

Modeling Tesla Stock Price Using Tweets

Tilina Alzaben, Jay Chung, Marion Haney, Divya Rao

04-21-2024

Executive summary

...

1 Introduction

1.1 Research Questions

- Do the frequency and sentiment of Tesla-related tweets have a relationship with Tesla's stock close price?
- Can we use frequency and sentiment of tweets mentioning Tesla to forecast Tesla's stock price?

1.2 Data

For this analysis, we utilized two datasets from Kaggle.

Twitter tweets: 80793 tweets, 25 companies Stock Market Data: Daily close price of stock for same 25 companies Companies are the top 25 "most watched" stocks from Yahoo Finance Filtered to only tweets that mention Tesla or TSLA stock 46% of the tweets in data mentioned TSLA (37422 tweets) Date range: 09-30-2021 to 09-30-2022 (excluding weekends).

...

2 Methods

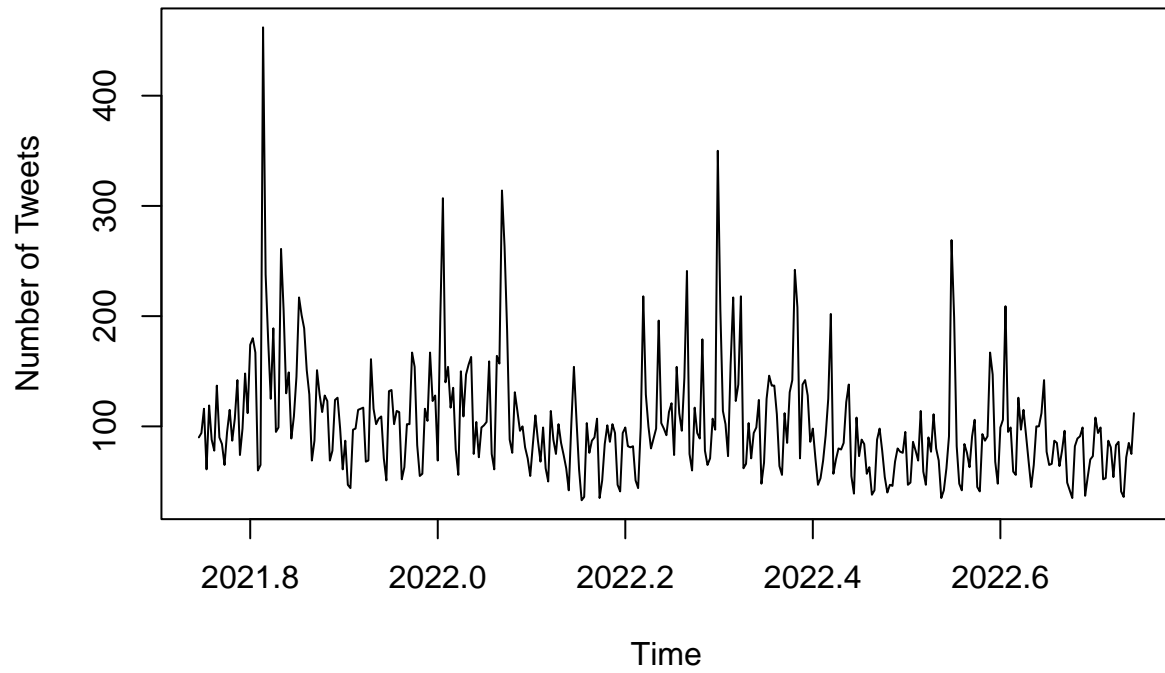
2.1 Exploratory Data Analysis

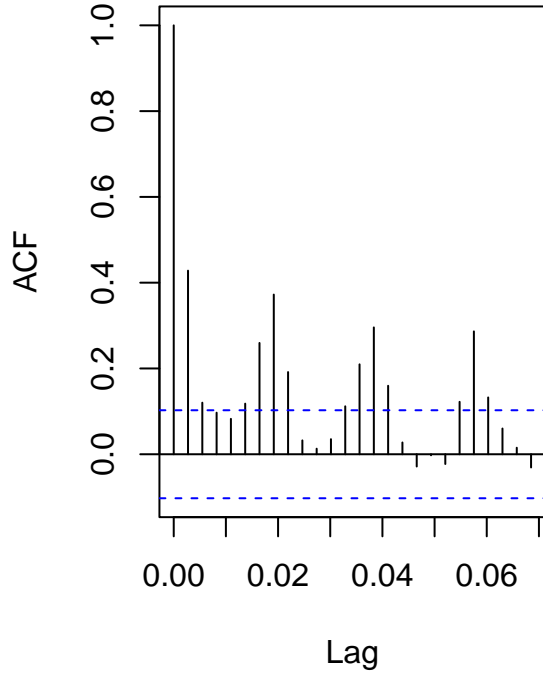
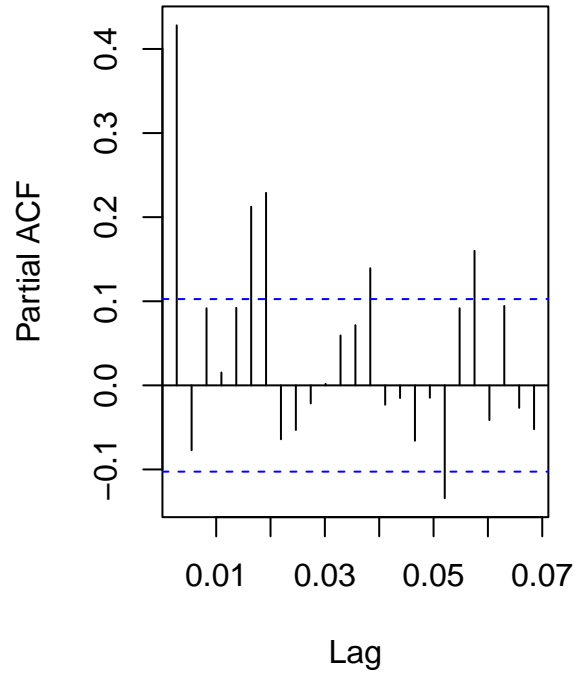
Exploratory Data Analysis (EDA) was conducted for deciding on reasonable data transformations and analysis, determining which type of model to use, and choosing the parameters for these models.

2.2 Number of Tesla Tweets Per Day

Our data included tweets from 09-30-2021 to 09-30-2022 (excluding weekends) which mentioned either Tesla's stock ticker or the company name itself. Our data included 37,422 tweets regarding Tesla over the year. Here, we view the number of tweets mentioning Tesla per day plotted over time along with the ACF and PACF plots of the time series.

Number of Tesla Tweets Per Day



ACF of Number of Tesla Tweets**PACF of Number of Tesla Tweets:**

The time series of the number of Tesla tweets per day shows non-stationary characteristics, with noticeable fluctuations and periodic spikes indicating changes in the mean and variance over time. The ACF plot specifically shows possible seasonality with four clusters of ACF spikes over the time period of the data. These spikes in the number of tweets mentioning Tesla appear to occur approximately every three months.

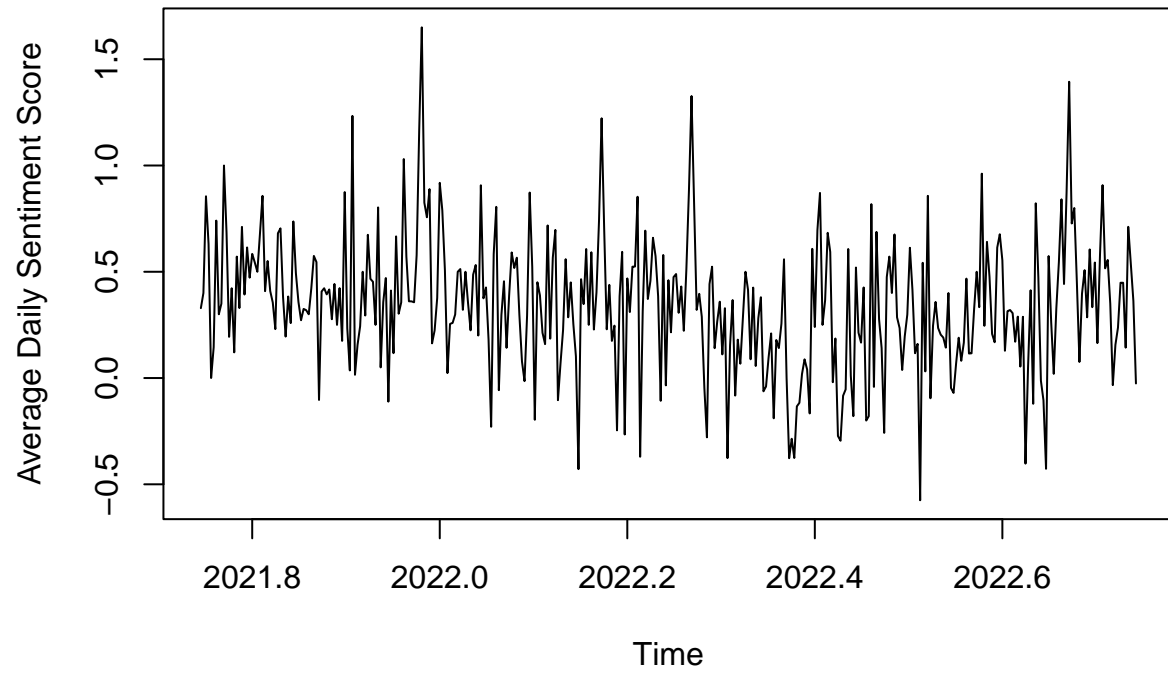
Since the number of Tesla tweets varies per day in a periodic manner, we decided to use averaging in our calculation of daily tweet sentiment. Thus, the fluctuating number of Tesla tweets is accounted for.

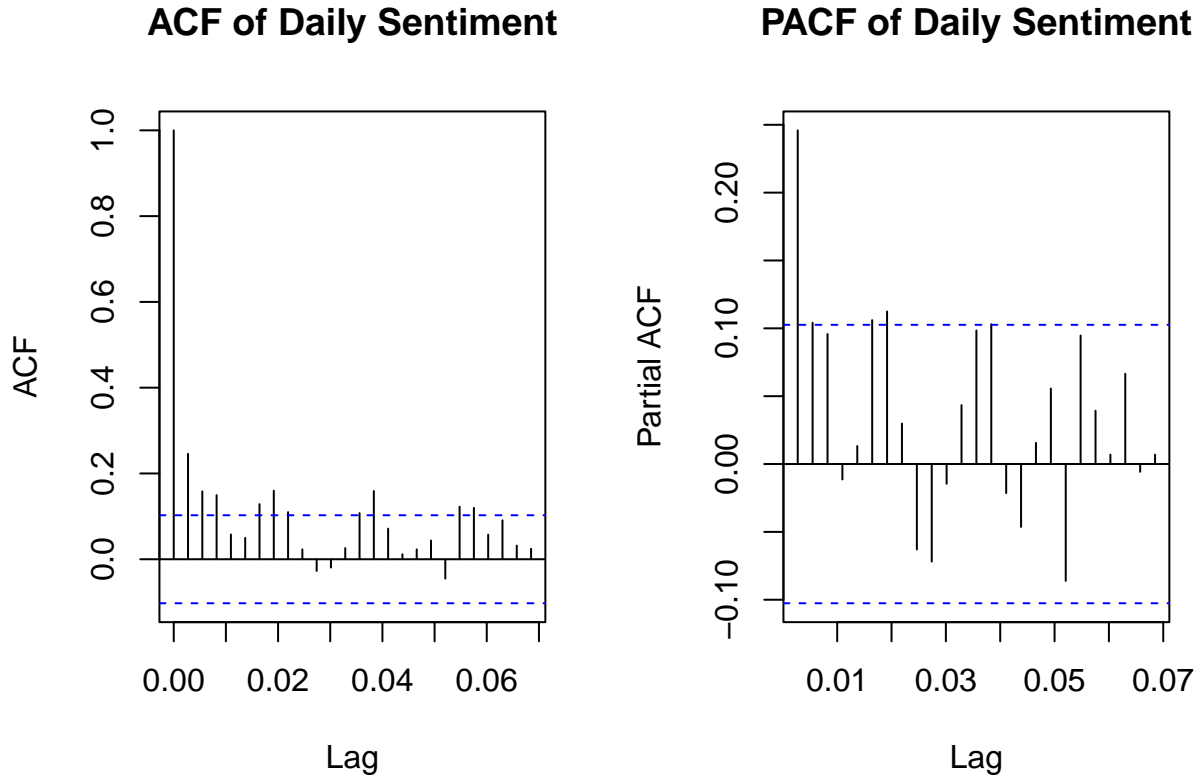
2.3 Daily Sentiment of Tesla Tweets

For our analysis, we calculated a daily score representing the average sentiment (positive or negative) for Tesla tweets posted on the given day. Our approach involved breaking each tweet into a bag of words representation, where each word is given a sentiment score of 1 if it is positive, -1 if it is negative, and 0 otherwise (neutral). These individual word scores were then summed into a sentiment score for each tweet. The sentiment scores per tweet ranged from -9 to 9. Finally, the individual tweet sentiment scores were averaged into a daily tweet sentiment score for the day. The daily sentiment scores ranged from -0.575 to 1.65.

Here, we can view the daily sentiment of tweets over the year of our data along with the ACF and PACF.

Daily Sentiment of Tesla Tweets



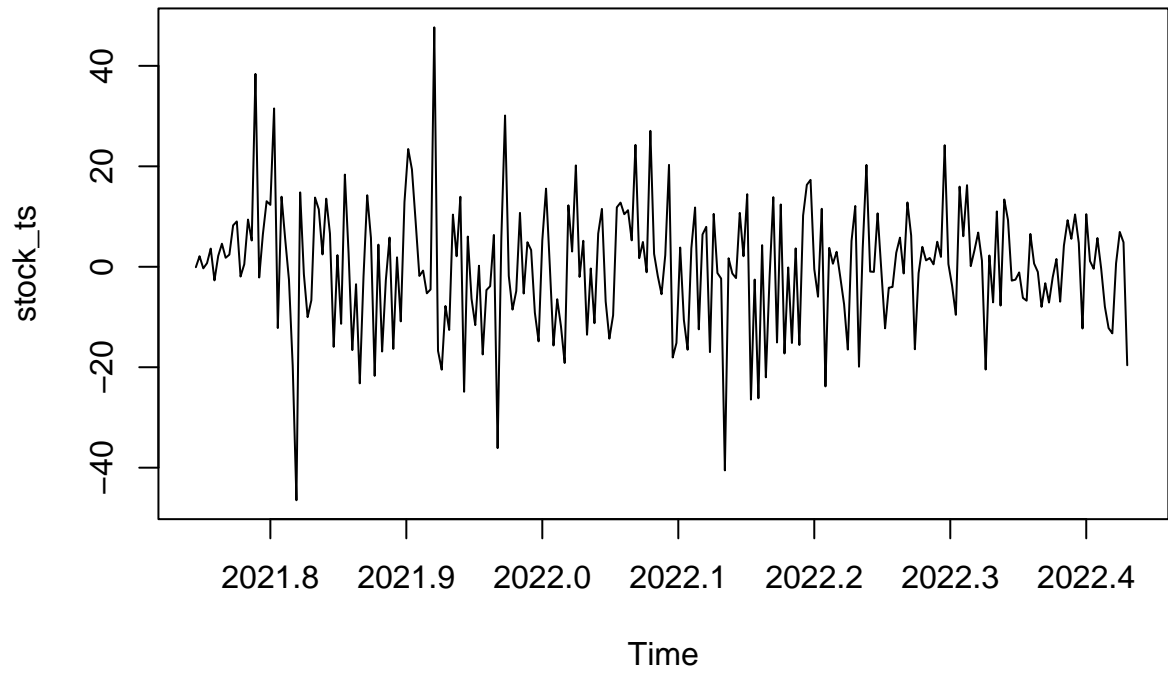


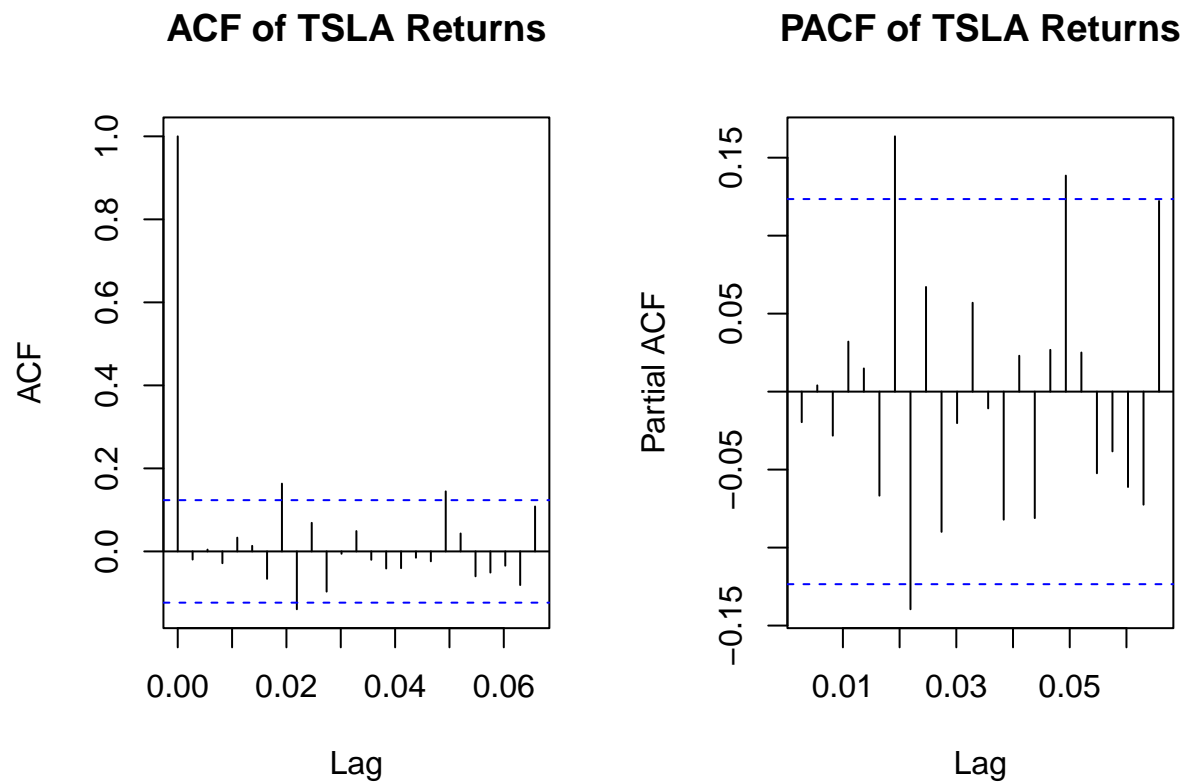
Similarly to the number of Tesla tweets per day, we observe non-stationary characteristics for the daily sentiment of Tesla tweets. we see four sections of ACF spikes approximately 3 months apart, suggesting seasonality. However, the magnitude of the spikes for daily sentiment of tweets is less severe than the magnitude of the corresponding spikes for the number of tweets per day. This is potentially due to the averaging of each individual tweet's sentiment score to produce a daily metric.

2.4 TSLA Stock Price and Returns

We calculated the returns of the TSLA stock by differencing the adjusted close price of TSLA to 1 degree.

TSLA Returns

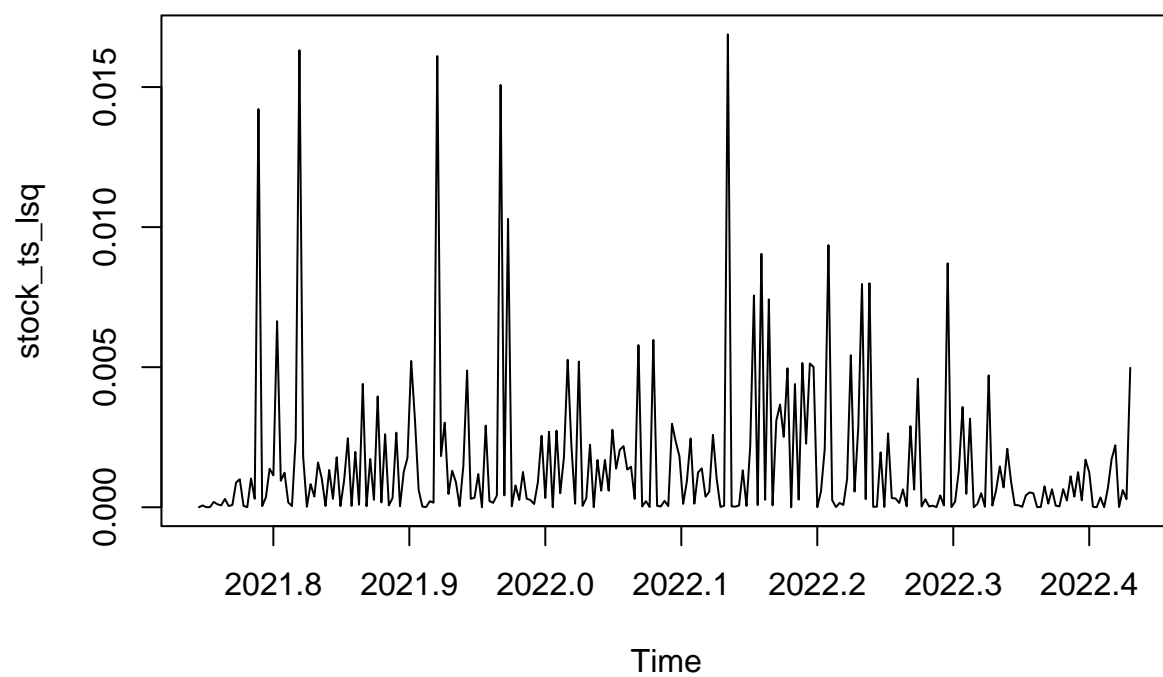


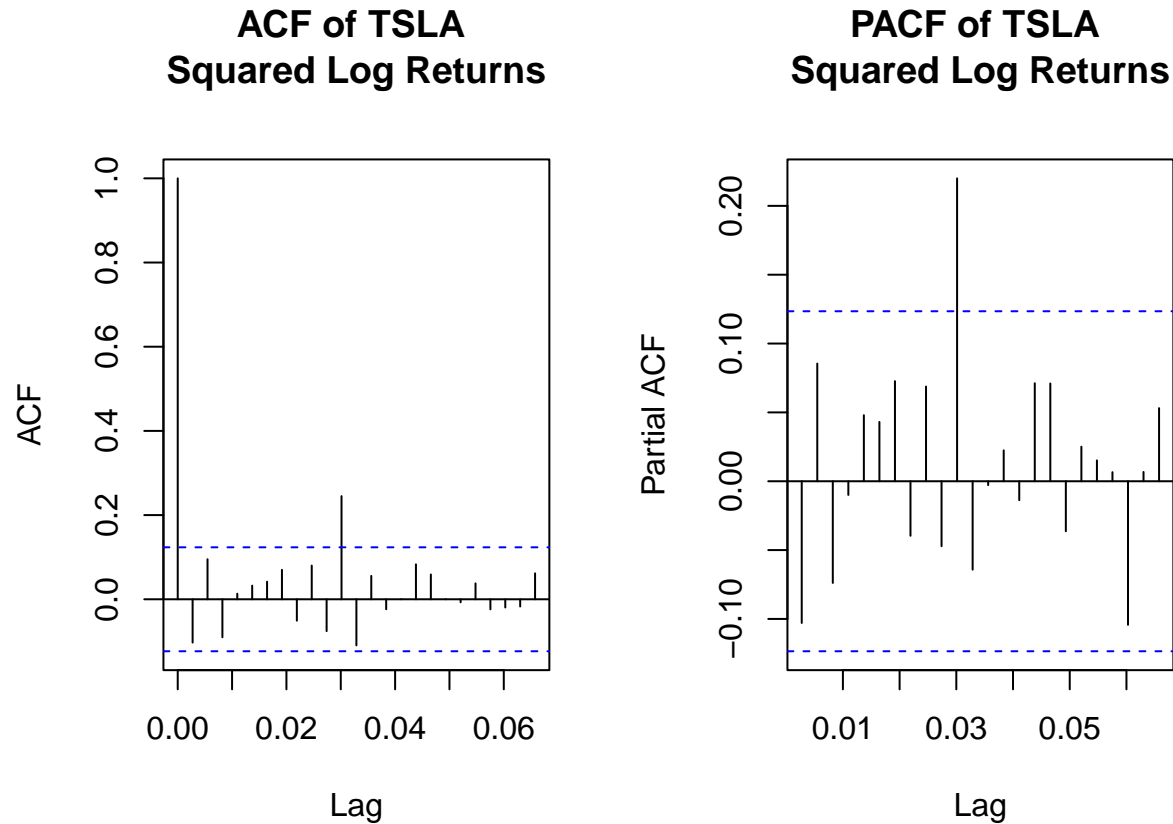


For TSLA returns, we see slight non-stationary characteristics. The ACF and PACF plots exhibit spikes outside the white noise error bounds at around lag 0.02 and 0.05. The time series plot of TSLA returns also exhibits some evidence of heteroskedasticity, with a truncated shape of the time series. The earlier returns exhibit a wider variance than the later returns.

Thus, we decided to square the logged returns to alleviate the slight non-stationary characteristics and heteroskedasticity.

TSLA Squared Log Returns





The square of logged TSLA returns exhibits only one ACF and PACF spike at around lag 0.03 and the plot of the square of logged TSLA returns no longer has a truncated shape.

So, we continued with using the square of logged TSLA returns to represent the measure of TSLA stock price in our model.

2.5 GARCH Model

Based on the non-constant variance and non-stationarity seen from the EDA for daily sentiment, we decided to use a GARCH model. A GARCH (generalized autoregressive conditionally heteroskedastic) model uses values of the past observations and variances to model the variance at time t . GARCH is used commonly with financial data because of the high volatility. Since we calculated the order of the model from running `auto.arima()`, we decided to fit a GARCH(1, 1) model to returns, the difference of the log of the adjusted closing price.

2.5.1 `auto.arima()`

The `auto.arima()` function was applied on the squared returns, which is the difference of the log of the adjusted closing price squared. This helped identify if there was autoregression or moving average pattern in the volatility of the returns and helps determine the order of the GARCH model.

```
## Series: stock$Returns^2
## ARIMA(1,0,1) with non-zero mean
##
## Coefficients:
```

```
##          ar1      ma1      mean
##      -0.6196  0.5105  0.0017
## s.e.    0.1958  0.2102  0.0002
##
## sigma^2 = 7.706e-06:  log likelihood = 1122.91
## AIC=-2237.83  AICc=-2237.66  BIC=-2223.72
##
## Training set error measures:
##              ME          RMSE          MAE          MPE          MAPE          MASE
## Training set -4.434692e-07  0.002759389  0.001733483 -17226.44  17257.83  0.7116065
##              ACF1
## Training set 0.01881578
```

2.6 VAR

We used a VAR (Vector Autoregression) model for feature selection for the two linear models. After running the VARselect() function, we found that the lagged values of l1 and l2 and lagged values of returns.

3 Results

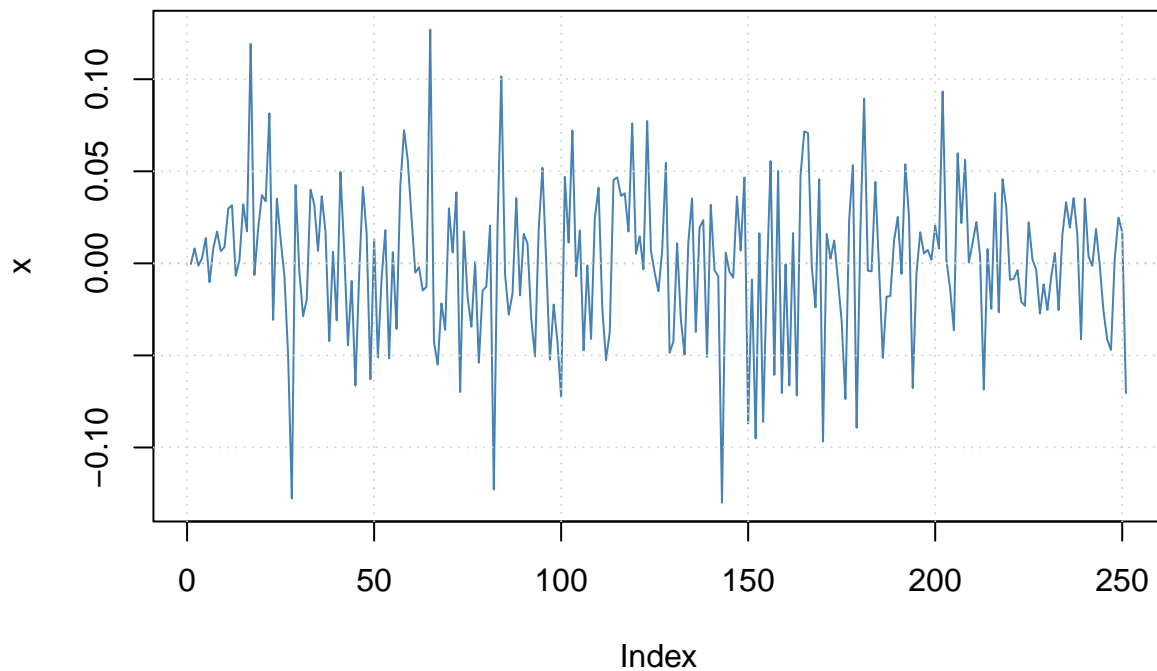
```
##
## Title:
##  GARCH Modelling
##
## Call:
##  garchFit(formula = ~garch(1, 1), data = stock_ts, trace = FALSE)
##
## Mean and Variance Equation:
##  data ~ garch(1, 1)
## <environment: 0x1369cad00>
##  [data = stock_ts]
##
## Conditional Distribution:
##  norm
##
## Coefficient(s):
##           mu          omega          alpha1          beta1
## 0.00034657 0.00007856 0.01438639 0.93899717
##
## Std. Errors:
## based on Hessian
##
## Error Analysis:
##      Estimate Std. Error t value Pr(>|t|)
## mu      3.466e-04 2.557e-03  0.136  0.892
## omega  7.856e-05 8.856e-05  0.887  0.375
## alpha1 1.439e-02 1.997e-02  0.720  0.471
## beta1  9.390e-01 5.494e-02 17.090 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Log Likelihood:
```

```

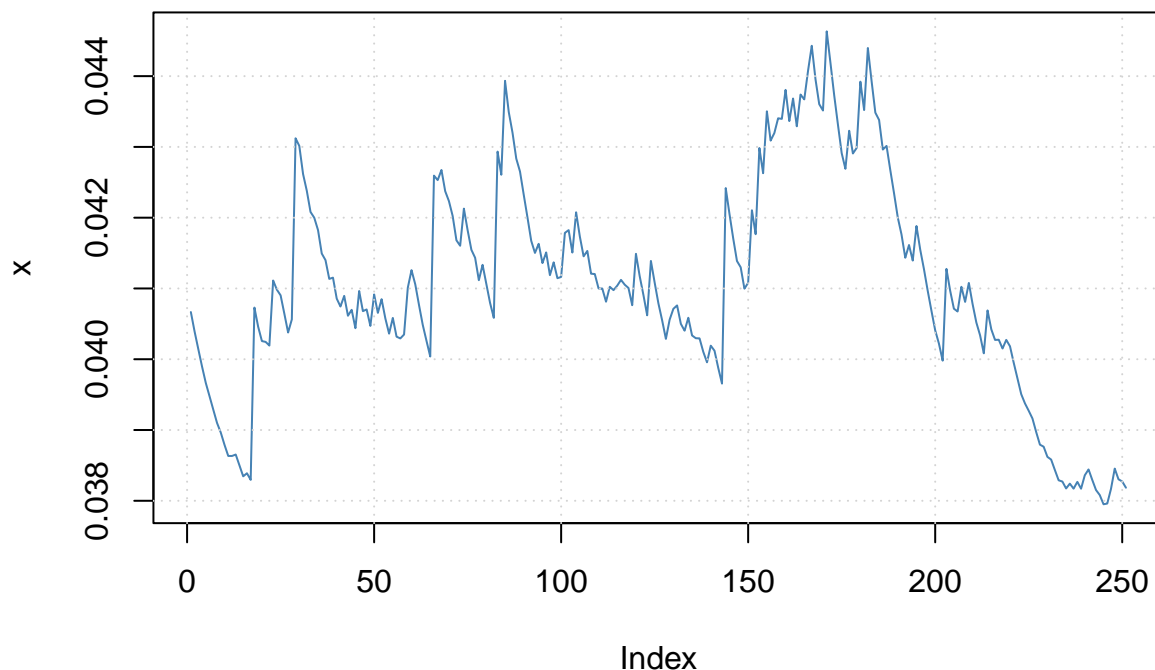
## 448.0723    normalized:  1.785149
##
## Description:
## Sun Apr 21 14:48:15 2024 by user:
##
##
## Standardised Residuals Tests:
##
##           Statistic      p-Value
## Jarque-Bera Test  R    Chi^2 11.8733436 0.002640804
## Shapiro-Wilk Test R      W    0.9837818 0.005887743
## Ljung-Box Test    R    Q(10) 17.7126029 0.060009507
## Ljung-Box Test    R    Q(15) 19.0600447 0.211025838
## Ljung-Box Test    R    Q(20) 27.7444493 0.115588155
## Ljung-Box Test    R^2  Q(10) 11.4358694 0.324582474
## Ljung-Box Test    R^2  Q(15) 30.6866295 0.009672734
## Ljung-Box Test    R^2  Q(20) 33.2456537 0.031704058
## LM Arch Test      R    TR^2  21.5516296 0.042863209
##
## Information Criterion Statistics:
##           AIC      BIC      SIC      HQIC
## -3.538425 -3.482243 -3.538922 -3.515816

```

Time Series



Conditional SD



```
##
## Call:
## lm(formula = Returns ~ daily_sentiment, data = combined_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.129903 -0.023383  0.001996  0.023012  0.126955
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.0006787  0.0037512   0.181   0.857
## daily_sentiment -0.0019186  0.0098397  -0.195   0.846
##
## Residual standard error: 0.04081 on 249 degrees of freedom
## Multiple R-squared:  0.0001527, Adjusted R-squared: -0.003863
## F-statistic: 0.03802 on 1 and 249 DF, p-value: 0.8456

##
## VAR Estimation Results:
## =====
## Endogenous variables: Returns, daily_sentiment
## Deterministic variables: const
## Sample size: 249
## Log Likelihood: 441.396
## Roots of the characteristic polynomial:
```

```

## 0.5583 0.3913 0.2375 0.2375
## Call:
## VAR(y = df, p = optimal.lags)
##
##
## Estimation results for equation Returns:
## =====
## Returns = Returns.l1 + daily_sentiment.l1 + Returns.l2 + daily_sentiment.l2 + const
##
##              Estimate Std. Error t value Pr(>|t|)
## Returns.l1      -0.021839   0.063820  -0.342   0.7325
## daily_sentiment.l1 -0.014608   0.010277  -1.421   0.1565
## Returns.l2        0.025587   0.065023   0.394   0.6943
## daily_sentiment.l2  0.020740   0.010126   2.048   0.0416 *
## const            -0.001562   0.004358  -0.358   0.7204
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.04078 on 244 degrees of freedom
## Multiple R-Squared: 0.02158, Adjusted R-squared: 0.005541
## F-statistic: 1.345 on 4 and 244 DF, p-value: 0.2537
##
##
## Estimation results for equation daily_sentiment:
## =====
## daily_sentiment = Returns.l1 + daily_sentiment.l1 + Returns.l2 + daily_sentiment.l2 + const
##
##              Estimate Std. Error t value Pr(>|t|)
## Returns.l1        1.31367    0.38965   3.371 0.000869 ***
## daily_sentiment.l1  0.16335    0.06275   2.603 0.009799 **
## Returns.l2         0.78679    0.39700   1.982 0.048618 *
## daily_sentiment.l2  0.15634    0.06183   2.529 0.012081 *
## const              0.18730    0.02661   7.040 1.94e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
##
## Residual standard error: 0.249 on 244 degrees of freedom
## Multiple R-Squared: 0.1198, Adjusted R-squared: 0.1054
## F-statistic: 8.302 on 4 and 244 DF, p-value: 2.708e-06
##
##
##
## Covariance matrix of residuals:
##              Returns daily_sentiment
## Returns      0.0016627    -0.0001715
## daily_sentiment -0.0001715    0.0619805
##
## Correlation matrix of residuals:
##              Returns daily_sentiment
## Returns      1.00000    -0.01689
## daily_sentiment -0.01689    1.00000

```

```
##
## Call:
## lm(formula = Returns ~ Returns_l1 + Returns_l2 + daily_sentiment +
##     daily_sentiment_l1 + daily_sentiment_l2, data = combined_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.125592 -0.025207  0.002669  0.024002  0.128883
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.001044   0.004789  -0.218   0.8277
## Returns_l1    -0.018204   0.065415  -0.278   0.7810
## Returns_l2     0.027764   0.065669   0.423   0.6728
## daily_sentiment -0.002767   0.010505  -0.263   0.7925
## daily_sentiment_l1 -0.014156   0.010439  -1.356   0.1763
## daily_sentiment_l2  0.021173   0.010278   2.060   0.0405 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.04085 on 243 degrees of freedom
## (2 observations deleted due to missingness)
## Multiple R-squared:  0.02186, Adjusted R-squared:  0.001734
## F-statistic: 1.086 on 5 and 243 DF, p-value: 0.3686

##
## Call:
## lm(formula = daily_sentiment ~ Returns + Returns_l1 + Returns_l2 +
##     daily_sentiment_l1 + daily_sentiment_l2, data = combined_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.87008 -0.12520  0.00037  0.14850  0.60489
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    0.18714   0.02666   7.019 2.22e-11 ***
## Returns       -0.10313   0.39161  -0.263  0.79250
## Returns_l1     1.31141   0.39049   3.358  0.00091 ***
## Returns_l2     0.78943   0.39788   1.984  0.04837 *
## daily_sentiment_l1 0.16185   0.06313   2.564  0.01096 *
## daily_sentiment_l2 0.15848   0.06247   2.537  0.01182 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2494 on 243 degrees of freedom
## (2 observations deleted due to missingness)
## Multiple R-squared:  0.12, Adjusted R-squared:  0.1019
## F-statistic: 6.63 on 5 and 243 DF, p-value: 8.299e-06
```

4 Discussion

Our goal for this project was to predict Tesla's stock price based on the sentiment of the tweets. For the first linear model, we are predicting returns based on past returns, current daily sentiment, past daily sentiment. We see only one star, which means none of the other features have a statistically significant coefficient other than L2 daily sentiment which has tiny coefficient of 0.02. This means there is a small correlation between returns and L2 daily sentiment. In the second linear model, we are predicting daily sentiment using returns, past returns, past daily sentiment. Almost all of the features are statistically significant value other than current returns. This means that we can predict the sentiment with an adjusted R^2 of 10%, these features account for 10% of the variability of daily sentiment. Past returns and past daily sentiment Granger causes current daily sentiment, i.e. past returns and past daily sentiment can be used to predict the daily sentiment for the current day. However, our original hypothesis was that we would be able to predict returns and stock price using past returns and daily sentiment, which is the opposite direction of Granger causality.

The findings from our model indicate a limited capacity for using tweet sentiment and frequency to predict Tesla's stock price. The significant but minimal correlation found in L2 daily sentiment suggests that while some predictive power exists, it is not strong enough to influence financial decisions effectively. The ability to predict daily sentiment from past stock prices and sentiment highlights potential areas for deeper exploration into the dynamics between social media and financial markets.

4.1 Conclusion

Stock price cannot be predicted based on the current or lagged 1 values of daily sentiment or lagged 1 or 2 values for returns. The coefficient for lagged 2 daily sentiment is statistically significant but too small to be useful. Daily sentiment cannot be predicted based on current stock price but CAN be explained by the lagged 1 or 2 values of the stock price or daily sentiment.

4.2 Future Work

As we move forward, our project will concentrate on enhancing the predictive models and broadening the analytical scope to better understand the relationship between social media activities and stock market behavior. Future efforts will include incorporating broader market performance indicators to understand Tesla's stock movements within the larger financial ecosystem. We will also explore longer temporal analyses to determine whether the effects of negative social media events on stock prices are only temporary or if they have a lasting impact. Additionally, we plan to incorporate tweet engagement metrics such as likes, retweets, and views, which may offer deeper insights into how social media dynamics can influence stock price movements. Through these improvements, we aim to build stronger predictive tools for financial market analysis, refining our approach to use the predictive power of social media sentiment more effectively.

...

Appendix

Notice how the appendix below gathers all the code blocks above and nicely pastes them together.

```
#####  
# STYLE EDITS: IGNORE THIS  
#####  
  
# normally you'll want to include this with the libraries at the beginning of your document  
knitr::opts_chunk$set(message = FALSE) # include this if you don't want markdown to knit messages  
knitr::opts_chunk$set(warning = FALSE) # include this if you don't want markdown to knit warnings  
knitr::opts_chunk$set(echo = FALSE) # set echo = FALSE to hide code from output  
suppressMessages(library(tidyverse))  
suppressMessages(library(lubridate))  
suppressMessages(library(tidytext))  
suppressMessages(library(textdata))  
suppressMessages(library(dplyr))  
suppressMessages(library(quantmod))  
suppressMessages(library(fGarch))  
# Gather stock tweets data  
tweets <- read.csv("/Users/marionhaney/Library/Mobile Documents/com~apple~CloudDocs/CMU/36671 Time Series/tweets.csv")  
# Retaining the year, month, day  
tweets$day <- as.Date(tweets$Date, "%Y-%m-%d %H:%M:%S")  
# Filter to just Tesla tweets  
tesla <- filter(tweets, tweets$Company.Name == "Tesla, Inc.")  
num_tweets_tesla <- data.frame(table(tesla$day))  
names(num_tweets_tesla) <- c("day", "num_tweets")  
# Time series for number of Tesla tweets per day  
tesla_ts <- ts(num_tweets_tesla$num_tweets,  
               start = c(2021, 273),  
               frequency = 365)  
plot(tesla_ts, main = "Number of Tesla Tweets \n Per Day",  
     ylab = "Number of Tweets")  
par(mfrow=c(1,2))  
acf(tesla_ts, main = "ACF of Number of Tesla Tweets")  
pacf(tesla_ts, main = "PACF of Number of Tesla Tweets")  
df <- read.csv("/Users/marionhaney/Library/Mobile Documents/com~apple~CloudDocs/CMU/36671 Time Series/tweets.csv")  
  filter(Stock.Name == 'TSLA') %>%  
  mutate(Tweet.ID = row_number()) %>%  
  dplyr::select(Tweet.ID, Date, Tweet)  
#dim(df)  
#names(df)  
  
df$Date <- ymd(substr(df$Date, 1, 10))  
tweets <- data.frame(df)  
  
# Sentiment analysis  
map_bing_sentiment <- function(sentiment) {  
  ifelse(sentiment %in% c("positive"), 1, ifelse(sentiment %in% c("negative"), -1, 0))  
}  
  
map_nrc_sentiment <- function(sentiment) {  
  nrc_positive_sentiments <- c("positive", "anticipation", "surprise", "trust", "joy")  
  nrc_negative_sentiments <- c("negative", "anger", "disgust", "fear", "sadness")
```



```

    ifelse(sentiment %in% nrc_positive_sentiments, 1,
           ifelse(sentiment %in% nrc_negative_sentiments, -1, 0))
  }

tweet_tokens <- tweets %>%
  unnest_tokens(word, Tweet)

sentiments <- get_sentiments("bing") %>% mutate(sentiment_score = map_bing_sentiment(sentiment))

tweets_sentiment <- tweet_tokens %>%
  inner_join(sentiments, by = "word", relationship = "many-to-many") %>%
  distinct(Tweet.ID, Date, word, .keep_all = TRUE)

tweets_sentiment_summary <- tweets_sentiment %>%
  group_by(Tweet.ID, Date) %>%
  summarise(sentiment_score = sum(sentiment_score, na.rm = TRUE), .groups = "drop")

daily_sentiment <- tweets_sentiment_summary %>%
  group_by(Date) %>%
  summarise(daily_sentiment = mean(sentiment_score))
sentiment_ts <- ts(daily_sentiment$daily_sentiment,
                  start = c(2021, 273),
                  frequency = 365)
plot(sentiment_ts, main = "Daily Sentiment of Tesla Tweets",
     ylab = "Average Daily Sentiment Score")
par(mfrow=c(1,2))
acf(sentiment_ts, main = "ACF of Daily Sentiment")
pacf(sentiment_ts, main = "PACF of Daily Sentiment")
df <- read.csv("/Users/marionhaney/Library/Mobile Documents/com~apple~CloudDocs/CMU/36671 Time Series/t
              filter(Stock.Name == 'TSLA') %>%
              dplyr::select(Date, Adj.Close)
df$Date <- as.Date(df$Date)
df$Returns <- c(diff(df$Adj.Close), NA)
stock_ts <- ts(df$Returns,
               start = c(2021, 273),
               frequency = 365)
df$Returns <- c(diff(log(df$Adj.Close)), NA)
stock_ts_lsqr <- ts(df$Returns^2,
                   start = c(2021, 273),
                   frequency = 365)

stock <- data.frame(df)
stock <- stock %>% na.omit()
# EDA TSLA returns
plot(stock_ts, main="TSLA Returns")
par(mfrow=c(1,2))
acf(stock_ts, main="ACF of TSLA Returns", na.action = na.pass)
pacf(stock_ts, main="PACF of TSLA Returns", na.action = na.pass)
plot(stock_ts_lsqr, main="TSLA Squared Log Returns")
par(mfrow=c(1,2))
acf(stock_ts_lsqr, main="ACF of TSLA \nSquared Log Returns", na.action = na.pass)
pacf(stock_ts_lsqr, main="PACF of TSLA \nSquared Log Returns", na.action = na.pass)
suppressWarnings(library(forecast))

```

```

arma_rt_squared <- auto.arima(stock$Returns^2, max.p = 5, max.q = 5, max.order = 10,
                             stationary = T, seasonal = F, trace = F,
                             stepwise = F, approximation = F)

summary(arma_rt_squared)
stock_ts <- ts(stock$Returns,
               start=c(2021,9),
               frequency=365)
garch_model <- garchFit(~ garch(1,1), data=stock_ts, trace=FALSE)
summary(garch_model)

par(mfrow=c(1,1))
plot(garch_model, which = 1)
plot(garch_model, which = 2)

par(mfrow=c(1,1))
combined_data <- left_join(stock, daily_sentiment, by = "Date")
#cor(combined_data$daily_sentiment, combined_data$Returns, use = "complete.obs")
model <- lm>Returns ~ daily_sentiment, data = combined_data)
summary(model)
suppressMessages(library(vars))
df <- combined_data[, c("Returns", "daily_sentiment")]
# VARselect
lag.select <- VARselect(df,
                        lag.max = 30,
                        type = "both")
optimal.lags <- lag.select$selection['AIC(n)']

# Fit the VAR model
var.model <- VAR(df, p = optimal.lags)

summary(var.model)
suppressMessages(library(dplyr))

combined_data <- combined_data %>%
  arrange(Date) %>%
  mutate(
    Returns_l1 = lag>Returns, 1),
    Returns_l2 = lag>Returns, 2),
    daily_sentiment_l1 = lag(daily_sentiment, 1),
    daily_sentiment_l2 = lag(daily_sentiment, 2)
  )

model <- lm>Returns ~ Returns_l1 + Returns_l2 + daily_sentiment + daily_sentiment_l1 + daily_sentiment_l2)
summary(model)
model <- lm(daily_sentiment ~ Returns + Returns_l1 + Returns_l2 + daily_sentiment_l1 + daily_sentiment_l2)
summary(model)

```