# RGR Stock Price Forecasting Project

**By Jack Wang**

# PROBLEM STATEMENT

**Stock prices** are hard to predict because they are not only affected by the performance of the underlying companies but also the expectations from the general public. As known, the stock price of firearm companies are highly correlated to the public opinions toward gun control.

My model intends to predict the stock price of one of the largest firearm company in the states, RGR (Sturm, Ruger & Co., firearm company), by using its historical stock price, public opinions toward gun control, and its financial reports to SEC.

# EXECUTIVE SUMMARY

**The goal** of my project is to build a time series regression model that predicts the stock price of RGR.

The data I am using would be historical stock price from **Yahoo Finance**, twitter posts scraped from **Twitter**, **subreddit** posts mentioned about gun control, and also the financial reports to **SEC**. I will do sentiment analysis on the text data and time series modeling on the historical stock price data. The model will be evaluated using MSE.

# DATA COLLECTION

# HISTORICAL STOCK PRICE DATA

Use **Yahoo Finance** to collect historical stock price.

- Open_price: the stock price when market opens

# TWITTER DATA

Use **twitterscraper** to collect historical tweets. Keyword is very specific - gun control.

- tweet_word_count_sum
- tweet_compound_score_sum
- tweets_sum
- tweet_word_count_mean
- tweet_compound_score_mean

/politics

- redd_pol_score_mean
- redd_pol_comment_mean
- redd_pol_compound_mean
- redd_pol_score_sum
- redd_pol_comment_sum
- redd_pol_post_count

/guns

- redd_gun_score_mean
- redd_gun_comment_mean
- redd_gun_compound_mean
- redd_gun_score_sum
- redd_gun_comment_sum
- redd_gun_post_count

# SEC DATA

Use **SEC** website to collect public reports.

- 10-K: company's annual report including financial statement
- 10-Q: company's quarter report including financial statement
- 8-K: company's unscheduled events that are important

# DATA DICTIONARY

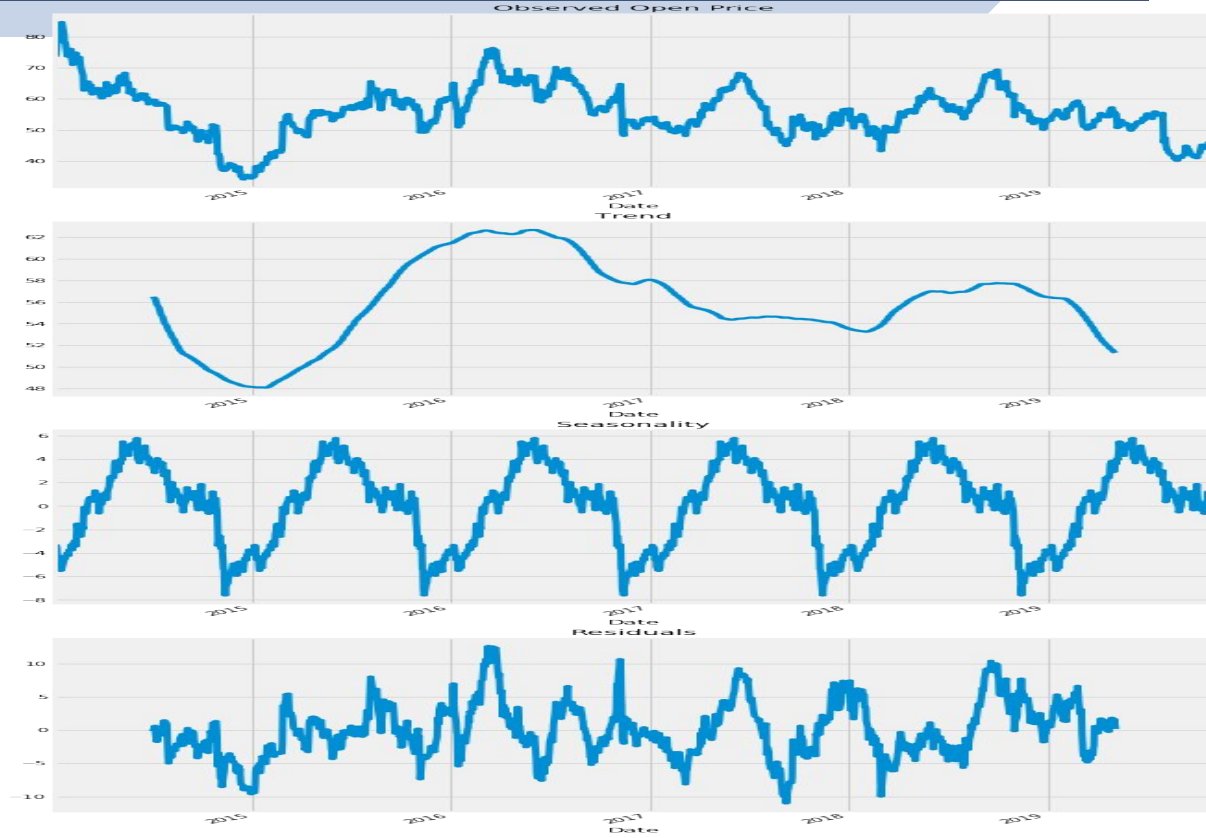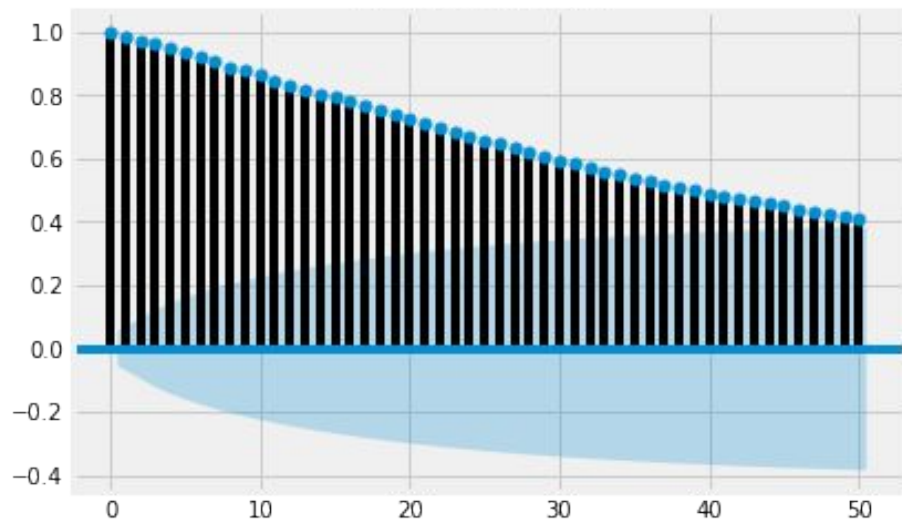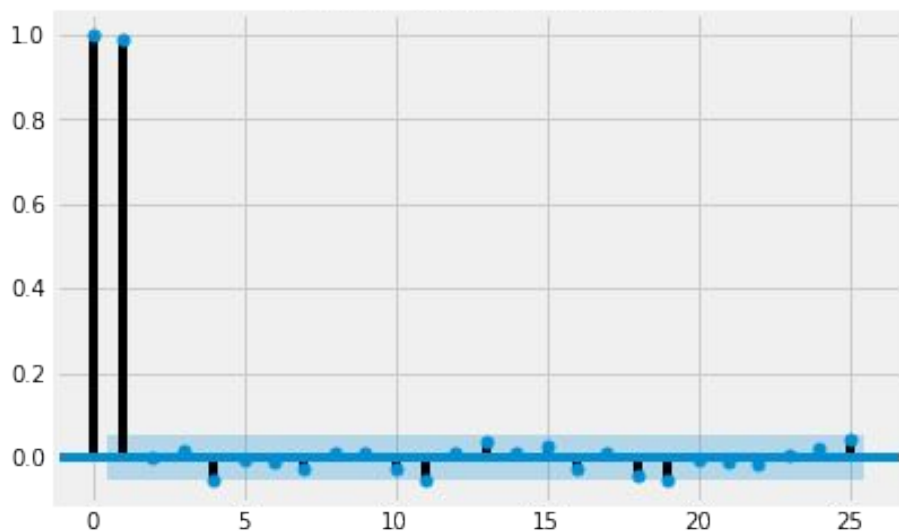| Column Name | Type | Data Collected | Description |
|---|---|---|---|
| open_price | float64 | Yahoo Finance | The open stock price for the given trade day |
| tweet_word_count_sum | int64 | Twitter | The total number of word counts of all tweets mentioning about gun control for the given date |
| tweet_compound_score_sum | float64 | Twitter | The total compound score (sentiment) of all tweets mentioning about gun control for the given date |
| tweets_sum | int64 | Twitter | The total counts of all tweets mentioning about gun control for the given date |
| tweet_word_count_mean | float64 | Twitter | The mean of word counts of all tweets mentioning about gun control for the given date |
| tweet_compound_score_mean | float64 | Twitter | The compound score mean of all tweets mentioning about gun control for the given date |
| redd_gun_score_mean | float64 | /guns subreddit | The subreddit score mean of all the threads on /guns for the given date |
| redd_gun_comment_mean | float64 | /guns subreddit | The comment number mean of all the threads on /guns for the given date |
| redd_gun_compound_mean | float64 | /guns subreddit | The compound score mean of all the threads on /guns for the given date |
| redd_gun_score_sum | float64 | /guns subreddit | The subreddit score sum of all the threads on /guns for the given date |
| redd_gun_comment_sum | float64 | /guns subreddit | The total number of comment counts of all the threads on /guns for the given date |
| redd_gun_post_count | float64 | /guns subreddit | The total number of thread counts of all the threads on /guns for the given date |
| redd_pol_score_mean | float64 | /politics subreddit | The subreddit score mean of all the threads on /politics for the given date |
| redd_pol_comment_mean | float64 | /politics subreddit | The comment number mean of all the threads on /politics for the given date |
| redd_pol_compound_mean | float64 | /politics subreddit | The compound score mean of all the threads on /politics for the given date |
| redd_pol_score_sum | float64 | /politics subreddit | The subreddit score sum of all the threads on /politics for the given date |
| redd_pol_comment_sum | float64 | /politics subreddit | The total number of comment counts of all the threads on /politics for the given date |
| redd_pol_post_count | float64 | /politics subreddit | The total number of thread counts of all the threads on /politics for the given date |
| 10-k | float64 | SEC | The RGR 10-K public reports from SEC |
| 10-q | float64 | SEC | The RGR 10-Q public reports from SEC |
| 8-k | float64 | SEC | The RGR 8-K public reports from SEC |

# EXPLORATORY DATA ANALYSIS

Open Stock Price

# MODELING

Autocorrelation

Partial Autocorrelation

# ARIMA (1, 0, 0)



ARIMA(1,0,0) Predictions

# SARIMA MODEL

I picked the PDQS using gridsearch over AIC score (Akaike Information Criteria).

It measures the simplicity & goodness of fit. The lower the **AIC score** the better.
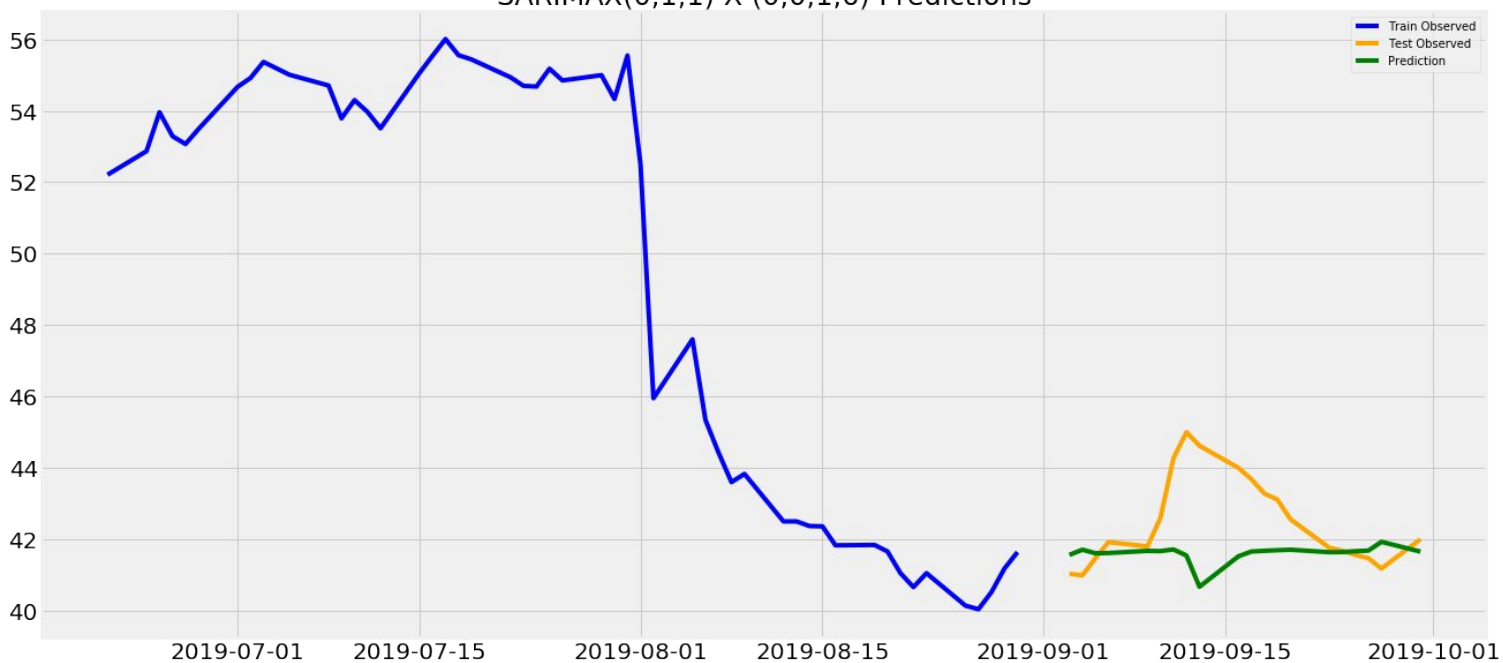
SARIMA(0,1,1)X(0,0,1,5) Predictions

SARIMAX(0,1,1) X (0,0,1,6) Predictions
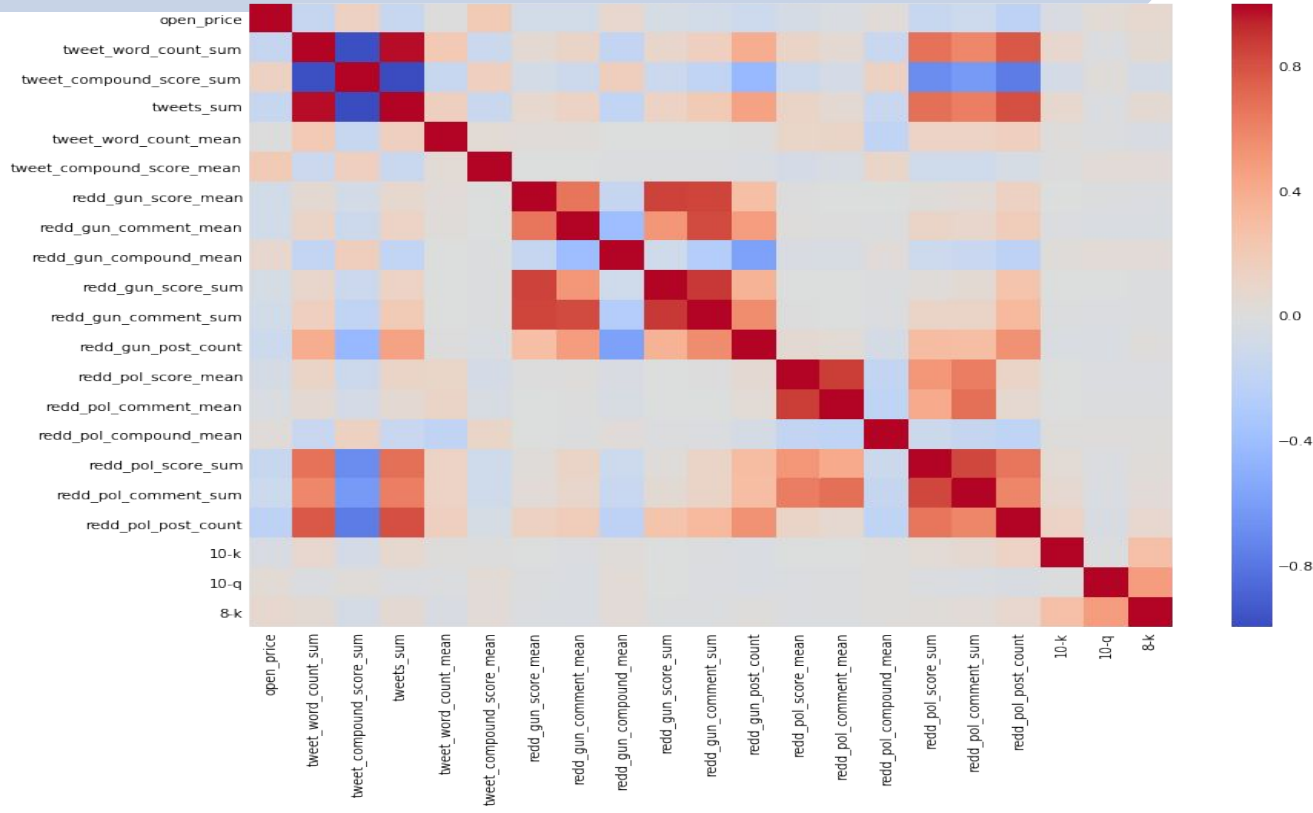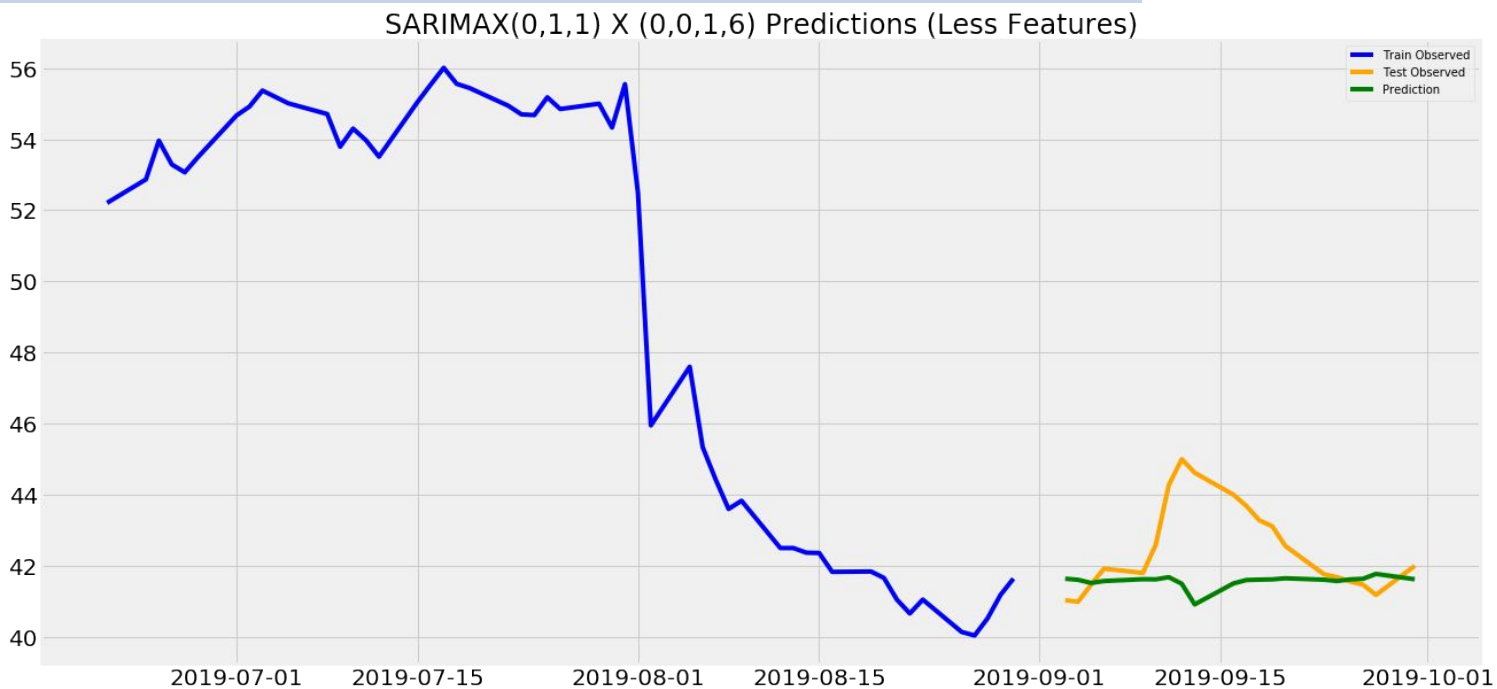
Correlation Heatmap

SARIMAX(0,1,1) X (0,0,1,6) Predictions (Less Features)

SARIMAX(0,1,1) X (0,0,1,6) Predictions (Percentage Difference)

# FLASK DEMO

# CONCLUSION & IMPROVEMENT

Stock prices are hard to predict.

# IMPROVEMENT

1. More data may help. I built the model on data from 2016 October - 2019 October because scraping Twitter data took a lot of time. So I would like to see how much the models will improve if we have more data.
2. More useful exogenous variables. I have included 20 features already, but some of them are not helping the model. It seems like the financial reports 10Q, 10k, and 8K don't help much.
3. Better models. I did not dig into the popular neural network models due to time limitation. I will definitely check out [LSTM] and some other models.

# REFERENCE

# REFERENCE

- Yahoo Finance
- SEC
- twitter
- twitterscraper
- A Guide to Time Series Forecasting with ARIMA in Python 3
- How to Scrap Reddit using pushshift.io via Python
- List of mass shootings in the United States
- Bug in ARIMA predict(): ValueError: Must provide freq argument if no data is supplied #3534
- Using AIC to Test ARIMA Models

# QUESTIONS