# STOCK MARKET PREDICTION USING SENTIMENT AND TIME SERIES ANALYSIS IN R

NAGARAJ Karunashree
VYAS Aryan
SZEPEK Maria Sofie

## Description:

This project studies the possibilities of forecasting stock market prices of firms using the sentiments captured via web scrapping. We have experimented with stock market price of Tesla using sentiment analysis and ARIMA model. An accuracy analysis was also carried out with a R- sqaured value of each of the model to evaluate how each of them faired in the forecasting. The aim is to help reduce participants in loss while investing using the twitter data. The stock data was pulled of the Yahoo Finance API. The sentiments were obtained off the sentences of tweets from twitter. Results obtained has proved that the ARIMA model has good R-Squared value for short term prediction.

The aim of this project is to forecast the price of the Tesla stock by using time series analysis by using the ARIMA model. It uses the advance computing of computers as a tool to extract data from the web and then parse it on tokens for further computation. The parsed text is further analyzed for sentiments thereby giving each a numeric value ranging from -1 to 1. The penultimate goal is to forecast the stock price of the previously mentioned firms by using the tweets and news data sentiments with the already existing livestock price. For this purpose, we have used historical data of opening price, closing price and highest volume traded for Tesla and Moderna. The stock price is affected by various factors like current trade scenarios, people's liking, the company's performance, pandemics, etc.

## Understanding the Objective

The aim of our project is to build a model that predicts the stock price of a given firm while it extracts words from twitter on related topics. Our goal is neither to make billions off the system nor waste billions too. But the objective is to help stock market investors by giving them a direction in taking a decision or not. Whether to buy/hold/sell a stock by providing the result in terms of visualizations.

## Data Collection

This is the process where we used a python script to scrape data off twitter. A sample of the python script used to scrape data is shown below. Also, we used Yahoo finance to get the stock price data for the corresponding interval of time.

```
1  from Scweet.scweet import scrape
2  from Scweet.user import get_user_information, get_users_following, get_users_followers
3
4  data = scrape(words=['amc'], since="2021-01-15", until="2021-08-09", from_account = 'reuters', interval=1, headless=False, display_type="Top", save_images=False, lan
5      resume=False, filter_replies=False, proximity=False)
```

*Fig. 1 Web Scrapping Script (Python)*

## Data Pre-Processing

This is the stage where the acquired data is processed into final datasets to work on. Cleaning the dataset is the focus in this stage. The overview of the dataset is shown below. There are many columns out of which we used the "Embedded text" column as the main feature for sentiment analysis as that contained the tweets of people. The data pre-processing was also used on the stock price data to make it ready and be combined into a value vector.



*Fig. 2 Tweets data (TESLA)*

## Data Preprocessing

To process the data, we use the ARIMA(p,d,q) model. Generally, stock investors use the auto regressive and moving average models to forecast the future trends. Highlights here would be estimation, forecasting and identification. These steps are repeated recursively until an optimal model is identified for prediction. R language provides auto.arima() method to forecast the time series data according to ARIMA(p,d,q) model.

## Forecasting Results

The process of predicting the future by relying upon the past and current data is called as forecasting. Various prediction techniques are used by stock analysts to predict the future stock trends value. The 'forecast' package offered in R was used to predict the future trends which took in values off the sentiment score and past historical stock price data. The "SAS" package of R was used for sentiment analysis thereby giving each sentence a sentiment score. All of these were fed into the ARIMA model which then forecasted the results for time series predictions. It also offered exponential smoothing and space models.

## View and Analyze Results

This is done to evaluate and view the outcome of the model. Screenshots of the evaluation and results are provided further in the section. Investors can view the results and graphs with a comparison view and then invest in the stock. They can use this as an "assistance" to buy/sell/hold a particular stock.

**AIRMA MODEL**



In ARIMA model, the identification is to be accomplished using auto co-relation function and partial auto co-relation function in order to identify p, d and q standards. For any realistic time, sequence generally p, d and q values vary between 0 and 2, but model estimation is executed for all probable combinations of p, d and q values.

**ARIMA() Function in R –**
Arima() function automates the inclusion of a constant. By default, d = 0 or the value of d = 1. A consant will be included if it improved the AICc value; for d > 1 the constant is always omitted. If allowdrift= FALSE is specified, then the constant is only allowed when d = 0.
In ARIMA model, the future value of a variable is a linear combination of past values and past errors, expressed as follows:

$$Y_t = \phi_0 + \phi_1 Y_{t-1} + \phi_2 Y_{t-2} + ... + \phi_p Y_{t-p} + \varepsilon_t - \theta_1 \varepsilon_{t-1} - \theta_2 \varepsilon_{t-2} - ... - \theta_q \varepsilon_{t-q}$$

*Fig. 4 Arima function Math*

Where, $Y_t$ is the actual value and $\square_t$ is random error at t, $\square_I$ and $\square_j$ are the coefficients, p and q

are integers which are often referred to as autoregressive and moving average reciprocal.

**Results and Conclusion**

**Word cloud**
The Word Cloud for Tesla tweets shows us the most common words used in users' tweets. We then calculate the sentiment score of these words thereby giving each a value ranging between 1 to -1. 1 being highly positive and -1 being highly negative.

```
tweet_Table <- tweet_data %>%
  unnest_tokens(word, stripped_text)
data(stop_words)
tweet_Table <- tweet_Table %>%
  anti_join(stop_words)
tweet_Table <- tweet_Table %>%
  count(word, sort = TRUE)
tweet_Table <-tweet_Table %>%
  filter(!word %in% c('replying', 'tesla', 'andothers', 'model', 'tsla', 'car', 'cars',
                'company', 'dont', 'teslas', 'im', 'didnt', 'musk', 'elon', 'gif',
                'time', 'day', 'battery', 'stock', 'drive', 'people', 'ive', 'doesnt',
                'cybertruck', 'texas', 'aug', 'jul', 'range', 'hey', 'nikola',
                'hes', 'thread', 'solar', 'vehicle', 'electric', 'ev', 'lol',
                'driving', 'autopilot', 'world', 'guy', 'month', 'fsd', 'app',
                'energy', 'video', 'money', 'jun', 'service', 'tesmaniancom',
                'companies', 'youtubecom', 'supercharger', 'delivery',
                'factory', 'youre', 'price', 'truck', 'berlin', 'california',
                'china', 'austin', 'ill', 'wait', 'call', 'jul', 'vehicles', 'update', 'share',
                'home', 'yeah', 'giga', 'ceo', 'told', 'guys', 'evs', 'road', 'tech',
                'week', 'tslaq', 'theyre', 'lot', 'twitter', 'answer'))
wordcloud2(tweet_Table,size=0.7,color='random-light', background Color="black")
```

*Fig. 5 Word Cloud (Tesla User Tweets)*

**Time Series Plot of the sentiment Time series**

Plotting the sentiment score w.r.t. the price shows us the correlation of how the stock price of TESLA has been greatly affected by the tweets.

```
plot_data <- subset (sentimentr_news_data, select = c(1))
plot_data <- cbind(plot_data, news_sentiment = sentimentr_news_data$Sentiment,
            tweet_sentiment = sentimentr_tweet_data$Sentiment,
            price = price_data$Returns)
rownames(plot_data) <- plot_data$Date
# Converting data into time-series
news_sentiment <- ts(plot_data$news_sentiment,start=decimal_date(ymd("2020-06-07")),
frequency=52)
tweet_sentiment <- ts(plot_data$tweet_sentiment,start=decimal_date(ymd("2020-06-07")),
frequency=52)
price <- ts(plot_data$price,start=decimal_date(ymd("2020-06-07")), frequency=52)
VAR_data <- cbind(news_sentiment, tweet_sentiment, price)
# Plotting time-series data
forecast::autoplot(VAR_data) +
  ggtitle("Time Series Plot of the Sentiment Time-Series") +
  theme(plot.title = element_text(hjust = 0.5))
```
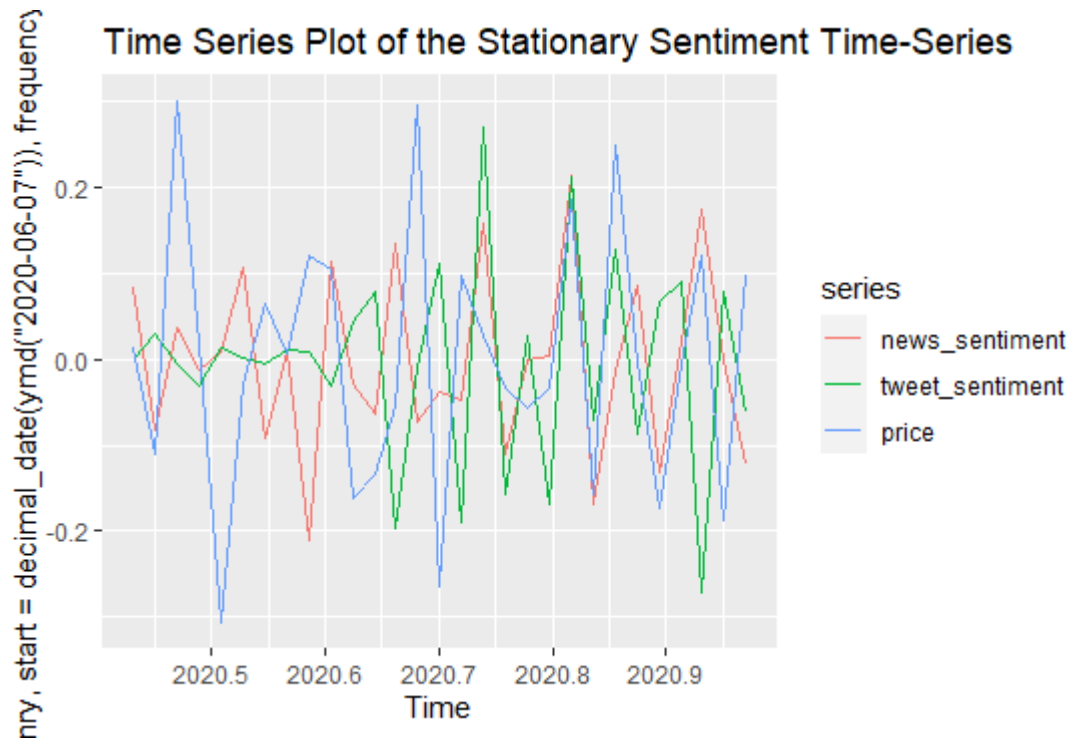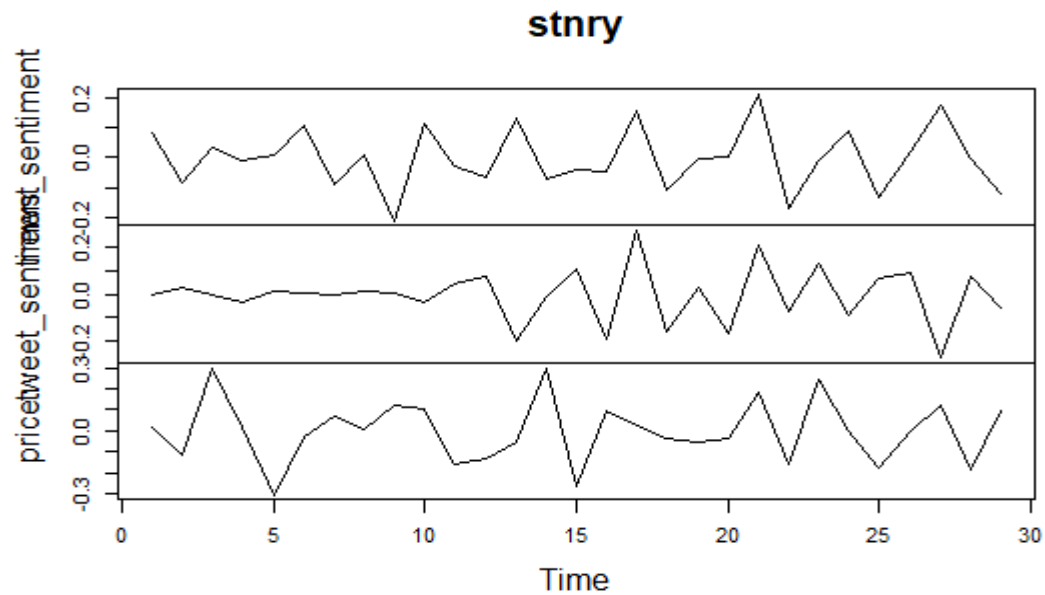
*Fig 6 Graphing the tweets and price (TESLA)*

**Choosing Best lag for our data**

choosing best lag for our data using Vector Auto-regressions VAR.

```
# Choosing best lag for our data
apply(VAR_data, 2, adf.test)
stnry = diffM(VAR_data)
apply(stnry, 2, adf.test)
plot.ts(stnry)
autoplot(ts(stnry, start=decimal_date(ymd("2020-06-07")), frequency=52)) +
  ggtitle("Time Series Plot of the Stationary Sentiment Time-Series")
VARselect(stnry, type = "none", lag.max = 6) # Highest lag order
var.a <- vars::VAR(stnry, lag.max = 6, ic = "AIC", type = "none")
```

*Fig 7 Plot Vector Auto-regressions*

**Forecasting News and tweet sentiments and the price**

The above figure shows us how the news and tweets of users for the word "Tesla" will be forecasted continuing the current and previous trends.

```
# Forecasting the time-series for 8 weeks ahead
fcast = predict(var.a, n.ahead = 8)
par(mar = c(2.5,2.5,2.5,2.5))
plot(fcast)
news_sentiment = fcast$fcst[1];
tweet_sentiment = fcast$fcst[2];
price = fcast$fcst[3];
x = news_sentiment$news_sentiment[, 1]
y = tweet_sentiment$tweet_sentiment[, 1]
z = price$price[, 1]
tail(VAR_data)
x = cumsum(x) + 0.17856389
y = cumsum(y) + (-0.003998803)
z = cumsum(z) + 0.04988139
par(mar = c(2.5,2.5,2.5,2.5))
```
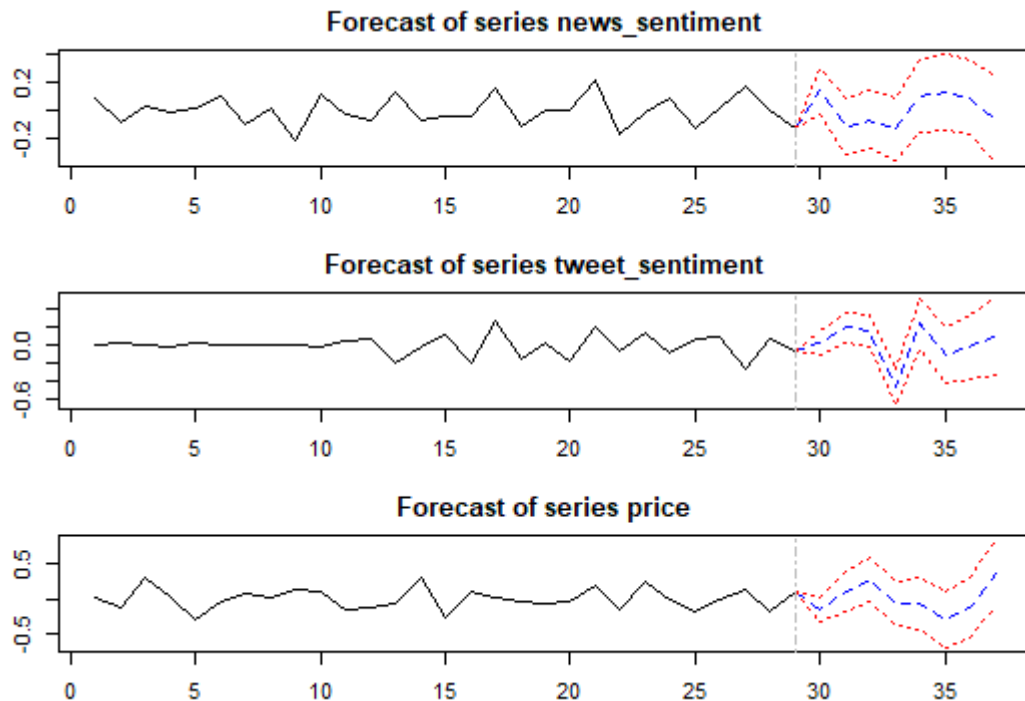
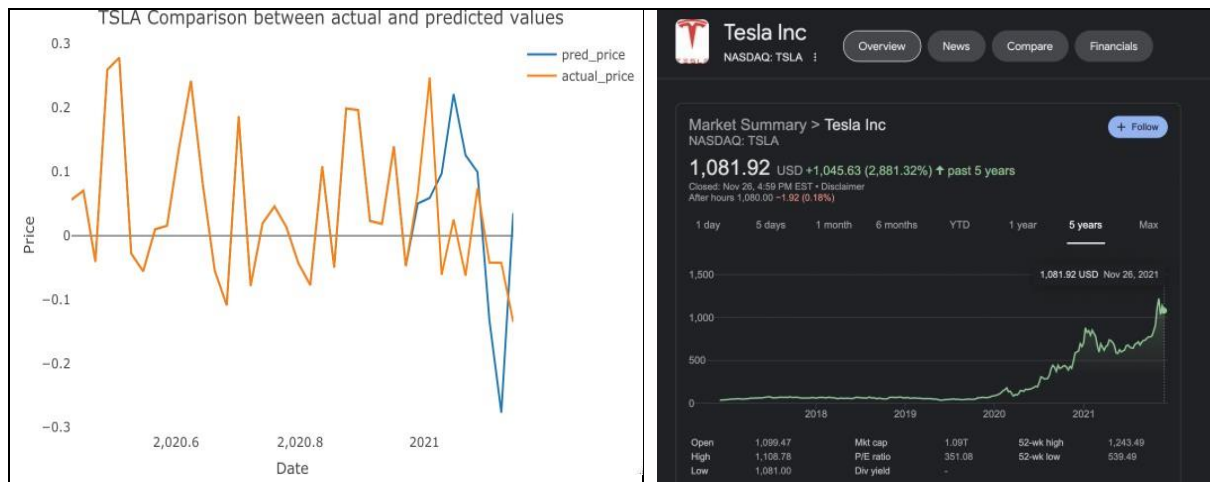*Fig 8 Forecasting News and tweet sentiments and the price*

**price stitched with predicted price**

As we can see that on the (LHS), the blue line indicates the predicted price of TESLA's stock by our model is close to the real time stock price of TESLA (RHS). Post 2021 start, the model was given to predict the stock and it did perfectly.

```
pred_price = ts(c(VAR_data[,3], z), start=decimal_date(ymd("2020-06-07")),
frequency=52)
pred_price_datframe <- as.data.frame(pred_price[1:38])
colnames(pred_price_datframe) <- c("z")
a = zoo(pred_price[1:38])
xyplot(a, grid=TRUE, panel = function(x, y, ...){
  panel.xyplot(x, y, col="red", ...)
  grid.clip(x = unit(31, "native"), just=c("right"))
  panel.xyplot(x, y, col="green", ...) })
```

## Conclusion

In this project we attempted to predict the stock market for TESLA using the historical stock price data and encompassing the impact of sentiments using tweets and news headlines of the same. The experimental results obtained demonstrated the potential of the ARIMA model in short term prediction. This could guide the investors in investing wisely on whether buy/sell/hold that stock.

The factors affecting the R-Squared value are regional investing trends, pandemics, stock's current value, etc.

Investors can take the risk of investment easily when they have an idea about the future value of the stocks.