

Final Project - MATH 2820 (L)

Contributors:

Grant Bowlds: Wrote the code for the data mining, sanitized the tweets and completed the analysis between volume and effect on Dow Jones.

Immanuel John Milton: Wrote much of the final proposal and final project. Looked into completing analysis on other world leaders, prevented by Twitter policies

James Tang: Wrote code to find the common words from President Trump. Analyzed the impact of tweets on the market to the market and on international countries

Introduction:

The main goal of this project is to quantify to the extent possible the relationship between the number of President Trump's tweets and the Dow Jones Industrial Average during the time from Trump's inauguration until 3/31/2020. The data used includes all of Trump's tweets from his inauguration until 3/31/20, as well as the Dow Jones Industrial average from that timeframe; data was collected and compared per day to grant the largest sample size to approach the population size. We also analyzed whether the content of the tweets was relevant to the performance of the market. Our hypothesis is that tweets about the market would have a positive effect on the market, while tweets about foreign powers would have a negative effect on the market. We do not believe that the number of tweets will directly affect the market in a measurable, quantifiable way.

Our R code and data sources are at the bottom of this document for reference. The other code is spread throughout this project.

Model:

Data Gathering

We began by gathering the relevant data. This began with the retrieval of the data that we need to use to test our hypothesis. We first gathered the Tweets that President Trump tweeted using Kaggle user Austin Reese's account that had scrapped all of President Trump's tweets from April 15th, 2019 until January 20th, 2020 and placed the contents of each of the tweets with the associated timestamp into a CSV file. A representative subset of that Excel file is included in the repository below.

```

13 ```{r}
14 library(ggplot2)
15 library(dplyr)
16 library(readr)
17 library(lubridate)
18 library(tidyr)
19 library(wesanderson)
20 library(ggrepel)
21
22
23 tweets <- read.csv(file = "C:/Users/grant/OneDrive/Desktop/MATH 2820L/Twitter-DJI/Data/trumptweets.csv")
24 dji <- read.csv(file = "C:/Users/grant/OneDrive/Desktop/MATH 2820L/Twitter-DJI/Data/^DJI.csv")
25 ~r

```

Our next step was to gather the other piece of data that was essential for this project. We needed the stock information from during this same time period. We decided to use the Dow Jones Industrial Average to best analyze the status of the stock market on that particular day. We once again used an outside resource: Yahoo Finance. Yahoo Finance included a CSV file of the market information of every day including every day of Trump's presidency, which included the timeframe that we were looking at analyzing. A representative subset of that Excel file is included in the repository below.

```

8 ```{r setup, include=FALSE}
9 knitr::opts_chunk$set(echo = TRUE)
10 library(dplyr)
11 library(ggplot2)
12 library(foreign)
13 trumptweets <- read.csv(url("https://raw.githubusercontent.com/grantbowlds/Twitter-DJI/master/Data/trumptweets_processed.csv"))
14 DJI <- read.csv(url("https://raw.githubusercontent.com/grantbowlds/Twitter-DJI/master/Data/%5EDJI.csv"))
15 ~r

```

Data Processing

The next step in this process was to process the data that we received from the two sources. We had gigabytes worth of data that we needed to pare down to be able to provide meaningful analysis. We first did this by only selecting the data from Trump's inauguration date (1/29/2017) to 3/31/2020. Then, using the Tidyverse package features, the list of dates were completed to integrate the dates that President Trump did not tweet. Weekend dates were also removed as the market, thereby Dow Jones, obviously is not open on weekends. A new column was added to the Dow Jones Industrial average CSV file to calculate the daily percent change from open to close of the market. These two data sets were then merged into a single data frame that could be used for visualization and analysis. Tidyverse packages were also used to run analysis and mine text from the tweets themselves. To test some of our hypotheses, dates, where the tweets did not meet our testing criteria, were not used.

```

26 dji <- dji %>%
27   mutate(perchange = (Close-Open)/Open *100)
28
29 dji <- dji %>%
30   mutate(perchange = (Close-Open)/Open *100)
31
32 tweets <- tweets %>%
33   mutate(date = ymd_hms(date))
34 tweets <- tweets %>%
35   mutate(date = round_date(date, "day"))
36
37 dji <- dji %>%
38   mutate(date = ymd(Date))
39
40
41 overlapDates <- interval(ymd(20170120), ymd(20200331))

```

```

43 tweets <- tweets %>%
44   filter(date %within% overlapDates)
45 dji <- dji %>%
46   filter(date %within% overlapDates)
47
48 tweetBD <- tweets %>% group_by(date) %>% tally()
49
50 tweetBD <- tweetBD %>%
51   mutate(date = ymd(date))
52
53 tweetBD %>%
54   mutate(date = as.Date(date))
55
56 tweetBD$n[tweetBD$n>5] <- 5
57 tweetBD$n <- as.integer(tweetBD$n)
58
59 dplr <- left_join(dji, tweetBD, by=c("date"))

```

```

61 right = replicate(804, 0)
62 right <- as.integer(right)
63
64 dplr$n <- coalesce(dplr$n, right)
65
66 nZero <- dplr %>%
67   filter(n == 0)
68 nOne <- dplr %>%
69   filter(n == 1)
70 nTwo <- dplr %>%
71   filter(n == 2)
72 nThree <- dplr %>%
73   filter(n == 3)
74 nFour <- dplr %>%
75   filter(n == 4)
76 fiveUp <- dplr %>%
77   filter(n >= 5)
78
79 suppressMessages(library("dplyr"))

```

We also used the tidytext to identify key features from the content of Trump's tweets. Using a sentiment analysis package, we were able to create a dummy variable, flagging each tweet as either positive or negative. We then mined our tweets to create a list of Trump's most commonly used words and bigrams, which we later utilized to flag tweets containing economic or international relations related terms. Grouping by day, we merged this dataset with our Dow Jones Industrial Average CSV file. We added columns containing the frequency of market mentions per day, the frequency of international mentions per day, and the sentiment ratio (positive tweets divided by total tweets) per day.

```

18 ```{r}
19 # Additional Processing
20 names(trumptweets)[10] <- "sentiment"
21 names(trumptweets)[11] <- "market_check"
22 names(trumptweets)[12] <- "intl_check"
23
24 text_results <- trumptweets %>%
25   mutate(Date = as.Date(date)) %>%
26   group_by(Date) %>%
27   summarise(market_mentions = sum(market_check),
28             intl_mentions = sum(intl_check),
29             sentiment_ratio =
30               length(sentiment[sentiment=="positive"])/ length(sentiment))
31 DJI <- DJI %>%
32   mutate(Date = as.Date(Date))
33
34
35 DJI_text <- right_join(DJI, text_results)
36 DJI_text <- na.omit(DJI_text)
37
38 DJI_text <- DJI_text %>%
39   mutate(net_market = Close - Open) %>%
40   mutate(market_groups = ifelse(market_mentions == 0, 0,
41                                 ifelse(market_mentions >= 1 & market_mentions < 5, 1, 2)),
42          intl_groups = ifelse(intl_mentions == 0, 0,
43                                ifelse(intl_mentions >= 1 & intl_mentions < 5, 1, 2)))
44 ```

```

Data Plots

In order to figure out Trump's specific impact on the market based on his comments on other countries and the market itself, we mined the most common words that Trump said. The relevant code is below with the results of the code.

```

11 # -----
12 # Exploratory Analysis
13 # -----
14
15 # Top Words
16 replace_reg <- "https?://[^\s]+&[&lt;|&gt;|\\bRT\\b"
17 # Top Single Words
18 words <- trumptweets %>%
19   mutate(text = str_replace_all(content, replace_reg, "")) %>%
20   unnest_tokens(word, content, token = "tweets")
21 words <- words %>%
22   anti_join(stop_words, by = "word") %>%
23   filter(!str_detect(word, "realDonaldTrump"))
24
25 words_count <- words %>%
26   group_by(word) %>%
27   count()
28 trump_words <- words_count %>%
29   arrange(-n)
30
31 # -----
32
33 # Top Word Pairs
34 bigrams <- trumptweets %>%
35   mutate(text = str_replace_all(content, replace_reg, "")) %>%
36   unnest_tokens(bigram, content, token = "ngrams", n = 2)
37 bigrams <- bigrams %>%
38   separate(bigram, into = c("first", "second"), sep = " ", remove = FALSE) %>%
39   anti_join(stop_words, by = c("first" = "word")) %>%
40   anti_join(stop_words, by = c("second" = "word")) %>%
41   filter(str_detect(first, "[a-z]") &
42          str_detect(second, "[a-z]")) %>%
43   filter(!str_detect(first, "[.]") &
44          !str_detect(second, "[.]")) %>%
45   filter(!str_detect(first, "http") & !str_detect(first, "donaldrump") &
46          !str_detect(second, "http") & !str_detect(second, "donaldrump"))
47
48 bigrams_count <- bigrams %>%
49   group_by(bigram) %>%
50   count()
51 trump_bigrams <- bigrams_count %>%
52   arrange(-n)

```

word <chr>	n <int>
trump	5059
president	2765
people	2444
donald	1819
country	1765
america	1621
time	1553
obama	1444
dont	1431
run	1144

1-10 of 10 rows

bigram <chr>	n <int>
donald trump	1526
fake news	582
white house	337
crooked hillary	329
witch hunt	319
hillary clinton	312
celebrity apprentice	258
president obama	221
trump tower	199
north korea	197

1-10 of 10 rows

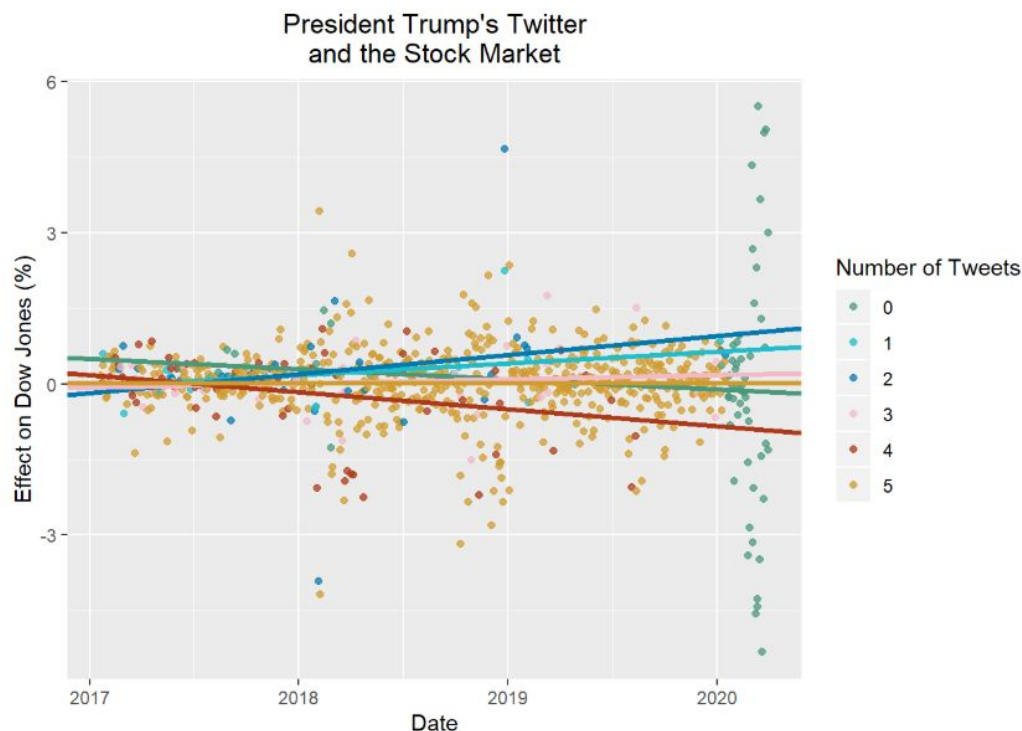
Next, we also created the linear models that combined the Dow Jones Industrial Average and Trump's tweets after the data had been sanitized. The relevant R code and the subsequent plot follows.

```

81 mZero=lm(perchange~date,data=nZero)
82 summary(mZero)
83 confint(mZero)
84 mOne=lm(perchange~date,data=nOne)
85 summary(mOne)
86 confint(mOne)
87 mTwo=lm(perchange~date,data=nTwo)
88 summary(mTwo)
89 confint(mTwo)
90 mThree=lm(perchange~date,data=nThree)
91 summary(mThree)
92 confint(mThree)
93 mFour=lm(perchange~date,data=nFour)
94 summary(mFour)
95 confint(mFour)
96 mFive=lm(perchange~date,data=fiveUp)
97 summary(mFive)
98 confint(mFive)

103 p <- ggplot(dplr, aes(y = perchange, x = date, color = factor(n))) +
104   theme(plot.title = element_text(hjust = 0.5)) +
105   geom_point(alpha = 0.75) +
106   geom_abline(intercept = 10.0730611, slope = -0.0005575, size = 1.2, color = "#49997c") +
107   geom_abline(intercept = -1.111e+01, slope = 6.437e-04, size = 1.2, color = "#1e8449") +
108   geom_abline(intercept = -1.784e+01, slope = 1.029e-03, size = 1.2, color = "#027ab0") +
109   geom_abline(intercept = -3.9308183, slope = 0.0002251, size = 1.2, color = "#f48b3d") +
110   geom_abline(intercept = 16.0250693, slope = -0.0009236, size = 1.2, color = "#ae3918") +
111   geom_abline(intercept = 1.803e-01, slope = -9.026e-06, size = 1.2, color = "#d19c2f") +
112   labs(title = "President Trump's Twitter\and the Stock Market", y="Effect on Dow Jones (%)", x = "Date",
113        color = "Number of Tweets")

```

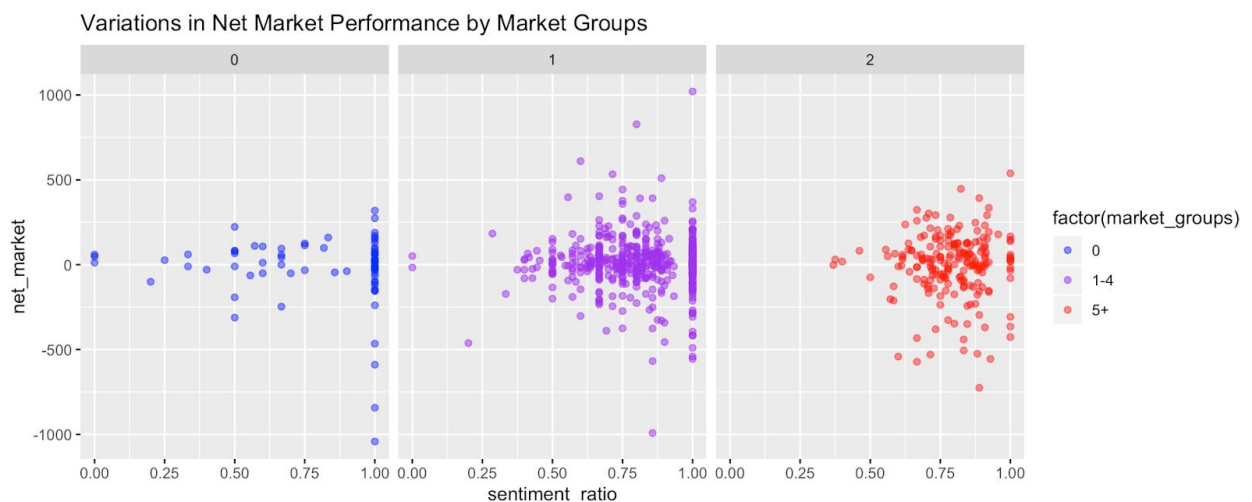


Next, we tried to determine if plots would show a relationship between the market performance or trading volume and President Trump's tweet mentions. We get similar results for international mentions groups vs market performance and market/international mentions groups vs. volume

```

47 ```{r}
48 # Market Groups --> Net Market
49 ggplot(DJI_text, aes(y = net_market, x = sentiment_ratio, color = factor(market_groups))) +
50   geom_point(alpha = 0.5) +
51   labs(title = "Variations in Net Market Performance by Market Groups") +
52   scale_color_manual(labels = c("0", "1-4", "5+"), values = c("blue", "purple", "red")) +
53   facet_wrap( ~ factor(market_groups))
54 # Market Groups --> Volume
55 ggplot(DJI_text, aes(y = Volume, x = sentiment_ratio, color = factor(market_groups))) +
56   geom_point(alpha = 0.5) +
57   scale_color_manual(labels = c("0", "1-4", "5+"), values = c("blue", "purple", "red")) +
58   facet_wrap( ~ factor(market_groups))
59
60 # Intl Groups --> Net Market
61 ggplot(DJI_text, aes(y = net_market, x = sentiment_ratio, color = factor(intl_groups))) +
62   geom_point(alpha = 0.5) +
63   scale_color_manual(values = c("blue", "purple", "red")) +
64   facet_wrap( ~ factor(intl_groups))
65 # Intl Groups --> Volume
66 ggplot(DJI_text, aes(y = Volume, x = sentiment_ratio, color = factor(intl_groups))) +
67   geom_point(alpha = 0.5) +
68   scale_color_manual(values = c("blue", "purple", "red")) +
69   facet_wrap( ~ factor(intl_groups))
70
71 # These plots all seem to suggest that there is no relationship between tweet mentions
72 # and market performance or trading volume
73 ```

```

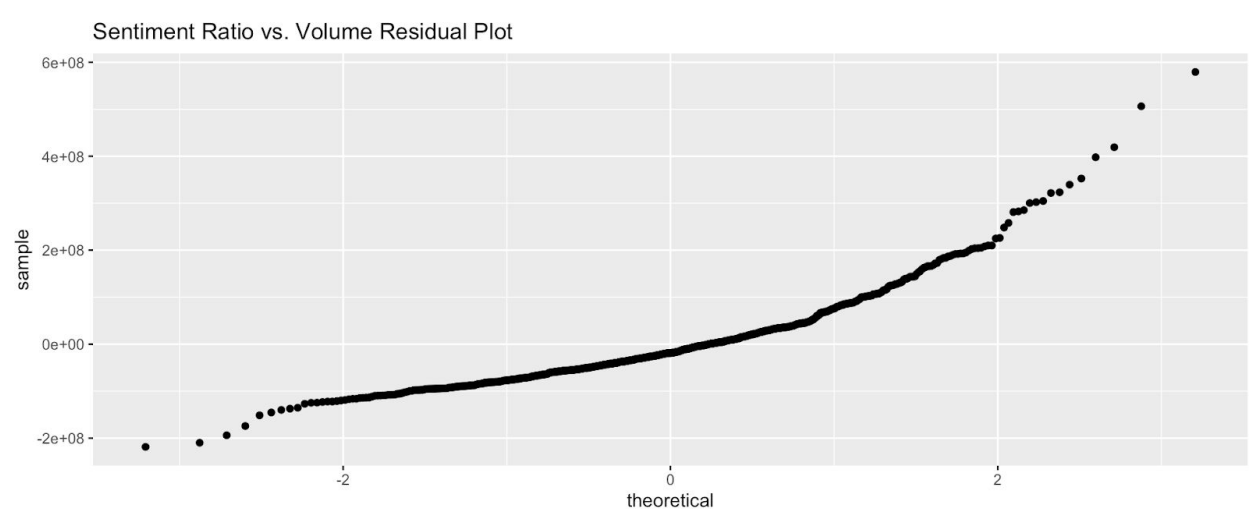
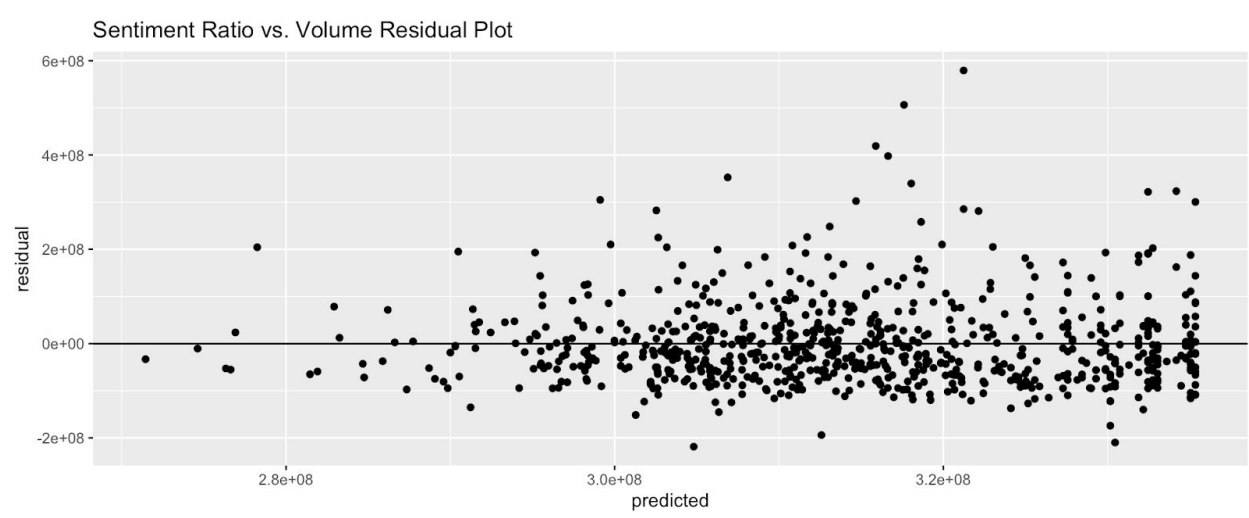
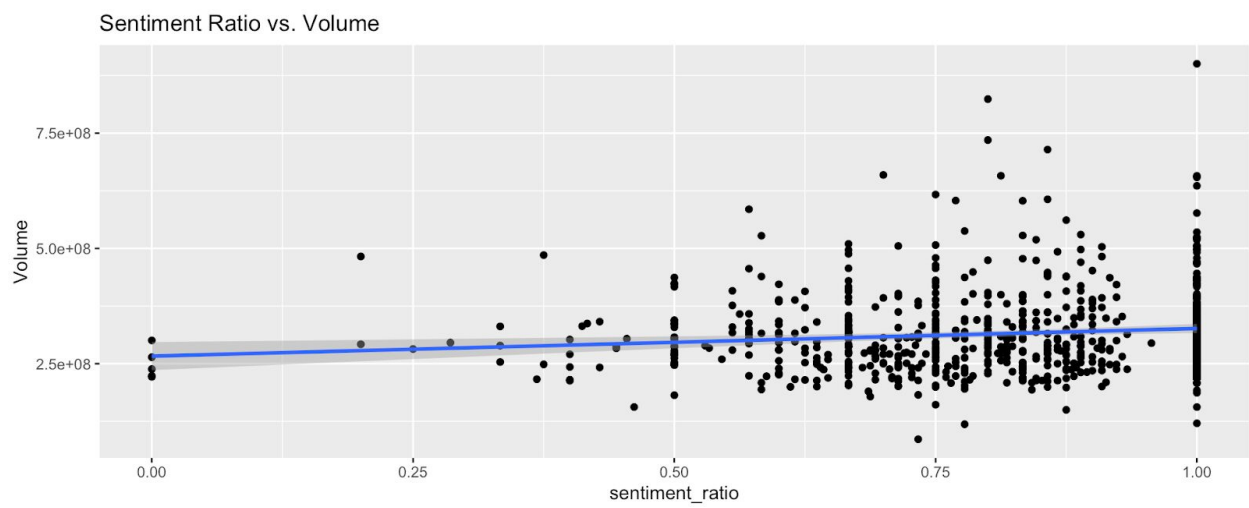


Next, we will build linear models and check their significance between the net market and the sentiment ratio with the market and international mentions from President Trump.

```
76 ```{r}
77 # Regression
78 # Net Market
79 model1 <- lm(net_market ~ sentiment_ratio + market_mentions + intl_mentions, data = DJI_text)
80 summary(model1)
81 # None are significant at 5% level
82
83 # Volume
84 # Remove Outlier
85 DJI_text <- DJI_text[-652,]
86 model2 <- lm(Volume ~ sentiment_ratio + market_mentions + intl_mentions, data = DJI_text)
87 # market mentions, intl mentions not significant
88 # sentiment is significant on volume at the 5% level, let's examine this further
```

We will now build those plots using the sentiment ratio and the tweeting volume to find any potential impact of the numbers of the tweets using the following R code, and the subsequent plot.

```
91 # Tweet sentiment is significant on volume at the 5% significance level
92 ggplot(DJI_text, aes(x = sentiment_ratio, y = Volume)) +
93   geom_point() +
94   labs(title = "Sentiment Ratio vs. Volume") +
95   geom_smooth(method="lm")
96
97 mod2_results <- data.frame(sentiment_ratio = DJI_text$sentiment_ratio,
98                           observed = DJI_text$Volume,
99                           predicted = model2$fitted.values,
100                          residual = model2$residuals)
101 model_volume <- lm(observed ~ sentiment_ratio, data = mod2_results)
102 summary(model_volume)
103 # Now significant at the 1% level
104
105 # residual plot
106 ggplot(mod2_results, aes(y = residual, x = predicted)) +
107   geom_point() +
108   labs(title = "Sentiment Ratio vs. Volume Residual Plot") +
109   geom_hline(yintercept = 0)
110 # qq plot
111 ggplot(mod2_results, aes(sample = residual)) +
112   geom_qq() +
113   labs(title = "Sentiment Ratio vs. Volume Residual Plot")
```



Data Analysis

From looking at the regression lines for the number of tweets less than the average per day (4.37), it is apparent that fewer tweets indicate a more positive slope for percent change of the Dow Jones. Regression and QQ plots were made for each trend line, and it appears all regression models, with the exception of N equal to zero, were normal. For an N equal to one and two, the 97.5% confidence interval was positive on both ends, indicating that we can say with 97.5% confidence that the market will go up when President Trump tweets one or two times. None of the other models had such significance, and R^2 increased with N. The code for the plots are below. The plot itself is above.

```
119 ggplot(mod_results, aes(y = residual, x = predicted)) +
120   geom_point() +
121   geom_hline(yintercept = 0)
122 ggplot(mod_results, aes(sample = residual)) +
123   geom_qq()

127 ggplot(mod_results1, aes(y = residual, x = predicted)) +
128   geom_point() +
129   geom_hline(yintercept = 0)
130 ggplot(mod_results1, aes(sample = residual)) +
131   geom_qq()

135 ggplot(mod_results2, aes(y = residual, x = predicted)) +
136   geom_point() +
137   geom_hline(yintercept = 0)
138 ggplot(mod_results2, aes(sample = residual)) +
139   geom_qq()

143 ggplot(mod_results3, aes(y = residual, x = predicted)) +
144   geom_point() +
145   geom_hline(yintercept = 0)
146 ggplot(mod_results3, aes(sample = residual)) +
147   geom_qq()

151 ggplot(mod_results4, aes(y = residual, x = predicted)) +
152   geom_point() +
153   geom_hline(yintercept = 0)
154 ggplot(mod_results4, aes(sample = residual)) +
155   geom_qq()

159 ggplot(mod_results5, aes(y = residual, x = predicted)) +
160   geom_point() +
161   geom_hline(yintercept = 0)
162 ggplot(mod_results5, aes(sample = residual)) +
163   geom_qq()
```

Our text analysis results indicate that market mention frequency and international mention frequency do not impact market performance or volume. The values for all three categories appear to be similar. However, it is important to note that sentiment ratio appears to increase variation in market performance. While sentiment ratio is not significant on market performance, sentiment ratio does appear to be significant on volume. A simple regression of sentiment ratio on observed volume shows that sentiment ratio is significant at the 1% level. The slope of this regression line is around $5.9e7$, indicating that a one unit increase in sentiment ratio will increase volume by around 59 million units. Additionally, we see that the residual and qq-plot represent an approximately normal distribution, satisfying our assumptions.

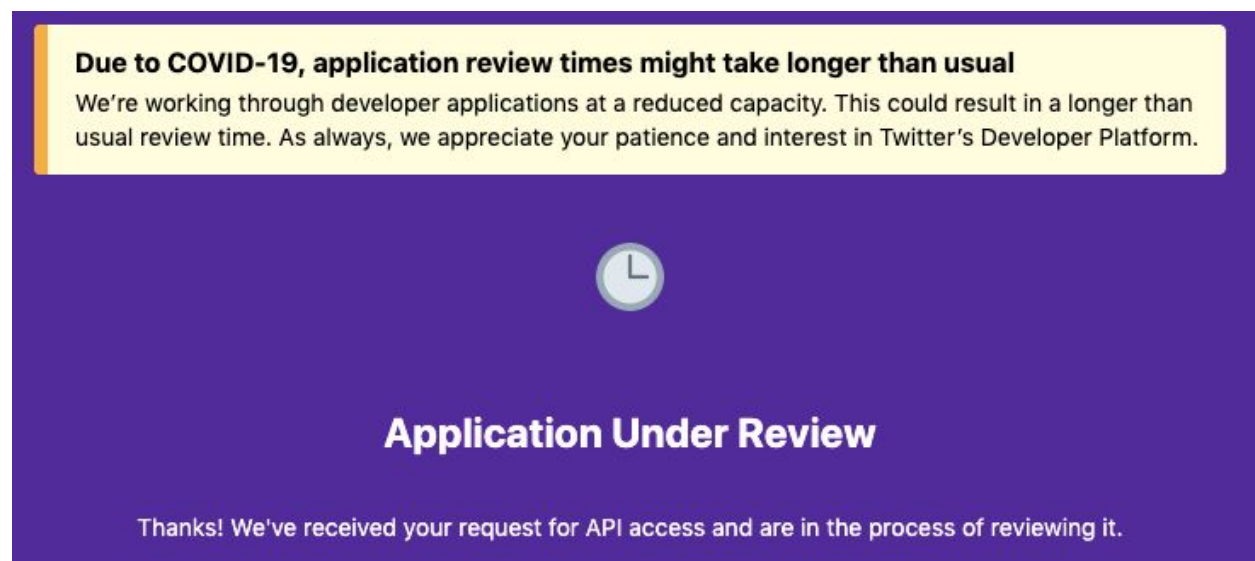
Conclusion:

When Donald Trump only tweets one to two times a day, there is a positive impact on the Dow Jones Industrial Average. When he tweets more than this, the R^2 value becomes so low that we cannot determine a trend. Since the market's value naturally increases with time due to inflation, and we can conclude that Donald Trump's tweets increase market variability, there is a negative correlation between the number of tweets and the success of the stock market on that given day.

Additionally, the content of Trump's tweets appears to impact the market but only slightly. Trump's mention of the economy or international issues does not impact the stock market. However, the more

positive his tweets are in a given day, the more volume the market receives and the greater the variability in market performance.

In terms of our information, we were limited in scope by the lack of data sources on other world leaders. We began looking into former presidents like President Obama and other world leaders but were unable to find a good source of information that included all their tweets with the appropriate time stamps. We attempted to utilize a tweet scraper to solve this problem using Python, however, Twitter requires passwords that only developers are able to access in order to utilize a script that grants access to Tweets. We applied to receive developer access but received the following image from Twitter. (Image underneath next paragraph) We would also have preferred to be able to conduct this research over a longer period of time that didn't include COVID-19 because that is an unprecedented variable that could confound much of this information, however, we found that the most complete reading of the data would include this information; in conclusion on this point, we would have preferred if the COVID-19 pandemic didn't occur, but since it did, we had to use that information.



What we would have liked to do in terms of this information would be to aggregate other world leaders and compare them to both the United States stock exchanges and their own countries. This would allow us to be able to tell if there is an exaggerated effect from President Trump's tweets from others in a similar position. This required the tweet scraper, but that was contingent on Twitter policies but that is a direction this project could go in the future. Another thing we would like to investigate would be whether President Trump only tweets on days that the market is doing well. This is very similar to our initial question of whether his tweets influence the market, but this would require much more data about stock prices at specific points of the day, but bears future investigation.

Keep an eye on your email.

- Be sure to watch the email address **ijohnmilton@gmail.com** as we may request more information to facilitate the review process in the coming days (be sure to check your spam folder as well).
- We review applications to ensure compliance with our [Terms of Service](#) and [Developer policies](#).
- We know that this application process delays getting started with Twitter's APIs. This information helps us protect our platform and serve the health of the public conversation on Twitter. It also informs product investments and helps us better support our developer community.
- You'll receive an email when the review is complete. In the meantime, check out our [documentation](#), explore our [tutorials](#), or

Data Sources:

R and HTML files: <https://github.com/grantbowlds/Twitter-DJI>

Data: <https://github.com/grantbowlds/Twitter-DJI/tree/master/Data>