', 'social distancing', 'delta', 'variant', 'omicron', 'booster', and 'virus'. Both Python and SQL were utilized to query and analyze the data. Additional data cleaning was performed, such as ensuring that the tweets to be analyzed were all English tweets and no duplicate information would be included. A couple date related columns were also added to the dataset to perform some further analysis. After the filtering was applied, there were a total of 207,906,311 tweets in the dataset.

Text Pre-Processing Methods

Before the model is performed, the data pre-processing must be done to clean the text data and make it suitable for running the machine learning algorithm. By pre-processing the data, it increases the overall accuracy and efficiency of the model. Various data-preprocessing techniques are included to prepare the model:

- **Tokenizer:** Splits sentences into words
- **DocumentAssembler:** Prepares data into a Document format that is processable by Spark NLP
- **Lemmetizer:** Brings all the words in the data to its lemma, or base form
- **StopWordsCleaner:** Filters stop words out to obtain meaningful words to describe the topic
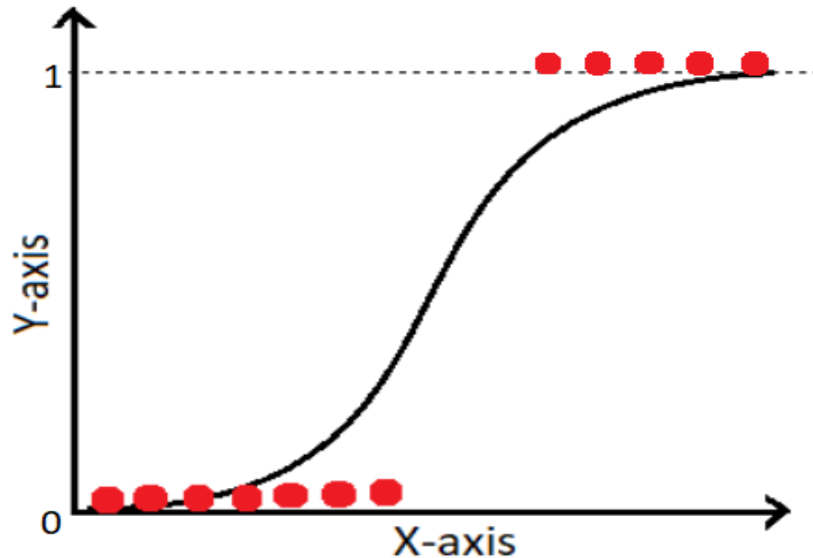- **Finisher:** Processes the data into a 'human readable' format from Spark NLP.

The model is automated through a machine learning pipeline. The ML pipeline helps in automating the machine learning workflow, as well as enabling the data to be transformed into a model that can be analyzed from the provided outputs.

Machine Learning Model Building: Logistic Regression

In this project, the modeling will be based on binary classification with two potential labels, real or fake. A common machine learning algorithm used for binary classification is Logistic Regression. Logistic Regression is a predictive analysis algorithm that is based on the concept of probability (Pant, 2019). The hypothesis of logistic regression tends to limit the cost function between the values 0 and 1.



*Figure 1 – Logistic Regression Diagram*

When using logistic regression, decision boundaries are used to separate the classes. This is when the classifier is expected to provide a set of outputs or classes based on probability when the inputs are passed through the prediction function, returning a probability score between 0 and 1.

Machine Learning Evaluation Metrics

One important evaluation metric for classification models is the confusion matrix. The confusion matrix is used to determine the performance of a classification model on a test

dataset, showing the relation between the predicted and actual values of the model. The matrix includes True Positive, True Negative, False Positive, and False Negative values. True Positive represents data points that have a true label which was correctly predicted by the model. True Negative represents the data points that have a false real label and have been correctly classified as false by the model. False Positive, also known as Type I Error, represents data points that have a false real label but has been predicted True by the model. False Negative, also known as Type 2 Error, represents data points that have a true real label and has been incorrectly predicted as False. True Positive and True Negative values indicate the datapoints in the model that have been correctly classified, while False Positive and False Negative values indicate data points that have been miss-classified by the model ("The Confusion Matrix: Unveiled", 2019).

Moving from the confusion matrix, some key metrics that are calculated are accuracy, precision, recall, and f-1 score:

- Accuracy = (TP + TN) / (TP + TN + FP + FN)
- Precision = TP / (TP + FP)
- Recall = TP / (TP + FN)
- F-1 = 2 * (Recall * Precision) / (Recall + Precision)

Accuracy is the simplest performance measure, and it is the ratio of correctly predicted observations to the total observations. Precision is the ratio of correctly predicted positive observations to the total predicted positive observations. A high precision score relates to a low False Positive rate. Recall is the ratio of correctly predicted positive observations to all the observations in actual class. Finally, F-1 Score is the weighted average of precision and recall, taking both false positives and false negatives into account ("Accuracy, Precision, Recall & F1 Score: Interpretation of Performance Measure", 2016).
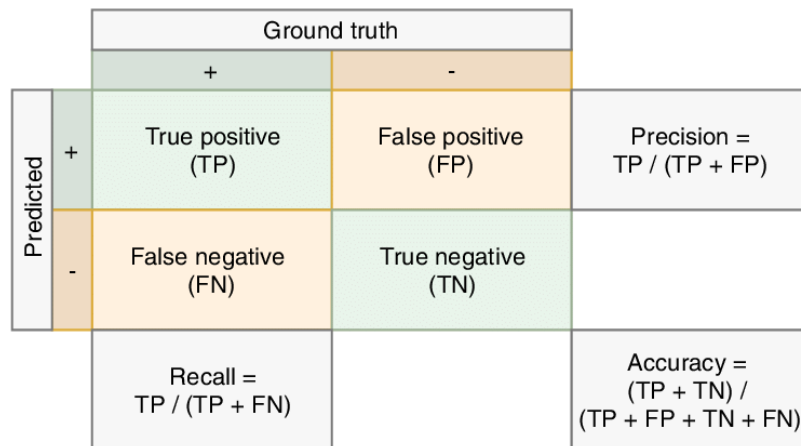
*Figure 2 – Confusion Matrix Diagram*

In a study done by Patwa et al. (2021), a binary classification task (Real vs Fake) was performed with a manually annotated dataset of 10,700 tweets and articles of real and fake news on COVID-19. Further analysis was done by running four machine learning models with the annotated dataset – Decision Tree, Logistic Regression, Gradient Boost, and Support Vector Machine. The SVM model performed the best, with an f1-score of 93.32%.

Obtaining Training Set and Data Pre-Processing for Machine Learning

To run the machine learning algorithm to predict fake news tweets, a training dataset needed to be obtained. In this project, the Pennsylvania State University's COVID-19 healthcare misinformation dataset (CoAID) was utilized. Within the dataset, it provided tweets that were manually classified as ClaimFake, ClaimReal, NewsFake, and NewsReal. Once these files were imported into Databricks, the dataset was then filtered so that each tweet would be classified as either Real or Fake. There was a total of 137,799 tweets in the training dataset that would be then used for modeling.
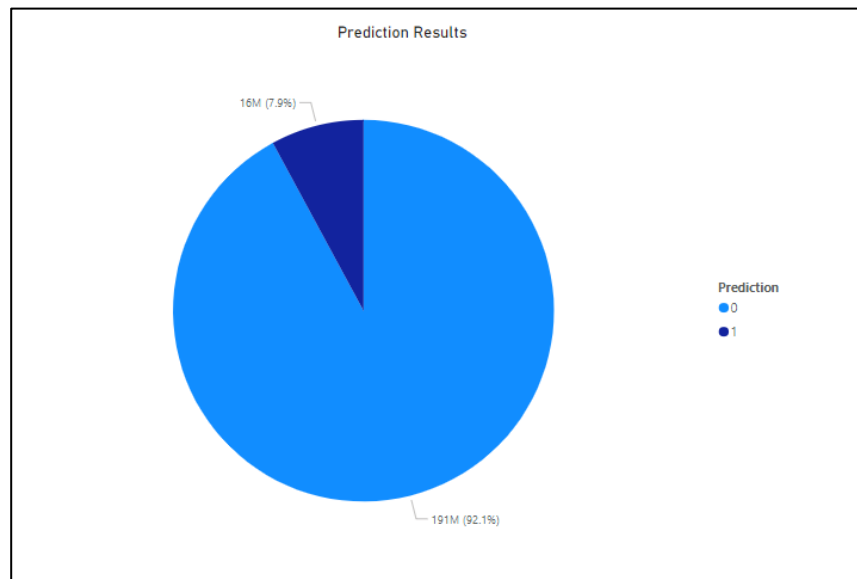
*Figure 3 – Sample Tweets from CoAID Training Dataset*

During the text-preprocessing stage, various steps were taken to process the data for modeling in the Spark Natural Language Processing (NLP) format. Once this was performed, a total of 133,606 tweets were available in the dataset.

## RESULTS

Model Testing

After pre-processing the data, the test dataset was split into 70% training and 30% testing. To perform the machine learning step, I utilized the binary logistic regression model and was able to achieve very good results. The accuracy score for the model was 0.99, meaning that 99% of the tweets were correctly classified from the model. Once the model was applied to the tweets, it was able to classify approximately 92% of the tweets that contained accurate information about COVID-19, while the remaining 8% of the tweets containing false information about the coronavirus. In this diagram below, the value 0 represents the classified 'Real' tweets, while the value 1 represents the classified 'Fake' tweets:

*Figure 4 – Model Prediction Results*

However, because the classification task is a highly imbalanced dataset with only 8% of the data as fake tweets, it is crucial to also look into the precision, recall, and f-1 scores to further analyze the validity of the model. The precision score was 0.91, so when the model classifies a tweet to be 'Real,' it is correct 91% of the time. The recall score was 0.92, so the model was able to correctly classify 92% of the 'Real' tweets. This shows that the model is ideal with both high precision and high recall. F-1 is the harmonic mean of precision and recall, and the model produced a score of 0.92. Because the score is high, this indicates a better prediction value.
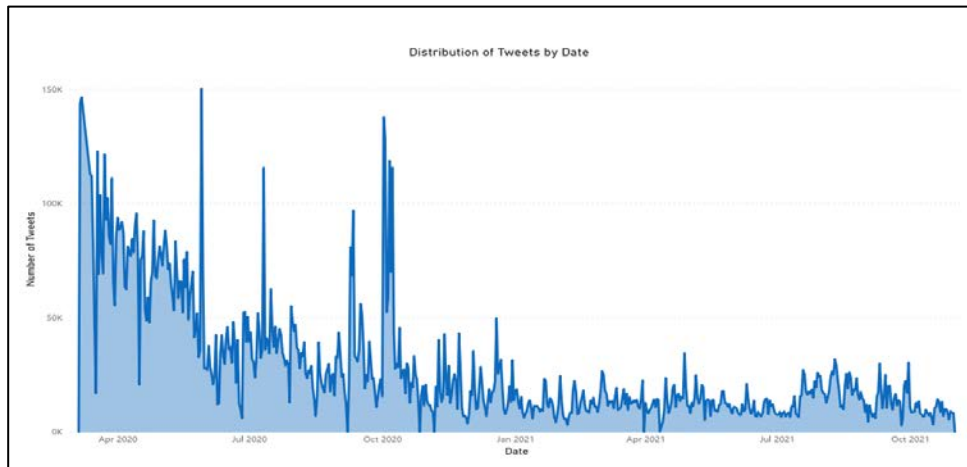
Analysis of Fake Tweets

First, a COVID-19 tweets time series was created to see the distribution of tweets from March 2020 to November 2021.

*Figure 5 – COVID-19 Tweets Time Series*

The two largest spikes in this chart are in the beginning of March, as well as in May 2020. This is potentially because of the increased Twitter discourse related to politics and decisions regarding COVID-19 and social distancing.

<u>User and Retweeter Analysis</u>

When looking into who is distributing fake news, it was interesting to note that of the top 10 users, only three user accounts were available on Twitter. Other accounts were either removed or suspended for spreading false information about COVID-19.



*Figure 6 – Twitter Users Distributing Fake News*

Of the Twitter users whose fake news tweets got retweeted the most, some notable users include Elon Musk and Joe Biden. Other users, like the top user 'ryanstruyk' were news reporters and correspondents, and bot accounts made an impact as well.
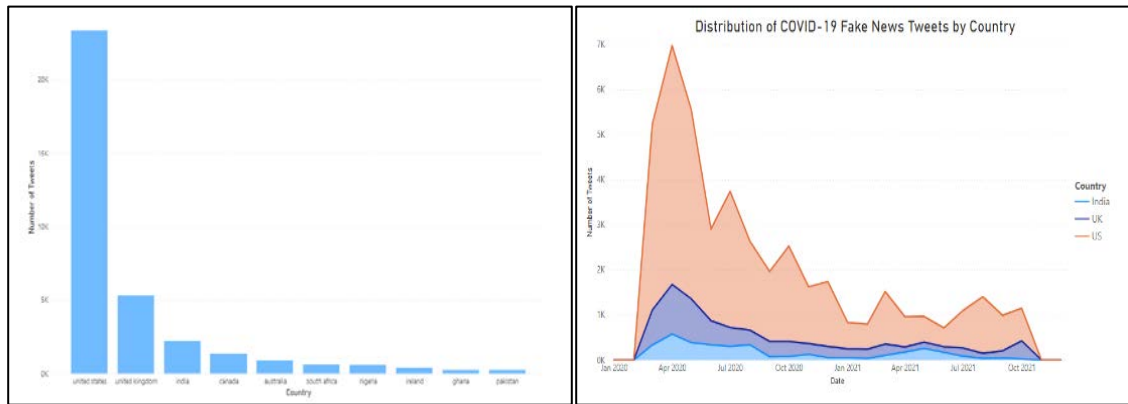
Top Words and Hashtags



*Figure 7 – Top Words and Hashtags in Fake News*

The two word clouds above showcase the top words and hashtags in fake news tweets. The sentiment in both these word clouds are very similar. Some top words to note are 'trump', 'biden', 'hoax', 'panic', and 'claims'. Of the top hashtags, some to take into account are #trump, #foxnews, #china, #politics, #trumpknew, and #chinesevirus. It was very common to see that among fake news tweets, the common themes included US politics, Donald Trump, China, social distancing, and vaccines.

Location Analysis

To analyze the data further, location was also examined. In this bar chart showcasing the top ten locations, United States, United Kingdom, and India were the top three.

*Figure 8 – COVID-19 Fake News by Country*

It can be assumed that because only English tweets were specified, there would be a more tweets from the United States. The United Kingdom and India are also large English-speaking counties, which is demonstrated in the charts above.

Now looking into the top five hashtags in fake news by country, one main difference was that the United States - in comparison to the United Kingdom and India - referenced Donald Trump and his responses to the COVID-19 crisis a lot more. Tweets from the United Kingdom and India were more focused on vaccinations and the pandemic in general.

| Country | Hashtag |
|---|---|
| United States | #trump, #pandemic, #quarantine, #cancelstudentdebt, #stayhome |
| United Kingdom | #lockdown, #vaccine, #nhs, #pandemic, #globalceasefire |
| India | #lockdown, #vaccine, #nhs, #pandemic, #globalceasefire |

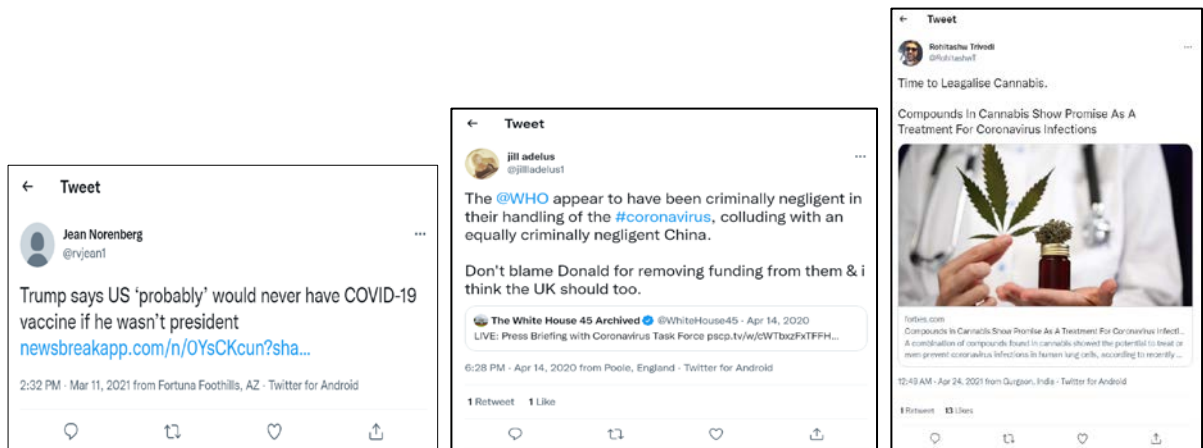*Figure 9 – Top 5 Hashtags in Fake News by Country*

*Figure 10 – Examples of Fake News in US (left), UK (middle), India (right)*

Above are some examples of fake news tweets by Country. While this classification task may be subjective, it is apparent that these tweets were correctly classified as fake news as they all go against the scientifically proven and factual information that is provided regarding the COVID-19 pandemic.

## CONCLUSION

The binary logistic regression model was able to accurately classify 92% of the tweets collected to have 'Real' information about COVID-19, while the remaining 8% of the tweets contained 'Fake' information about COVID-19. United States tweets tended to be higher in frequency in comparison to other countries. Tweets with false information regarding COVID-19 were commonly related to politics, the pandemic, social distancing, and vaccinations, which was clearly shown in the most keywords and hashtags. Donald Trump was the most referenced in fake news tweets, especially in the United States and United Kingdom. Occasionally, there were anti-Asian sentiment within tweets, most likely due to the fact that the virus first started spreading in Wuhan, China. Finally, it was very interesting to see that many fake news tweets and Twitter user accounts were currently unavailable, which is also highly due to Twitter's recent online misinformation policy.

Some limitations included that processing the data was a bit time consuming due to the large datasets, and there were occasional issues with the Databricks cluster. In addition, all the

tweets were English-based. Regarding future implications, it would be interesting to analyze more data based on a larger time frame. Also, trying other machine learning models and NLP modeling techniques to further analyze the data and look for any similarities and/or differences in prediction accuracy. In this study, only the United States, United Kingdom, and India were analyzed; therefore, comparing tweets from more locations as well as analyzing non-English tweets would be interesting to consider. Finally, looking into bot accounts and its impact on COVID-19 tweets could also be examined.

# **REFERENCES**

"Accuracy, Precision, Recall & F1 Score: Interpretation of Performance Measures." *Exsilio Solutions*, 9 Sept. 2016, https://blog.exsilio.com/all/accuracy-precision-recall-f1-score-interpretation-of-performance-measures/.

Cinelli, M., Quattrociocchi, W., Galeazzi, A., Valensise, C. M., Brugnoli, E., Schmidt, A. L., … Scala, A. (2020). The COVID-19 social media infodemic. *Scientific Reports*, *10*(1). https://doi.org/10.1038/s41598-020-73510-5

Cuomo, R. E., Purushothaman, V., Li, J., Cai, M., & Mackey, T. K. (2020). Sub-national longitudinal and geospatial analysis of COVID-19 Tweets. *PLOS ONE*, *15*(10). https://doi.org/10.1371/journal.pone.0241330

Evanega, S., Lynas, M., Adams, J., & Smolenyak, K. (2020). Coronavirus misinformation: quantifying sources and themes in the COVID-19 'infodemic' (Preprint). *The Cornell Alliance for Science*.

Kouzy, R., Abi Jaoude, J., Kraitem, A., El Alam, M. B., Karam, B., Adib, E., … Baddour, K. (2020). Coronavirus Goes Viral: Quantifying the COVID-19 Misinformation Epidemic on Twitter. *Cureus*, *12*(3). https://doi.org/10.7759/cureus.7255

Li, Y., Twersky, S., Ignace, K., Zhao, M., Purandare, R., Bennett-Jones, B., & Weaver, S. R. (2020). Constructing and Communicating COVID-19 Stigma on Twitter: A Content Analysis of Tweets during the Early Stage of the COVID-19 Outbreak. *International Journal of Environmental Research and Public Health*, *17*(18). https://doi.org/10.3390/ijerph17186847

Pant, A. (2019, January 22). *Introduction to logistic regression*. Medium. Retrieved March 23, 2022, from https://towardsdatascience.com/introduction-to-logistic-regression-66248243c148

Patwa, P., Sharma, S., Pykl, S., Guptha, V., Kumari, G., Akhtar, M. S., ... & Chakraborty, T. (2021, February). Fighting an infodemic: Covid-19 Fake News Dataset. In *International Workshop on Combating On line Hostile Posts in Regional Languages during Emergency Situation* (pp. 21-29). Springer, Cham.

Rubin, V. L., Conroy, N., Chen, Y., & Cornwell, S. (2016, June). Fake news or truth? using satirical cues to detect potentially misleading news. In *Proceedings of the second workshop on computational approaches to deception detection* (pp. 7-17).

Rufai, S. R., & Bunce, C. (2020). World leaders' usage of Twitter in response to the COVID-19 pandemic: a content analysis. *Journal of Public Health*, *42*(3), 510–516. https://doi.org/10.1093/pubmed/fdaa049

Shahi, G. K., Dirkson, A., & Majchrzak, T. A. (2021). An exploratory study of COVID-19 misinformation on Twitter. *Online Social Networks and Media*. https://doi.org/10.1016/j.osnem.2020.100104

"The Confusion Matrix: Unveiled." *Medium*, Towards Data Science, 9 Dec. 2019, https://towardsdatascience.com/the-confusion-matrix-unveiled-2d030136be40.