# Do The New York Times articles affect the Tweets about COVID-19 vaccines?

## An examination of attention and sentiments

Qinyue Hao

Columbia University
Quantitative Methods in the Social Sciences
Master's Thesis
May 10, 2022

# Table of Contents

# Abstract

Inspired by the agenda-setting theory, this study aims to find out if mainstream news media articles about the COVID-19 vaccines affect Twitter posts. Using data aggregated at the platform and date levels, linear regression models were used to test the association in two dimensions: the number of posts as a measure of attention and the average sentiment of content as a measure of attitude. Time series analysis was also performed to detect the time dependency in the number and sentiment of Tweets. The results showed no substantive association between the news articles and tweets but revealed a weekly pattern of tweet counts and sentiments as well as a dependency of the tweet sentiments on the previous day. This means that the New York Times articles do not influence the attention level and attitudes on Twitter about the COVID-19 vaccines. Therefore, it may be more effective to share useful information and updates directly on Twitter to promote the vaccines. In addition, when the discussion is emotionally directed in one direction, it usually lasts for some time. This feature can be leveraged to monitor people's attitudes, predict future direction, and take actions accordingly.

# Introduction

Vaccine hesitancy has become an especially important issue during the COVID-19 pandemic. Positive discussions in the media may help promote vaccinations. The media ecosystem today includes two parts – mainstream mass media and social media. The content

on mainstream media is determined by the media practitioners, while on social media, every

user can produce and share information and opinions. Given the increasing interaction

between the two types of media, I wonder whether mainstream media could affect how

people tweet about COVID-19 vaccines. If so, in what aspects? Does it affect how much

attention people pay to the vaccines, or further affect people's attitude towards them? To test

such relationships, this research takes The New York Times as an representative of

mainstream media, Twitter as an representative of social media, and compares the content of

the COVID-vaccine-related news articles in The New York Times with the relevant tweets.

## Literature Review

There has been much research on the effect of mass communication, and some scholars

found evidence that mass communications have the power to "cause learning and change the

attitudes and opinions of the audiences"[1]. A theory called Agenda setting has been developed

over decades, which explained how the mass media could act on people's minds and affect

how they think of an issue in multiple ways.

- Development of the agenda-setting theory

In the history of communication, the agenda-setting theory originated from the Chapel

Hill Study. It is a small-scale audience survey McCombs and Shaw conducted on the eve of

the U.S. election in 1968, which revealed the significant impact of mass media on the public

---

[1] Wilbur Schramm, "The Effects of Mass Communications: A Review," *Journalism Quarterly* 26, no. 4 (December 1, 1949): 397–409, https://doi.org/10.1177/107769904902600403.

agenda[2]. Later, they published a paper called "The Agenda-Setting Function of the Mass Media" in 1972, where the concept and the theoretical framework of this theory formally came into being. Then during the 1972 presidential election, McCombs and Shaw further tested the causality of the agenda setting effect in the Charlotte study. The results of the Charlotte study provide some evidence of causality, namely that the media (or at least newspapers) are influential in shaping the public agenda, rather than the other way around[3].

The next breakthrough in the agenda-setting theory happened in 1997 when McCombs and Shaw developed the agenda-setting theory to the second level. They proposed the theory of "attribute agenda-setting" in a paper on Spanish elections, where they included the influence of issue attributes and expanded the scope of agenda-setting theory[4].

The rapidly changing media ecosystem in the 21st century posed challenges to the agenda-setting theory. At the time of the Chapel Hill study, newspapers were at the center of the American public opinion arena, and television was in its take-off phase; half a century later, traditional media were suffering from decline, and the rapidly emerging Internet journalism industry was experiencing a power shift from portals to social media. Technological updates in the media ecosystem have posed new challenges to this classic theory and question whether the agenda-setting theory still has strong explanatory power. At that time, McCombs and Lei Guo proposed the third level of agenda-setting theory. Their study found that the agenda-setting effect still exists in the new era -- the attribute network

---

[2] MAXWELL E. McCOMBS and DONALD L. SHAW, "THE AGENDA-SETTING FUNCTION OF MASS MEDIA*," *Public Opinion Quarterly* 36, no. 2 (January 1, 1972): 176–87, https://doi.org/10.1086/267990.

[3] Donald Lewis Shaw and Maxwell E McCombs, *The Emergence of American Political Issues: The Agenda-Setting Function of the Press* (St. Paul: West Pub. Co., 1977).

[4] Maxwell McCombs et al., "Candidate Images in Spanish Elections: Second-Level Agenda-Setting Effects," *Journalism & Mass Communication Quarterly* 74, no. 4 (December 1, 1997): 703–17, https://doi.org/10.1177/107769909707400404.

portrayed by the media has a substantial impact on the public opinion – they called it

"Network Agenda Setting"(NAS).[5] Therefore, social network analysis has become an

important tool for public agenda and media effect research since then.

- Intermedia Agenda setting between mainstream media and Twitter

As the media ecosystem becomes more complicated with various media platforms

coming into play, a theory called "Intermedia agenda-setting" was developed to explain the

message flow between various media platforms and the impacts on the content they have on

each other. Some former research focused on the intermedia agenda-setting effects between

mainstream media and Twitter. It is found that while Twitter and the mainstream media agree

on the importance of many issues, Twitter also focuses on issues overlooked by the

mainstream media, such as news on environmental issues and gender equality.[6] Coway et al.

examined the intermedia agenda-setting effects with time series methods, on the Twitter feed

about the primary candidates in the 2012 presidential election, together with that about the

Republican and Democratic parties, in comparison with the articles on the most influential

newspapers in the United States. They discovered a symbiotic relationship emerged between

the Twitter agenda and traditional news, with different intensities and lags between releases

for different topics.[7]

[5] Lei Guo and Maxwell McCombs, "Network Agenda Setting: A Third Level of Media Effects," n.d., 20.
[6] Ingrid Rogstad, "Is Twitter Just Rehashing? Intermedia Agenda Setting between Twitter and Mainstream Media," *Journal of Information Technology & Politics* 13, no. 2 (April 2, 2016): 142–58, https://doi.org/10.1080/19331681.2016.1160263.
[7] Bethany A. Conway, Kate Kenski, and Di Wang, "The Rise of Twitter in the Political Campaign: Searching for Intermedia Agenda-Setting Effects in the Presidential Primary," *Journal of Computer-Mediated Communication* 20, no. 4 (July 1, 2015): 363–80, https://doi.org/10.1111/jcc4.12124.

Some studies have further looked into the direction of this agenda-setting effect. An intermedia study on elections found that "Slow" newspapers often precede other media's coverage[8]. Su and Borah explored the intermedia influence on the media content on the topics of climate with special attention to directionality. They found Twitter's and newspapers' agendas influence each other, and while Twitter drives the newspaper's agenda in terms of breaking news, newspapers influence ongoing discussions on Twitter during periods of non-breaking news.[9]

In conclusion, the intermedia agenda setting has been wildly proven to be true from multiple perspectives, and they also suggest mainstream media agenda often precede social media agenda, yet there are interactions with a time lag, and social media can sometimes in turn affect mainstream media agendas. However, the agenda-setting theory focused mainly on politics and public agenda, yet it is possible that such powerful effects may also exist in other areas, for instance, public health. It is claims more attention especially during the current pandemic – it would be of much help if mainstream media could help convey the right messages and knowledge about the COVID-19 pandemic and promote vaccinations. This research aims to evaluate the power of the influence of mainstream media and test whether the agenda agenda-setting can expand from politics to the field of public health, and be leveraged for good.

---

[8] Raymond A. Harder, Julie Sevenans, and Peter Van Aelst, "Intermedia Agenda Setting in the Social Media Age: How Traditional Players Dominate the News Agenda in Election Times," *The International Journal of Press/Politics* 22, no. 3 (July 1, 2017): 275–93, https://doi.org/10.1177/1940161217704969.

[9] Yan Su and Porismita Borah, "Who Is the Agenda Setter? Examining the Intermedia Agenda-Setting Effect between Twitter and Newspapers," *Journal of Information Technology & Politics* 16, no. 3 (July 3, 2019): 236–49, https://doi.org/10.1080/19331681.2019.1641451.

# Data Source

- News Source: The New York Times official website

According to the Reuters Institute, The New York Times is the most popular newspaper among US digital news consumers. Readership statistics show that 39% of Americans who subscribed to a digital edition of a print newspaper in 2020 paid for The New York Times. Besides, The New York Times Twitter account is one of the most followed news accounts on Twitter, with more than 52.5 million followers[10].

- Social Media content Source: Twitter

Twitter is one of the most popular platforms among US social media users. About 23% of US adults are users of Twitter, and 46% of Twitter users in the US use it every day[11]. The Selected Metrics and Financials report of Twitter shows that the monetizable Daily Active Users on Twitter is 38 million, by the end of 2021.

The majority of users of both Twitter and The New York Times are Democratic. According to Pew Research Center, more than 90% of people listing The New York Times as their source of news are Democrats[12]. Another study by Pew Research Center found that among the 10% of Twitter users who send the majority of the daily tweets, 69% are Democrats[13]. Therefore, I assume it's reasonable to observe and compare these two platforms.

---

[10] As of April 10th, 2022. The New York Times is the second most followed mainstream media on Twitter, only after CNN (@cnnbrk and @cnn), who does not provide an official API.
[11] "• Twitter by the Numbers (2022): Stats, Demographics & Fun Facts," February 22, 2022, https://www.omnicoreagency.com/twitter-statistics/.
[12] "25 New York Times Readership Statistics [The 2021 Edition]," *Letter.Ly* (blog), March 14, 2021, https://letter.ly/new-york-times-readership-statistics/.
[13] "Pew Finds 69% of Twitter's Most Prolific Users Are Democrats," *VentureBeat* (blog), October 15, 2020, https://venturebeat.com/2020/10/15/pew-finds-69-of-twitters-most-prolific-users-are-democrats/.

## Data Collection

I collected all The New York Times news articles' metadata using The New York Times official API[14] and the full text of articles using the Newspaper package[15], with the keywords "COVID vaccine". The dataset covers the period from 2020/01/01 and 2022/03/26 and contains the text and metadata of 10425 articles about COVID-19 vaccines. Because of the limits of the Twitter API and for the interest of time, I had to narrow down the time frame so that the number of requests does not exceed the API limit. To achieve that, I explored the time trend of the number of news articles, which indicates the attention The New York Times paid to COVID-19 vaccines.
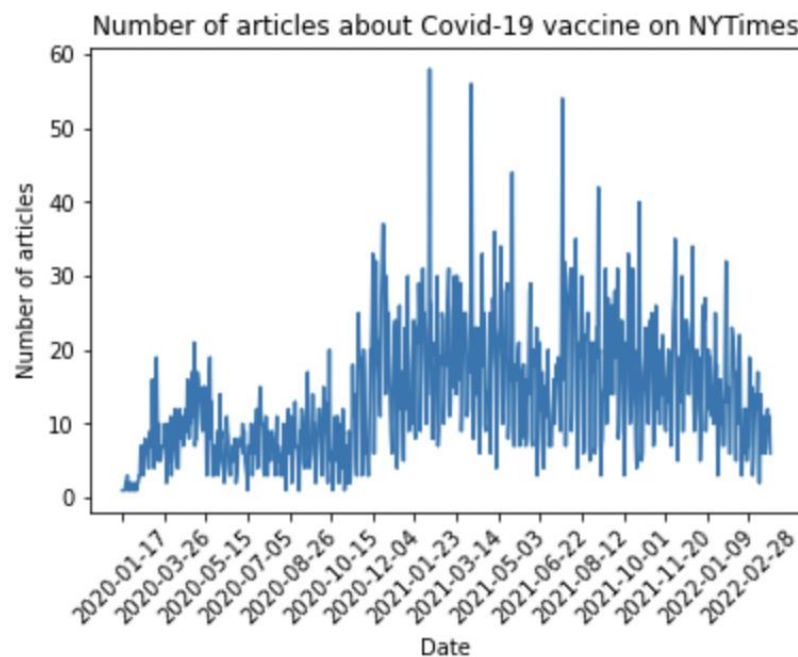


*Figure 1. The number of COVID-19-vaccine-related news articles in The New York Times*

---

[14] "Article Search | Dev Portal," accessed April 10, 2022, https://developer.nytimes.com/docs/articlesearch-product/1/overview.

[15] Lucas Ou-Yang, *Newspaper3k: Article Scraping & Curation*, Python, 2022, https://github.com/codelucas/newspaper.

A reasonable time frame for this study should have enough variation in the number of news articles. Figure 1 above shows that there are some sudden rises, which are good to be seen as cut points – something big must have happened on those dates. Therefore, before and after those dates, by comparing the news articles and tweets about the COVID-19 vaccine, I may be able to observe whether they changed in the same ways to make causal inferences.

To further determine a single cut point, I calculated the difference between each date and 1-7 dates before the date and selected those dates all of whose differences with former days are outliers (> Q3+1.5IQR), as shown in the table below.

*Table 1. The difference in the number of articles with 1-7 days before*

|  | abstract | diff1 | diff2 | diff3 | diff4 | diff5 | diff6 | diff7 | total_diff |
|---|---|---|---|---|---|---|---|---|---|
| **2020-03-12** | 16.0 | 16.0 | 8.0 | 9.0 | 16.0 | 7.0 | 16.0 | 12.0 | 84.0 |
| **2020-04-23** | 16.0 | 16.0 | 16.0 | 7.0 | 4.0 | 16.0 | 16.0 | 5.0 | 80.0 |
| **2020-04-28** | 17.0 | 17.0 | 7.0 | 6.0 | 17.0 | 17.0 | 9.0 | 17.0 | 90.0 |
| **2020-05-05** | 17.0 | 17.0 | 9.0 | 17.0 | 17.0 | 9.0 | 17.0 | 10.0 | 96.0 |
| **2020-05-12** | 15.0 | 15.0 | 6.0 | 15.0 | 5.0 | 15.0 | 3.0 | 15.0 | 74.0 |
| **2020-05-19** | 16.0 | 16.0 | 7.0 | 12.0 | 16.0 | 16.0 | 13.0 | 16.0 | 96.0 |
| **2020-06-01** | 14.0 | 14.0 | 14.0 | 9.0 | 14.0 | 11.0 | 14.0 | 5.0 | 81.0 |
| **2020-09-23** | 14.0 | 14.0 | 5.0 | 14.0 | 8.0 | 14.0 | 9.0 | 14.0 | 78.0 |
| **2020-10-13** | 20.0 | 12.0 | 18.0 | 16.0 | 10.0 | 7.0 | 11.0 | 5.0 | 79.0 |
| **2020-11-16** | 19.0 | 16.0 | 13.0 | 5.0 | 8.0 | 8.0 | 1.0 | 10.0 | 61.0 |
| **2021-02-10** | 58.0 | 38.0 | 44.0 | 44.0 | 48.0 | 36.0 | 33.0 | 36.0 | 279.0 |
| **2021-04-01** | 56.0 | 28.0 | 39.0 | 37.0 | 44.0 | 45.0 | 35.0 | 31.0 | 259.0 |
| **2021-04-29** | 36.0 | 12.0 | 9.0 | 17.0 | 19.0 | 30.0 | 11.0 | 16.0 | 114.0 |
| **2021-05-20** | 44.0 | 22.0 | 28.0 | 27.0 | 34.0 | 36.0 | 15.0 | 16.0 | 178.0 |
| **2021-07-19** | 24.0 | 15.0 | 13.0 | 9.0 | 5.0 | 8.0 | 6.0 | 7.0 | 63.0 |
| **2021-07-20** | 54.0 | 30.0 | 45.0 | 43.0 | 39.0 | 35.0 | 38.0 | 36.0 | 266.0 |
| **2021-08-30** | 23.0 | 16.0 | 17.0 | 2.0 | 4.0 | 6.0 | 2.0 | 7.0 | 54.0 |
| **2021-09-01** | 42.0 | 19.0 | 19.0 | 35.0 | 36.0 | 21.0 | 23.0 | 25.0 | 178.0 |
| **2021-09-09** | 30.0 | 15.0 | 16.0 | 22.0 | 23.0 | 27.0 | 10.0 | 11.0 | 124.0 |
| **2021-09-24** | 31.0 | 12.0 | 8.0 | 3.0 | 6.0 | 17.0 | 17.0 | 5.0 | 68.0 |
| **2021-10-07** | 33.0 | 12.0 | 12.0 | 16.0 | 25.0 | 30.0 | 12.0 | 13.0 | 120.0 |
| **2021-10-20** | 40.0 | 14.0 | 28.0 | 36.0 | 33.0 | 27.0 | 20.0 | 23.0 | 181.0 |
| **2021-12-10** | 30.0 | 17.0 | 11.0 | 11.0 | 14.0 | 25.0 | 20.0 | 9.0 | 107.0 |
| **2022-01-19** | 25.0 | 13.0 | 15.0 | 10.0 | 14.0 | 9.0 | 6.0 | 13.0 | 80.0 |
| **2022-02-01** | 32.0 | 12.0 | 22.0 | 25.0 | 19.0 | 22.0 | 18.0 | 16.0 | 134.0 |
| **2022-03-01** | 19.0 | 12.0 | 12.0 | 14.0 | 7.0 | 7.0 | 7.0 | 13.0 | 72.0 |

In the end, I chose to set the time frame to a month, February 2021 and collected 654,788 tweets and their metadata in this month. According to the table above, there were large variations in the number of COVID-19-vaccine-related news articles in The New York Times, especially on the date 2021-02-10, which allows us to observe how the tweets react to the changes.

## Theory and Hypotheses

This study aims to explore and test the relationship between the content of The New York Times articles and the content from the Tweets, with the example topic of COVID-19 vaccines. Based on the media agenda-setting theory, I assume that there is a relationship between the accountant of The New York Times articles and tweets about the same topic and that they tend to change in the same direction. There are two major hypotheses: (1) the level of attention The New York Times pay to COVID-19 vaccines affects the level of attention are users pay to the same topic;(2) the sentiment in The New York Times articles about COVID-19 vaccines affects the sentiment peoples have towards them on Twitter. To be more specific, people tweet more about COVID-19 vaccines when The New York Times writes more about them, and when The New York Times writes about the vaccines in a more positive way, the discussions about them tend to become more positive on Twitter.

## Variable descriptions

As described in the hypotheses, there are four major concepts I want to measure in this study: the attention The New York Times has on COVID-19 vaccines, the attention Twitter users have on COVID-19 vaccines, the emotion The New York Times has towards COVID-19 vaccines, and the emotion Twitter users have towards COVID-19 vaccines. I measured those four concepts in the following way: First, the level of attention on COVID-19 vaccines on each platform is measured with many documents that contain the keyword "COVID vaccine"; Second, the sentiment toward COVID-19 vaccines on each platform is measured based on the content and, more specifically, the words contained in each relevant documents,

using a specific natural language processing tool named VaderSentiment, which calculate the sentiment score by searching for and comparing the number of positive and negative words.

In the end, I have four continuous numeric variables for those two concepts ready for modeling. The measurements are on the date level, therefore for February 2021, there are 28 observations in total, one for each date in this month. Additional to the four main concepts and their proxies, I also take into account the time series pattern in them. To account for it the trend over time I've also included a categorical variable, the day of the week, to control for the possible cyclic pattern of those measures in a week.

## Methodology

- Natural Language Processing

To measure the sentiment of the relevant documents, I preprocessed the text in the documents with Natural Language Processing Techniques in the following steps:

1. Cleaned the text with Regular Expressions to keep only words of English letters only;

2. Removed the English "stop words", which are words in English that are so common that carry little information.

3. Lemmatized the words with the "WordNetLemmatizer" from the Natural Language Toolkit (NLTK)[16] open-source library;

---

[16] Steven Bird, Ewan Klein, and Edward Loper, *Natural Language Processing with Python: Analyzing Text with the Natural Language Toolkit*, 1st edition (Beijing ; Cambridge Mass.: O'Reilly Media, 2009).

4. Conducted sentiment analysis with the VADER sentiment analysis tools[17], and use the "compound score" as the proxy for the overall sentiment of documents. The "compound score" is the normalized sum of all the lexicon ratings that can vary between -1(most extreme negative) and +1 (most extreme positive). I chose to use VADER because it is a parsimonious rule-based model designed for sentiment analysis on social media text, and the majority of the data is social media text – tweets. In addition, it has been widely used in social media sentiment studies, especially for analyzing tweets[18].

- Linear Regression Modeling and Time Series Analysis

To model the relationship between two continuous numeric variables, I started with the simplest linear regression model – the ordinary list squares (OLS) model. Then I added some control variables, replaced the original variable with its quadratic or cubic forms to account for possible higher-order relationships, and added interaction terms to the model to see whether those newly added variables could improve the model fit. In addition, I also conducted a time series analysis with the ARIMA models to account for the long-term and short-term serial autocorrelation in the dependent variable that affects the model results. In

---

[17] C. Hutto and Eric Gilbert, "VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text," *Proceedings of the International AAAI Conference on Web and Social Media* 8, no. 1 (May 16, 2014): 216–25.

[18] Anton Borg and Martin Boldt, "Using VADER Sentiment and SVM for Predicting Customer Response Sentiment," *Expert Systems with Applications* 162 (December 30, 2020): 113746, https://doi.org/10.1016/j.eswa.2020.113746; Dr Rajesh Bose, P. S. Aithal, and Sandip Roy, "Survey of Twitter Viewpoint on Application of Drugs by VADER Sentiment Analysis among Distinct Countries," SSRN Scholarly Paper (Rochester, NY: Social Science Research Network, March 15, 2021), https://papers.ssrn.com/abstract=3805424; J. Garay, R. Yap, and M. J. Sabellano, "An Analysis on the Insights of the Anti-Vaccine Movement from Social Media Posts Using k-Means Clustering Algorithm and VADER Sentiment Analyzer," *IOP Conference Series: Materials Science and Engineering* 482 (March 2019): 012043, https://doi.org/10.1088/1757-899X/482/1/012043.

the end, I performed the Granger causality test to detect the causal relationship between the content of The New York Times articles and tweets. Granger Causality[19] test detects causality by testing whether a time series is useful for forecasting another. It is essentially an F-test for the joint significance of the lags of the independent variable, net of the lags of the dependent variable, which is a more general form of the cross-lagged correlation analysis technique in "the Charlotte study" McCombs and Shaw used to detect causality when they developed the agenda-setting theory.

## Results

- Descriptive statistics – the number of News articles and Tweets

I initially explored the data with the descriptive statistics and boxplots. As shown in Table 2 below, in February 2021, on average, there are about 19 New York Times news articles and 23,382 tweets about COVID-19 vaccines every day. The minimum number of daily New York Times news articles about COVID-19 vaccines is 7 and the maximum number is 58. The minimum number of daily Tweets about COVID-19 vaccines is about 6264 and the maximum number is 31324. Overall, there are many more tweets than news articles about cove it 19 vaccines.

*Table 2: Descriptive statistics of the number of News articles and Tweets each day*

|  | count_n | count_t |
|---|---|---|
| count | 28 | 28 |

---

[19] C. W. J. Granger, "Investigating Causal Relations by Econometric Models and Cross-Spectral Methods," *Econometrica* 37, no. 3 (1969): 424–38, https://doi.org/10.2307/1912791.

| | | |
|---|---|---|
| **mean** | 19.39286 | 23381.8929 |
| **std** | 10.32917 | 6264.46475 |
| **min** | 7 | 5250 |
| **25%** | 12.25 | 20331.75 |
| **50%** | 19 | 24845 |
| **75%** | 24.25 | 27734 |
| **max** | 58 | 31324 |

It is also worthwhile looking into the distribution of the sentiment in tweets within each date. The box plot below shows the distribution of sentiment scores each day, which suggests that the news sentiments around February 13th, 2021 were very special – the sentiment score is more positive, and more clustered around the median value, compared to other dates. It may suggest that on that date The New York Times was reporting an important event about COVID-19 vaccines similarly and positively. Besides, a decrease in the sentiment score is observed on February 15th, with the median sentiment score back to 0 (neutral), which may suggest The New York Times has finished reporting the specific positive event about COVID-19 vaccines and resumed the neutral reporting style.
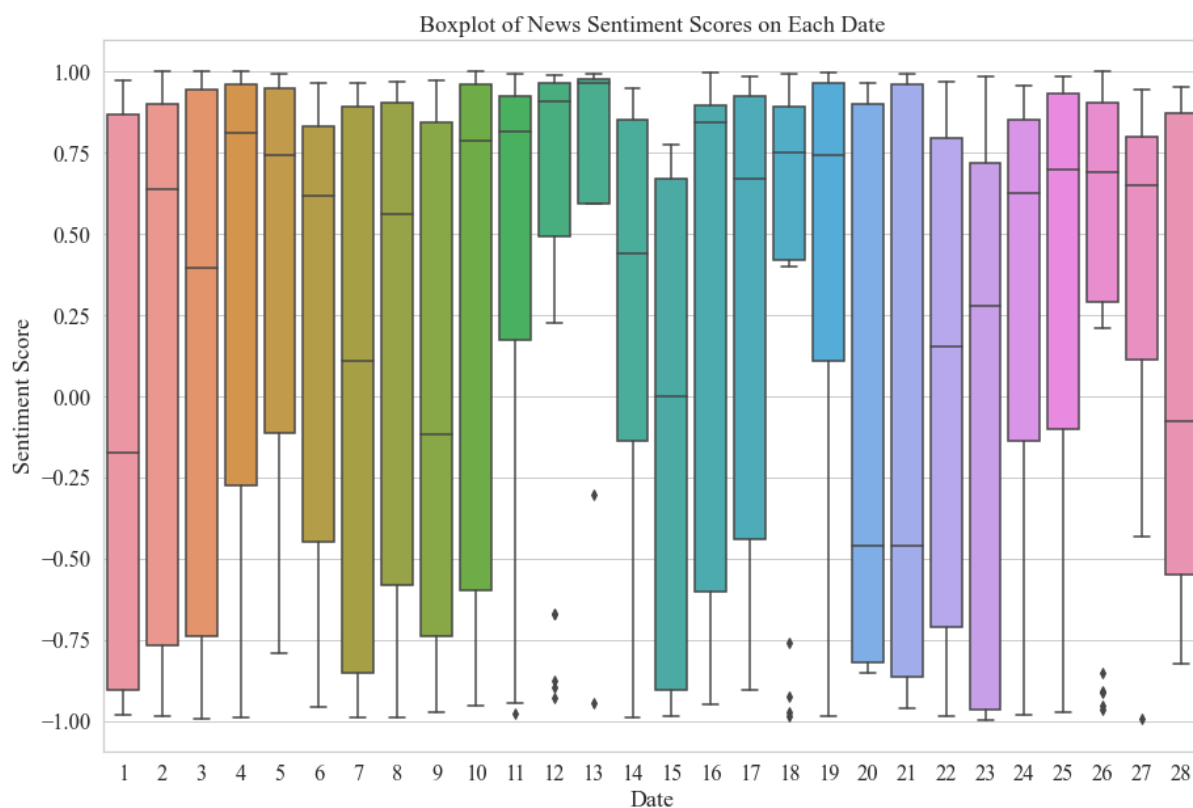
*Figure 2: News Sentiment Score Distribution on each date*

- Descriptive statistics – the sentiment of News articles and Tweets

As shown in Table 3 below, in February 2021, on average, the average sentiment score in

The New York Times news articles about COVID-19 vaccines is 0.22 while the average

sentiment score in tweets is much lower – 0.07. Both the average sentiment score is positive,

which means both New York Times journalists and editors and Twitter users had positive

attitudes toward COVID-19 vaccines. The New York Times articles write more positive

about the vaccines while for Twitter users the attitude is closer to neutral. Although on

average, positive sentiments were observed on both platforms, the text sentiment varies

across time. For The New York Times articles, the average sentiment score for each day

varies from -0.17 to 0.56, while for tweets there was much less variation -- the minimum is

0.04, which is still positive, and the maximum value is 0.13, mildly positive.

*Table 3: Descriptive statistics of the average sentiment in News articles and Tweets each day*

| | sentiment_n | sentiment_t |
|---|---|---|
| **count** | 28 | 28 |
| **mean** | 0.222853 | 0.07678 |
| **std** | 0.196539 | 0.023407 |
| **min** | -0.176511 | 0.043603 |
| **25%** | 0.051759 | 0.063142 |
| **50%** | 0.232983 | 0.072201 |
| **75%** | 0.377745 | 0.093325 |
| **max** | 0.565725 | 0.13359 |

As for the tweets about COVID-19 vaccines, the sentiment scores did not vary as much

as the news. However, we can still observe the same trend in the sentiments on The New

York Times and Twitter -- an increase in sentiment around February 13[th] and a decrease on

February 15[th]. It may be indicating that there exists a relationship between the sentiment in
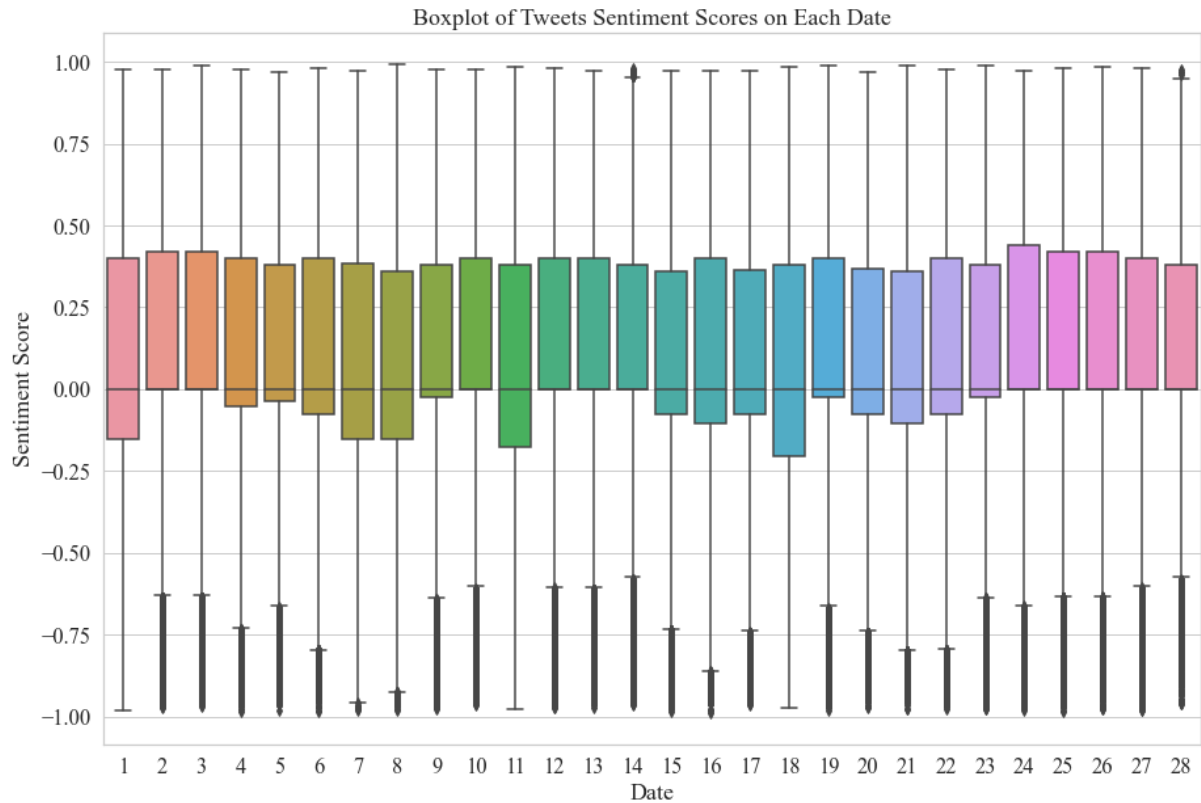
the news and tweets.

*Figure 3: Tweets Sentiment Score Distribution on each date*

- Exploratory data analysis

As shown in figure 4, the number of news articles and tweets have approximately the same time trend, and weekly patterns are observed in both variables. In light of this cyclic pattern, there is a need to control the day of the week when modeling the relationship between the news and tweets, to rule out the impact of this time trend.
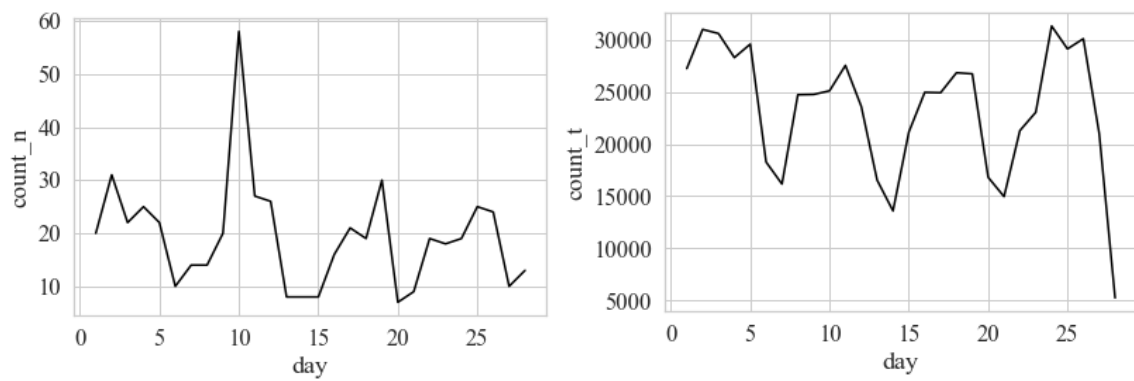
*Figure 4: The Number of The New York Times articles(left) and Tweets(right) overtime*

Before modeling, I visualized the possible bivariate relationship between the number of news and tweets to see the strength, the direction, and the order of the relationship. Figure 5 shows that the relationship between the number of tweets and news articles is overall positive and there is a chance that the relationship is quadratic/cubic rather than linear.
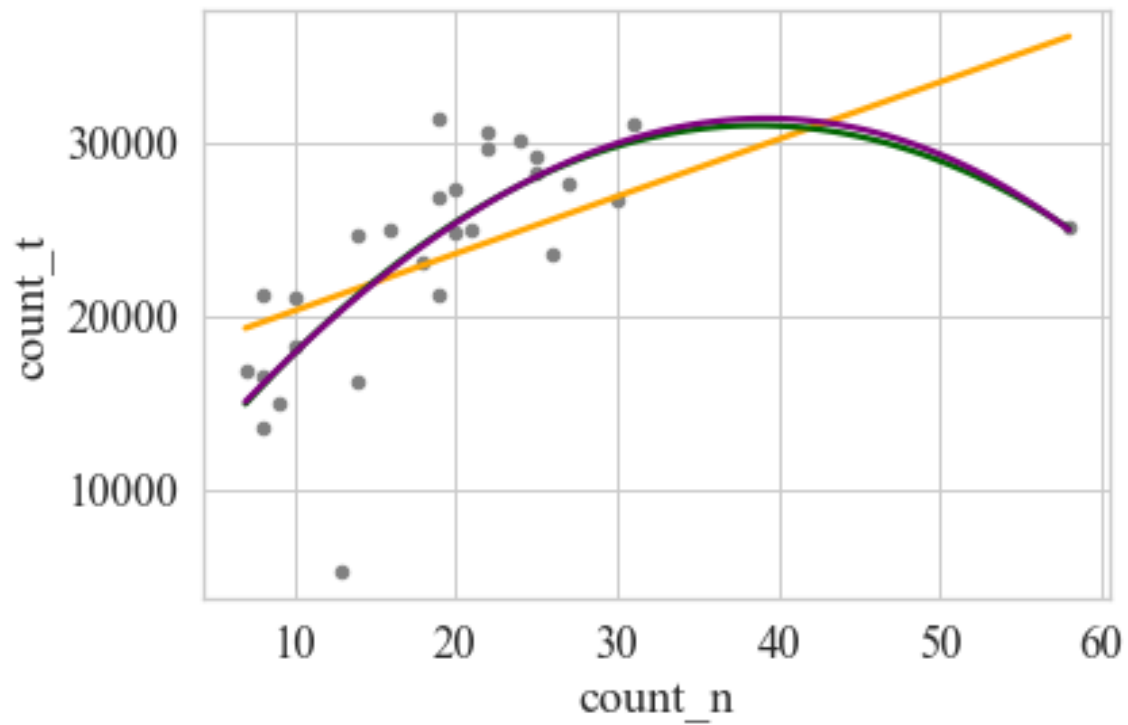
*Figure 5: The number of Tweets by the number of News articles*

As shown in figure 6, different from the counts, the sentiment scores of news articles and tweets do not change at a very similar pace and need further exploration into details. Same as the number of articles, the sentiment in New York Times articles about COVID-19 vaccines also has a weekly trend, which should be controlled for in the modeling steps.
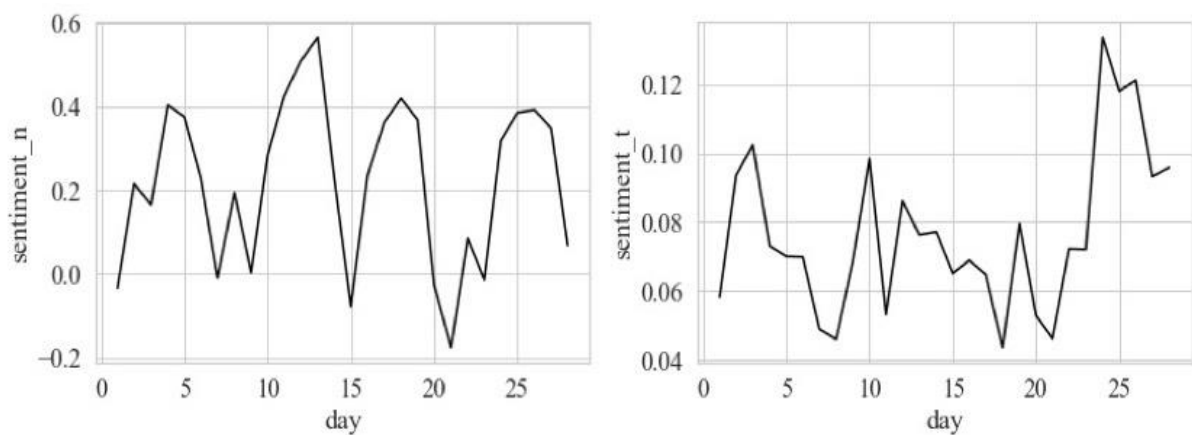


*Figure 6: The Average Sentiment Score of The New York Times articles (left) and tweets (right) overtime*
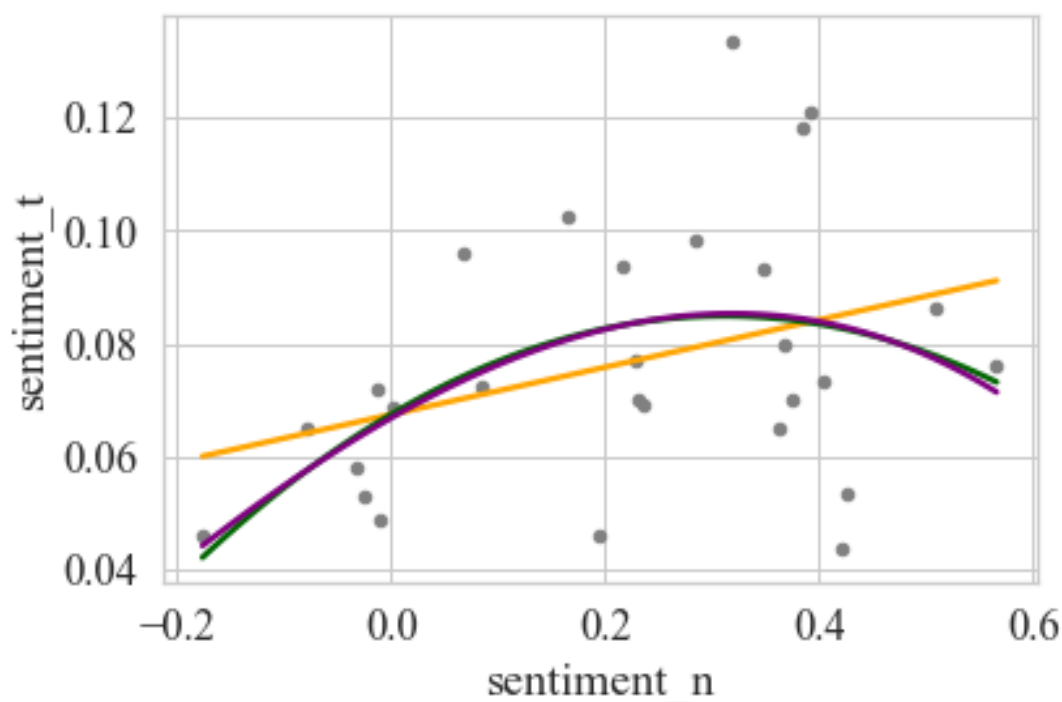
*Figure 7: The sentiment of Tweets by the number of News articles*

As shown above, although the data points are more scattered, the relationship between

the sentiments in the tweets and news articles is also positive. It is also possible that a

quadratic/cubic line rather than a straight line is a better fit for the relationship, while the

quadratic/cubic fitting lines are not much different.

- Comparison between the content of News and Tweets

To get a general understanding of the content of The New York Times and Twitter posts,

I made a comparison between the frequently used keywords in the news and tweets with

word clouds as shown in the two figures below. After text cleaning and lemmatization, I

tokenized the corpora and got the frequency of each word, and sized the words in the word

cloud pictures with their frequencies.

*Figure 8: Frequently used words in New York Times articles in Feb 2021*



*Figure 9: Frequently used words in Tweets in Feb 2021*

Given the two figures above, The New York Times seems to have a focus on the people,

location, and time in reporting COVID-19 vaccine-related events. It is worth noting that the

word "like" appears very frequently in The New York Times posts which may indicate a

positive sentiment towards the COVID-19 vaccines. This is also in line with the overall

positive sentiment in The New York Times articles we discovered above. As for the

frequently used keywords in tweets, it seems that people tweet more about vaccination experience. They mention a lot about receiving the COVID-19 vaccine shots, especially the first dose, and they also mention the major vaccine manufacturers such as Pfizer and Johnson& Johnson. Comparing the word cloud pictures, we can see that news articles and tweets may have different focuses when they mention COVID-19 vaccines – News articles write more on a more general level and report new big events about COVID-19 vaccines, while on Twitter, people discuss vaccine-related experiences on the individual level.

**Results**

  I aggregate the data into a date-level time-series structure to match and merge the news data set and the tweets data set, then modeled the relationship between the number of news and tweets, as well as the relationship between the average sentiments in the news and tweets respectively, in response to each of the hypotheses.

- Models on the number of news articles and tweets

*Table 4: Linear Regressions Predicting the Number of Tweets*

| | Model 1 | Model 2 | Model 3 |
|---|---|---|---|
| Intercept | 17000*** (2186.38) | 28170*** (2768.50) | 27790*** (1642.04) |
| count_n | 329.13** (99.90) | -25.72 (87.77) | |
| day_of_the_week[T.Monday] | | -4167.14 (2474.10) | -4108.08[+] (2285.82) |
| day_of_the_week[T.Saturday] | | -9779.33** (2733.68) | -9614.38*** (2295.24) |

| | | | |
|---|---|---|---|
| day_of_the_week[T.Sunday] | | -15390***(2632.76) | -15270*** (2292.98) |
| day_of_the_week[T.Thursday] | | 402.42 (2308.52) | 395.79 (2272.74) |
| day_of_the_week[T.Tuesday] | | -1677.06 (2334.76) | -1653.42 (2274.49) |
| day_of_the_week[T.Wednesday] | | 595.74 (2338.36) | 1093.89 (2393.00) |
| count_n_3 | | | -0.0161 (0.02) |
| Observations | 28 | 28 | 28 |
| Adj. R-squared | 0.267 | 0.729 | 0.737 |
| Durbin-Watson | 1.319 | 1.321 | 1.342 |

$^+p < .1$. $^*p < .05$. $^{**}p < .01$. $^{***}p < .001$.

In model 1, the simple OLS model showed a positive relationship between the number of news and tweets. The number of tweets increase by about 329, on average, for every one more The New York Times article published on that date. However, this relationship is not reliable. The results of model 2 and model 3 above provide evidence against my Hypothesis 1 that the number of news affects the number of tweets on the same date. Controlling for the day of the week, there is no substantial relationship between the number of tweets and the number of news, neither its original form nor its cubic form. Instead, the results revealed a weekly trend. Controlling for the cubic form of the number of news, on average, the number of tweets mentioning the COVID-19 vaccines on Monday is about 4108 less than on Friday (the reference day of the week), that on Saturday is 9617 less than on Friday, and that on Sunday is 15390 less than on Friday. It indicates that people tweet substantially less about the COVID-19 vaccines on weekends.

- Models on the sentiment of news articles and tweets

*Table 5 : Linear Regressions Predicting the Sentiment of Tweets*

|  | **Model 1** | **Model 2** | **Model 3** |
|---|---|---|---|
| Intercept | 0.067*** (0.006) | 0.075*** (0.018) | 0.0265 |
| sentiment_n | 0.042*** (0.022) | -0.017 (0.02) | 0.031 (0.028) |
| day_of_the_week[T.Monday] |  | -0.012 (0.016) | -0.007 (0.017) |
| day_of_the_week[T.Saturday] |  | -0.009 (0.021) | -0.024$^{+}$ (0.013) |
| day_of_the_week[T.Sunday] |  | -0.017 (0.016) | -0.011 (0.016) |
| day_of_the_week[T.Thursday] |  | -0.003 (0.019) | -0.037* (0.013) |
| day_of_the_week[T.Tuesday] |  | 0.015 (0.016) | 0.004 (0.015) |
| day_of_the_week[T.Wednesday] |  | 0.033 (0.035) | 0.019 (0.013) |
| lag_sentiment_t |  |  | 0.6965*** (0.18) |
| Observations | 28 | 28 | 27 |
| Adj. R-squared | 0.090 | 0.086 | 0.452 |
| Durbin-Watson | 1.058 | 0.610 | 1.802 |

$^{+}$p < .1. *p < .05. **p < .01. ***p < .001.

Model 1 shows that the sentiment is positively related to the sentiment in news on the same date. As the average sentiment score of news increases by one, the average sentiment score of tweets increases by about 0.042. However, when the weekly trend is controlled for,

as in Model 2, that relationship becomes not statistically significant. The explanatory power

of Model 1 and Model 2 are both very low, which indicates that the number of The New

York Times articles about COVID-19 vaccines is not a good predictor of the number of

Tweets about COVID-19 vaccines. Besides, the Durbin-Watson statistics of the model

residuals are also very low, which indicates that there is autocorrelation in the residuals from

Model 1 and Model 2. To detect and control for the autocorrelation in the model residuals, I

conducted Time Series Analysis on the residuals of Model 2 with the ARIMA models.

"ARIMA" being the acronym for "Auto Regressive Integrated Moving Average",

ARIMA models are a collection of models used for providing explanations of the current

time series data based on historical data, including the lagged values of the variable and the

lagged forecast errors. With the same logic, the model can also be used to make future

predictions. Auto Regression (AR) is a model that uses the lagged observations of the

variable. Integrated (I) is a differencing process on the observations to make the time series

stationary. Moving Average (MA) is a model that includes the residuals of the previous

forecasts to make current forecasts.

The result of the auto-ARIMA model showed that the ARIMA (1,0,0) model fits the

residuals of my Model 2 the best. The "Auto Regressive" (AR) term should be 1, which

means 1 lag of the dependent variable should be added into the model. The Integrated (I)

term should be 0 because the residuals of Model 2 are already stationary, and therefore do not

need any differencing. The Moving Average (MA) is also 0, indicating no lagged forecast

errors are needed in the model, in addition to the lagged dependent variable, to rule out the

serial correlation. In light of these results, I decided to include 1 lag of the sentiment in

Tweets, which would control for the time series dependency and leave the rest of the Model 2

residuals just white noise – as shown in Figure 10, given the autocorrelation line never

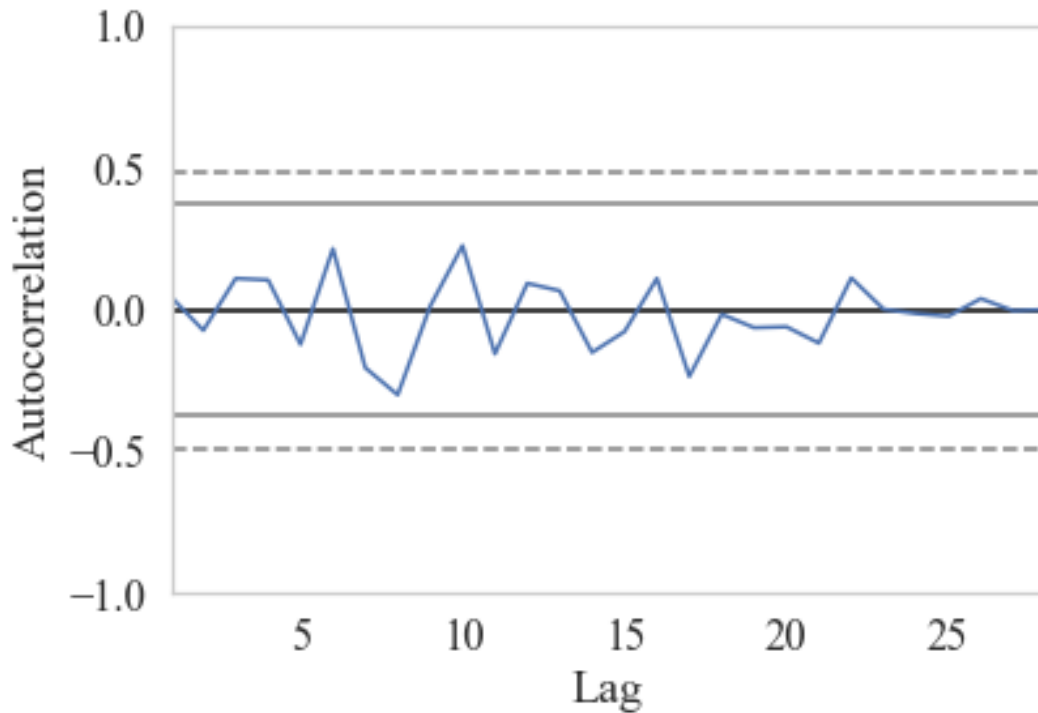crossed the margin lines with any number of lags.



*Figure 10: The residuals left after the ARIMA(1,0,0) model*

With the lagged tweets sentiment score included in Model 3, the model fit improve a lot,

and a highly significant relationship was observed between the tweet sentiment on the date

and the previous date. All else being equal, for every unit increase in the average sentiment

score of the previous day's tweets, the average sentiment score of the tweets increases by an

average of approximately 0.7. Besides, the tweets are on average more negative on Thursdays

and Saturdays, compared to those on Fridays.

Further testing on the causality between news and tweets is conducted with the Granger

Causality test to determine whether one time series helps forecast another. However, none of

28

the Granger Causality tests results turned out to be statistically significant, which means we cannot reject the null of no causality from the number of posts on The New York Times and Twitter, or from the sentiment of posts on The New York Times and Twitter.

## Discussion

In my initial hypotheses, I suspected that news articles may affect tweets on the COVID-19 vaccine topics, as the agenda setting theory suggested in the political science field. However, this "agenda setting effect" did not show in this study between The New York Times articles and the Twitter posts, neither in terms of the count of posts nor the sentiment orientations of the posts. Instead, the day of the week accounts seems to account for some of the variation in the discussion frequency and the sentiment towards COVID-19 vaccine-related topics. People have a tweeting habit of being more active on Fridays, than Wednesdays, Thursdays, and Weekends. They also sent more positive tweets on Fridays, compared with on Thursdays and weekends. This may be due to the fact that Friday is a commonly agreed happy day, followed by a two-days off for most people, and they started relaxing from the heavy work. Another factor that plays a significant role in explaining the sentiment toward COVID-vaccine-related topics is the Twitter public opinion on them on the date before – the overall attitude towards COVID-19 vaccines persist and strongly affect people's attitude on the next day.

It is worth noting that the sentiment of the tweets towards COVID-19 vaccines is not dependent on the news sentiment, but dependent on that of the previous day. It implies that if there is a need to advocate COVID-19 vaccines and relieve the resistance and antipathy to it

on Twitter, publishing news articles to demonstrate their nice property may not effectively

reach Twitter users and sway their attitudes. Instead, it is more advisable to leverage the

discussions inside the Twitter platform, such as inviting influential doctors or medical

scholars to demonstrate the advantages of COVID-19 vaccines. According to the time series

pattern, once a positive discussion starts, it can last for days if combined with constant

nudges in the same direction.

The absence of evidence for my hypothesized relationship between news and tweets also

implies that the "intermedia agenda setting" theory may be hard to generalize into other fields

such as health-related issues like COVID-19 vaccines, though it's been proven in the political

field empirically for multiple times. However, it does not completely deny the possibility that

relationships exist between mainstream news media and social media. This study has some

limitations. First, The New York Times and Twitter are selected in this study, because of

their popularity and their well-developed public APIs. Other media platforms more or less

restrict the retrieval of their data for analysis, which makes their data inaccessible to me. The

two platforms in this study may not be the best representatives of mainstream and social

media platforms for comparison, because their user groups do not predominantly match. Set

aside the realist concerns, it is ideal to include the content from all the mainstream media

platforms and all the social media platforms. Future studies can continue working in this

direction and aim to pick or construct more comparable samples. Second, there may be

multiple ways to operationalize the concepts of interest. For instance, the concept of

"attention" on certain issues can be measured in various ways. I chose to use the number of

publications on the media platforms (news articles or tweets), but the frequency of weighted

frequency of the phrase "COVID-19 vaccine" and its synonyms can also measure "attention". It may be worth some experiments to determine which one best measure the concept.

In summary, this study detected no prominent intermedia effect, and studies may be conducted in the future to include the other media platforms and use other possible measures to see whether the findings of this study persist.

## Conclusion

There is no substantial intermedia effect between the COVID-vaccine content on The New York Times and Twitter, neither in terms of the number of posts nor the sentiment orientations. Instead, time series dependency was found on the sentiment of tweets. For one thing, the average number and sentiment of tweets change substantially in a week. For another, the mean tweet sentiment on a day has an impact on that of the next day. Therefore, instead of merely publishing updates on the mainstream news media, the promotion of COVID-19 vaccines on Twitter may require some engagement in the discussions on the social media platform itself. Once the discussions are emotionally driven in a direction, the impact tends to last for some time, which can be leveraged to monitor people's attitudes and act accordingly.

## References

"• Twitter by the Numbers (2022): Stats, Demographics & Fun Facts," February 22, 2022.

        https://www.omnicoreagency.com/twitter-statistics/.

Letter.ly. "25 New York Times Readership Statistics [The 2021 Edition]," March 14, 2021.

https://letter.ly/new-york-times-readership-statistics/.

"Article Search | Dev Portal." Accessed April 10, 2022.

https://developer.nytimes.com/docs/articlesearch-product/1/overview.

Bird, Steven, Ewan Klein, and Edward Loper. *Natural Language Processing with Python:*

*Analyzing Text with the Natural Language Toolkit*. 1st edition. Beijing ; Cambridge

Mass.: O'Reilly Media, 2009.

Borg, Anton, and Martin Boldt. "Using VADER Sentiment and SVM for Predicting

Customer Response Sentiment." *Expert Systems with Applications* 162 (December 30,

2020): 113746. https://doi.org/10.1016/j.eswa.2020.113746.

Bose, Dr Rajesh, P. S. Aithal, and Sandip Roy. "Survey of Twitter Viewpoint on Application

of Drugs by VADER Sentiment Analysis among Distinct Countries." SSRN Scholarly

Paper. Rochester, NY: Social Science Research Network, March 15, 2021.

https://papers.ssrn.com/abstract=3805424.

Conway, Bethany A., Kate Kenski, and Di Wang. "The Rise of Twitter in the Political

Campaign: Searching for Intermedia Agenda-Setting Effects in the Presidential

Primary." *Journal of Computer-Mediated Communication* 20, no. 4 (July 1, 2015):

363–80. https://doi.org/10.1111/jcc4.12124.

Garay, J., R. Yap, and M. J. Sabellano. "An Analysis on the Insights of the Anti-Vaccine

Movement from Social Media Posts Using k-Means Clustering Algorithm and

VADER Sentiment Analyzer." *IOP Conference Series: Materials Science and*

*Engineering* 482 (March 2019): 012043. https://doi.org/10.1088/1757-899X/482/1/012043.

Granger, C. W. J. "Investigating Causal Relations by Econometric Models and Cross-Spectral Methods." *Econometrica* 37, no. 3 (1969): 424–38. https://doi.org/10.2307/1912791.

Guo, Lei, and Maxwell McCombs. "Network Agenda Setting: A Third Level of Media Effects," n.d., 20.

Harder, Raymond A., Julie Sevenans, and Peter Van Aelst. "Intermedia Agenda Setting in the Social Media Age: How Traditional Players Dominate the News Agenda in Election Times." *The International Journal of Press/Politics* 22, no. 3 (July 1, 2017): 275–93. https://doi.org/10.1177/1940161217704969.

Hutto, C., and Eric Gilbert. "VADER: A Parsimonious Rule-Based Model for Sentiment Analysis of Social Media Text." *Proceedings of the International AAAI Conference on Web and Social Media* 8, no. 1 (May 16, 2014): 216–25.

McCOMBS, MAXWELL E., and DONALD L. SHAW. "THE AGENDA-SETTING FUNCTION OF MASS MEDIA*." *Public Opinion Quarterly* 36, no. 2 (January 1, 1972): 176–87. https://doi.org/10.1086/267990.

McCombs, Maxwell, Juan Pablo Llamas, Esteban Lopez-Escobar, and Federico Rey. "Candidate Images in Spanish Elections: Second-Level Agenda-Setting Effects." *Journalism & Mass Communication Quarterly* 74, no. 4 (December 1, 1997): 703–17. https://doi.org/10.1177/107769909707400404.

Ou-Yang, Lucas. *Newspaper3k: Article Scraping & Curation*. Python, 2022.

    https://github.com/codelucas/newspaper.

VentureBeat. "Pew Finds 69% of Twitter's Most Prolific Users Are Democrats," October 15,

    2020. https://venturebeat.com/2020/10/15/pew-finds-69-of-twitters-most-prolific-

    users-are-democrats/.

Rogstad, Ingrid. "Is Twitter Just Rehashing? Intermedia Agenda Setting between Twitter and

    Mainstream Media." *Journal of Information Technology & Politics* 13, no. 2 (April 2,

    2016): 142–58. https://doi.org/10.1080/19331681.2016.1160263.

Schramm, Wilbur. "The Effects of Mass Communications: A Review." *Journalism Quarterly*

    26, no. 4 (December 1, 1949): 397–409.

    https://doi.org/10.1177/107769904902600403.

Shaw, Donald Lewis, and Maxwell E McCombs. *The Emergence of American Political*

    *Issues: The Agenda-Setting Function of the Press*. St. Paul: West Pub. Co., 1977.

Su, Yan, and Porismita Borah. "Who Is the Agenda Setter? Examining the Intermedia

    Agenda-Setting Effect between Twitter and Newspapers." *Journal of Information*

    *Technology & Politics* 16, no. 3 (July 3, 2019): 236–49.

    https://doi.org/10.1080/19331681.2019.1641451.