

A candlestick chart on a dark background. The chart shows price movement with green candles for upward movement and red candles for downward movement. Two trend lines are drawn: a blue line sloping downwards from the top left to the top right, and an orange line sloping upwards from the bottom left to the top right. A blue arrow points downwards at the intersection of the two lines, and a yellow arrow points upwards slightly below that intersection.

Capstone Project 2

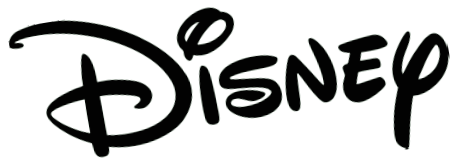
Twitter Sentiment Analysis and Stock Prediction



Introduction

Problem Statement: Analyze stock market movements using Twitter sentiment analysis to find the correlation between “public sentiment” and “market sentiment”


The price stocks trade at seem to be determined more by the human perception of the stock. Behavioral economics states that the emotions and moods of individuals affect their decision making process. Twitter can be used to gauge the public sentiment and possibly predict stock price movements. This study will focus on 4 individual stocks: Netflix (\$NFLX), Disney (\$DIS), Amazon (\$AMZN), and Google (\$GOOGL).

The Netflix logo, featuring the word "NETFLIX" in white, bold, sans-serif capital letters on a red rectangular background.The Disney logo, featuring the word "Disney" in its signature black script font.The Amazon logo, featuring the word "amazon" in a black, lowercase, sans-serif font with a curved orange arrow underneath it.The Google logo, featuring the word "Google" in its multi-colored, sans-serif font.

Data Collection

Stored all my data into a SQLite Database

Twitter Data

Used Twitter API Standard Search to collect tweets from 5-15-19 to 6-26-19. The following are the queries used: 

- Netflix: '@Netflix OR \$NFLX OR Netflix'
- Disney: '@Disney OR @ESPN OR @ABCnetwork OR @Pixar OR @Marvel OR \$DIS'
- Amazon: '@Amazon OR @PrimeVideo OR @awscloud OR @TwitchPrime OR @Alexa OR @WholeFoods OR \$AMZN'
- Google: '@Google OR @Android OR @Waymo OR \$GOOGL'

Tweets Dataframe
created_at
tweet
follower_count
pos_sent
neu_sent
neg_sent
compound_sent
sentiment
Company

Twitter Sentiment Analysis

1. Clean the tweet so that it processes more accurately

- a. Regular Expression library use to locate text strings and remove them
- b. String to be removed:

- i. User mention
- ii. Hyperlinks
- iii. Hashtag sign
- iv. “RT”

'RT @Google: Toy Story is back. See the latest Toy Story 4 trailer #WithALittleHelp from Google → <https://t.co/np6XbygVvi> <https://t.co/Hnpmy...>'



' Toy Story is back See the latest Toy Story 4 trailer WithALittleHelp from Google '

2. VADER (Valence Aware Dictionary and sEntiment Reasoner) as sentiment tool

- a. Lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in social media.
- b. Provides positive, negative, and neutral score. Then computes a compound score which sums the valence scores of each word in the lexicon, adjusted according to the rules, and then normalized to be between -1 (most extreme negative) and +1 (most extreme positive).
 - i. Compound Score > 0.05 = Positive Tweet
 - ii. Compound Score < -0.05 = Negative Tweet
 - iii. $-0.05 < \text{Compound Score} < 0.05$ = Neutral Tweet

Data Collection

Stock Data

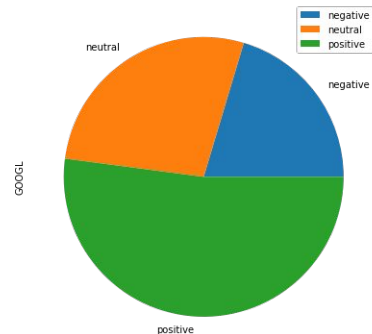
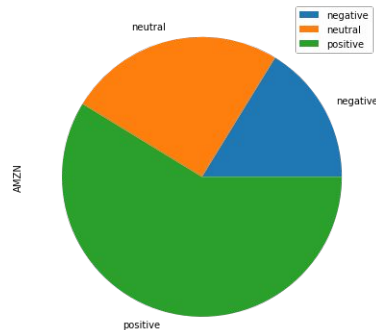
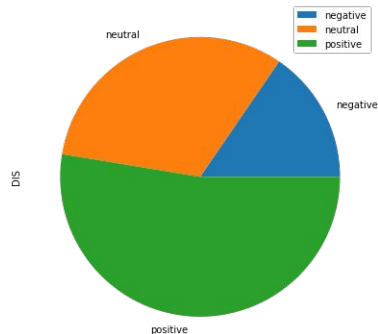
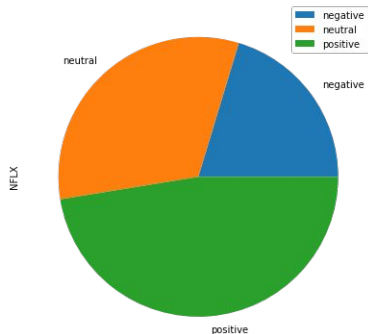
Used Alpha Vantage API to retrieve stock data for all 4 stocks from the date range 5-15-19 to 6-26-19. The data is missing weekends and holidays since the market is not open on those days.

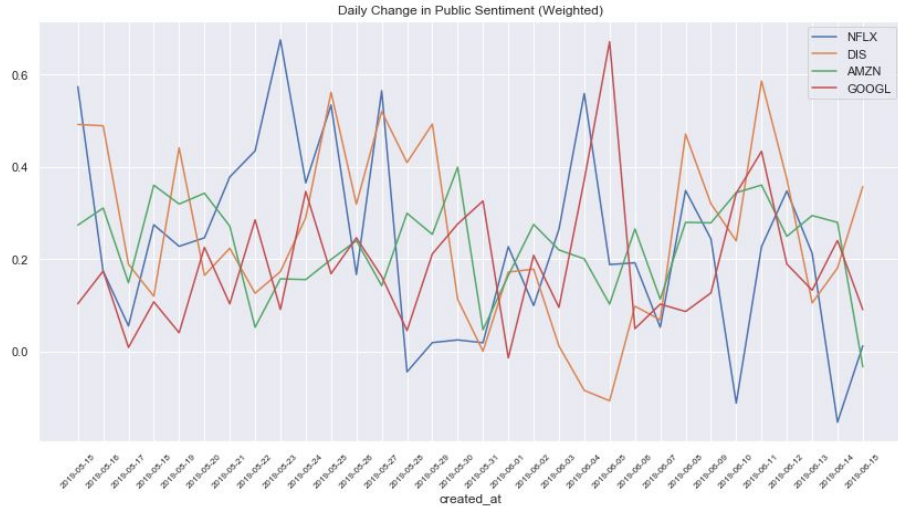
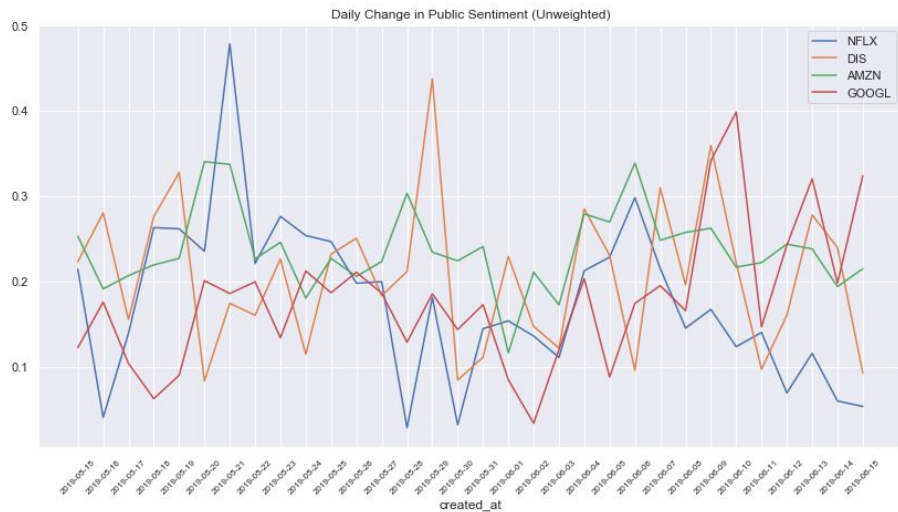
The information for the weekends and holidays were filled in using forward-filling linear interpolation.

Stock Dataframe
date
open
high
low
close
volume
company

Exploring Twitter Data

Company	Total Tweets	Avg Daily Tweets	Average Compound Sent
Netflix	64,055	2,002	0.18
Disney	78,404	2,450	0.21
Amazon	79,707	2,491	0.24
Google	89,728	2,804	0.18





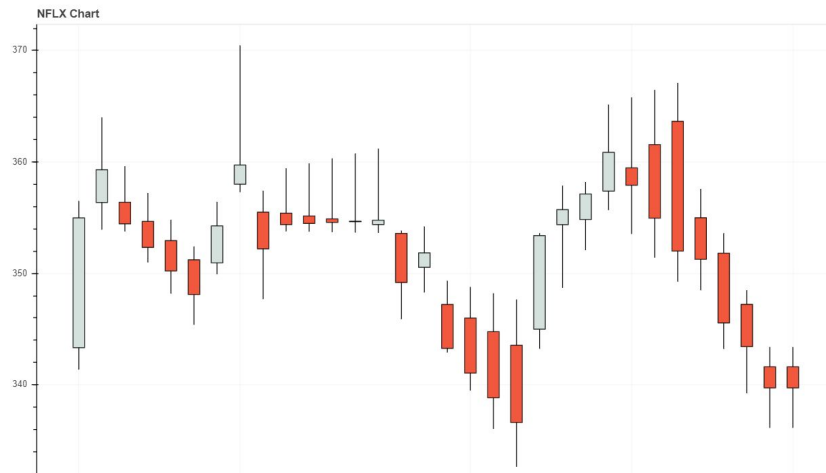
Daily Change in Public Sentiment

Twitter accounts that have more followers have more influence on the community since their tweet reaches more people. To incorporate this, we have given weight to each tweet according to proportion of followers they have compared to the total amount of followers for each day.

Exploring Stock Data

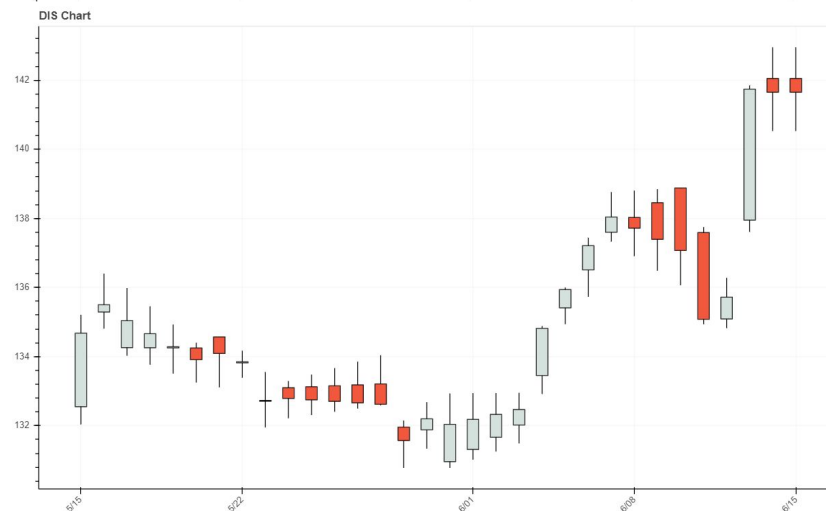
Netflix

- Average Close Price: \$350.97
- Average Daily Volume: 5,380,332



Disney

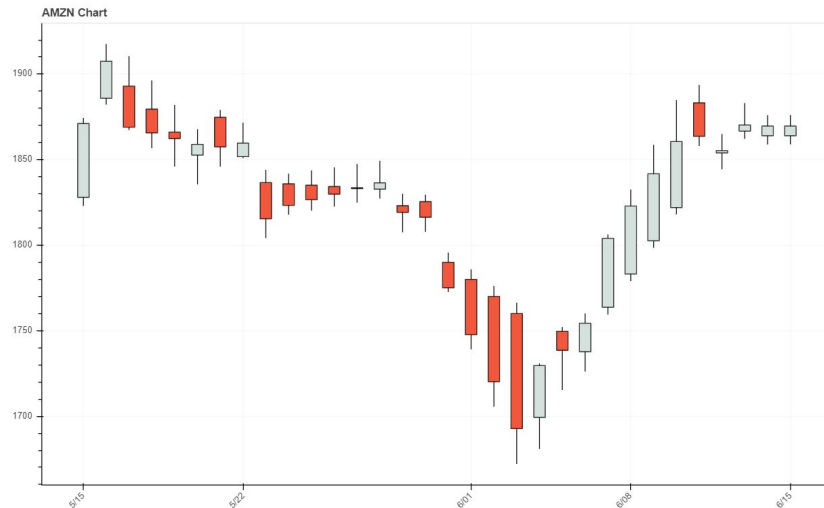
- Average Close Price: \$134.97
- Average Daily Volume: 7,934,529



Exploring Stock Data

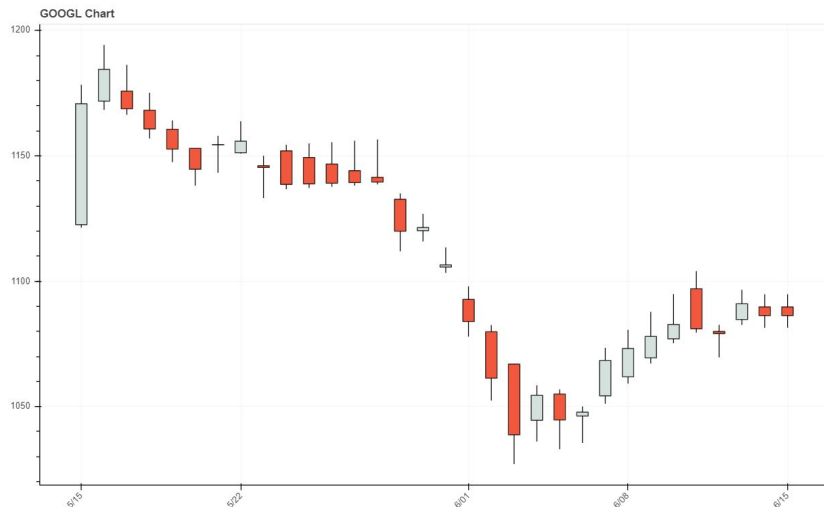
Amazon

- Average Close Price: \$1,824.00
- Average Daily Volume: 4,322,092

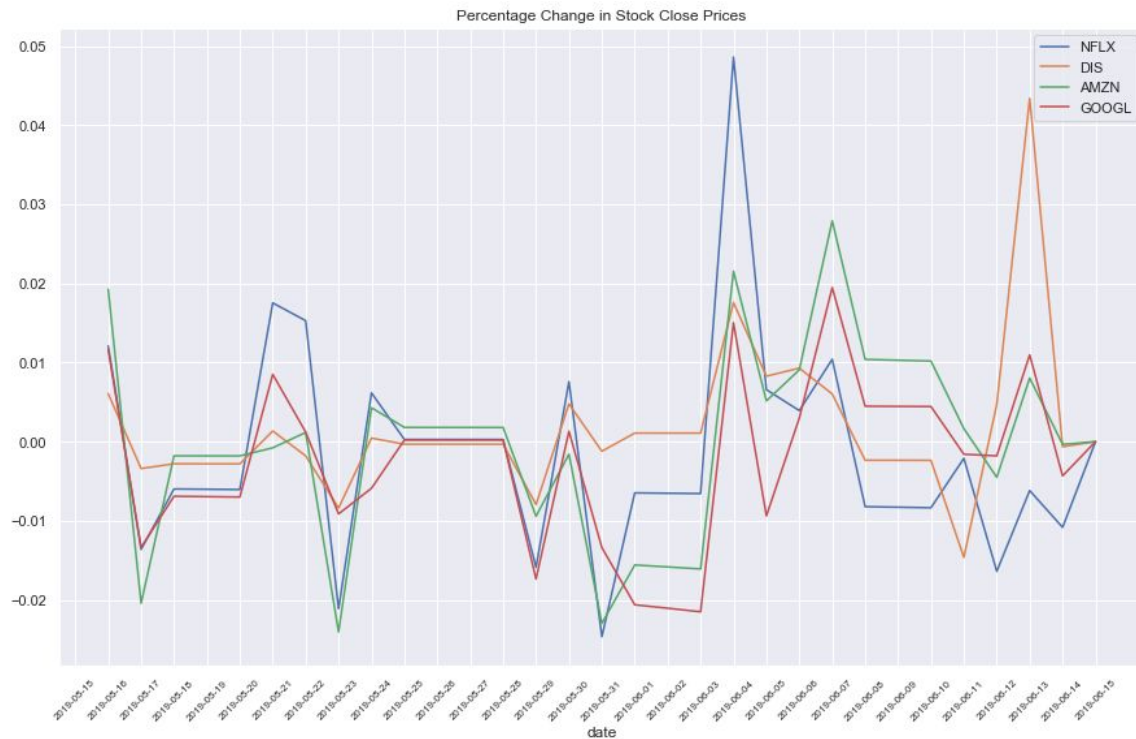


Google

- Average Close Price: \$1,110.56
- Average Daily Volume: 1,702,416



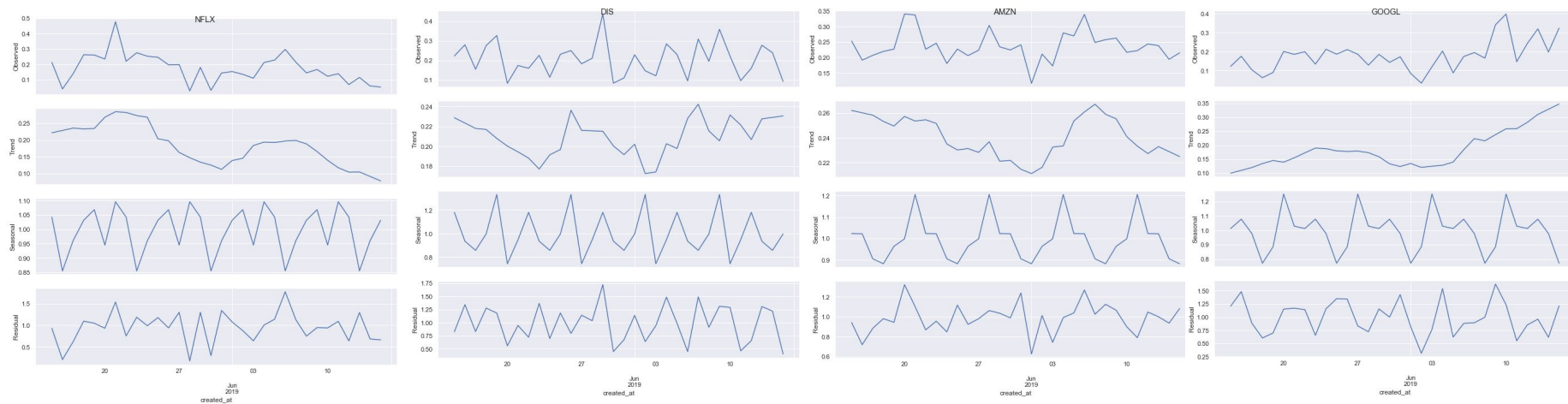
Daily Percentage Change in Stock Close Prices



Standardized Stock Price

Netflix and Disney are sold at way lower prices than Amazon and Google; therefore, harder to compare.
→ percentage change in price

Time Series Analysis: Seasonal Decomposition



Multiplicative Time Series: **Value = Base Level x Trend x Seasonality x Error**

After decomposing the seasonality from the time series, it is apparent that the public sentiment for Netflix has decreased over time. Google, on the other, has an upward trend in public sentiment. The other two have a more balanced trendline.

Time Series Analysis: Granger Causality

This tests whether the time series in the second column Granger causes the time series in the first column.

H₀: The public sentiment for the company does NOT Granger cause the movement in stock price for that company.

H₁: The public sentiment for the company Granger causes the movement in stock price for that company.

```
=====
No Causality

NFLX closing stock price and NFLX twitter sentiment showed NO causality, count: 31
DIS closing stock price and DIS twitter sentiment showed NO causality, count: 31
AMZN closing stock price and AMZN twitter sentiment showed NO causality, count: 31
-----

Causality

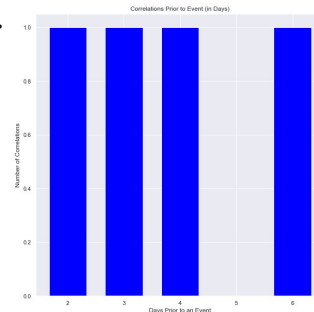
GOOGL closing stock price and GOOGL twitter sentiment showed causality, count: 31

~~~~~

Causality Count: 1
No Causality Count: 3

std 1.479019945774904
avg 3.75
Percent showing causality: 0.25
```

The only Google's public sentiment and market sentiment time series rejected the null. There was a correlation shown between the public sentiment and market sentiment for 2,3,4 and 6 days prior to the stock price movements.



Machine Learning: Feature Engineering

Twitter Data Set

Original Dataframe

created_at	tweet	follower_count	neg_sent	neu_sent	pos_sent	compound_sent	sentiment	Company
2019-05-15	RT @MileyCyrus: Black Mirror Out June 5th @net...	537	0.000	1.000	0.000	0.0000	neutral	NFLX
2019-05-15	RT @LaurenGerman: Let's all give a HUGE ROARIN...	2	0.000	0.733	0.267	0.7027	positive	NFLX
2019-05-15	@netflix This shit got me weak😓😓 talkin thru t...	259	0.535	0.465	0.000	-0.8481	negative	NFLX
2019-05-15	RT @MileyCyrus: Black Mirror Out June 5th @net...	762	0.000	1.000	0.000	0.0000	neutral	NFLX
2019-05-15	RT @LaurenGerman: Let's all give a HUGE ROARIN...	6	0.000	0.733	0.267	0.7027	positive	NFLX

Each row in the twitter data set represents one tweet (~ 2,500 tweets per day over the span of 31 days) → Aggregate data

- Weighted average compound sentiment score (weight = number of followers for the tweet / total number of followers for the day)
- Proportion of tweets classified as negative, positive, and neutral. To prevent multicollinearity, the proportion of tweets classified as neutral has been left out.

Resulting Dataframe



date	w_avg_sent	percent_neg	percent_pos
2019-05-15	0.573347	0.083511	0.454787
2019-05-16	0.173842	0.409941	0.400855
2019-05-17	0.055269	0.231081	0.442793
2019-05-18	0.274295	0.152960	0.565753
2019-05-19	0.227649	0.173720	0.584442

Machine Learning: Feature Engineering

Stock Data Set

Original Dataframe

date	open	high	low	close	volume	Company
2019-05-15	343.34	356.500000	341.390	354.990000	6.340118e+06	NFLX
2019-05-16	356.37	364.000000	353.935	359.310000	6.441463e+06	NFLX
2019-05-17	356.39	359.620000	353.785	354.450000	4.725448e+06	NFLX
2019-05-18	354.67	357.219333	350.990	352.336667	4.690804e+06	NFLX
2019-05-19	352.95	354.818667	348.195	350.223333	4.656159e+06	NFLX

To perform time series forecasting, time lags must be introduced into the model. The stock data from the previous day will be used to predict the stock growth of the current day.

Fundamental Analysis

- Percentage change in S&P index price

Standardizing Variables

- Volume
- Percentage change in stock growth instead of close price

Resulting Dataframe



date	volume_l1	stock_growth	stock_growth_l1	sp_growth
2019-05-15	NaN	NaN	NaN	NaN
2019-05-16	0.791431	0.012169	NaN	NaN
2019-05-17	0.884759	-0.013526	0.012169	0.008895
2019-05-18	-0.695522	-0.005962	-0.013526	-0.005837
2019-05-19	-0.727426	-0.005998	-0.005962	-0.002250

Machine Learning: Building Predictive Model

Final Dataset: Merging Twitter Data and Stock Data

Target Variable	Predictor Variables
<ul style="list-style-type: none">Stock Percentage Change	<ul style="list-style-type: none">Stock Percentage Change (1 day lag)Standardized Volume (1 day lag)S&P 500 Percentage Change (1 day lag)Weighted Compound Score (2 day lag)Proportion of Negative Tweet (2 day lag)Proportion of Positive Tweet (2 day lag)

Training Data:

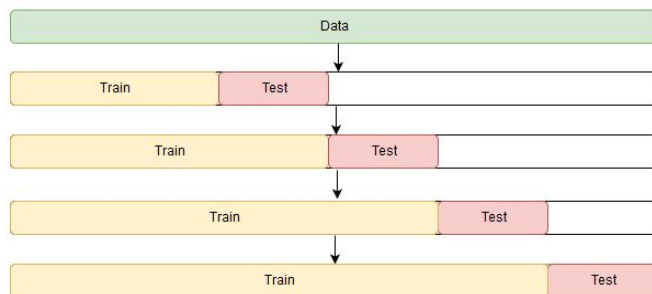
05/15/2019 - 06/15/2019 → 30 rows

Testing Data:

06/16/2019 - 06/26/2019 → 10 rows

Algorithm: Regression

Cross Validation: TimeSeriesSplit Method (5-fold)



Machine Learning: Building Predictive Model

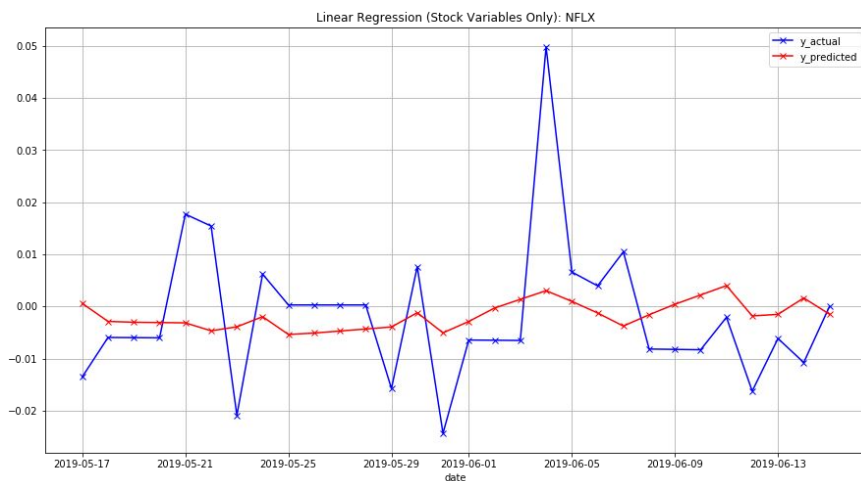
Linear Regression Stock Variables Only (finds the best fit line with the smallest prediction error throughout)

Predictor Variables: Stock Growth (1 day lag), Standardized Volume (1 day lag), S&P 500 Percentage Change (1 day lag)

Target Variable: Stock Growth

NFLX

- RMSE: 0.258
- MAPE: 54.94 %
- Direction: 44 %



Machine Learning: Building Predictive Model

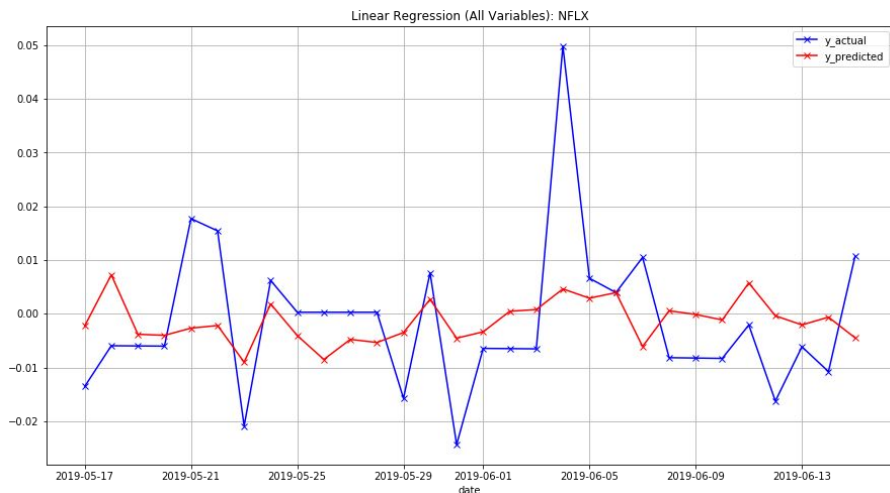
Linear Regression All Variables

Predictor Variables: Stock Percentage Change (1 day lag), Standardized Volume (1 day lag), S&P 500 Percentage Change (1 day lag), Weighted Compound Score (2 day lag), Proportion of Negative Tweet (2 day lag), Proportion of Positive Tweet (2 day lag)

Target Variable: Stock Growth

NFLX

- RMSE:0.040
- MAPE: 32.93 %
- Direction: 32.0 %



Machine Learning: Building Predictive Model

Performance on Test Data & Coefficients Interpretation

Growth in S&P 500 for the previous day seemed to have the highest positive effect while stock growth for the previous day seemed to have the most negative effect on stock growth. All the sentiment related variables have a positive effect on stock growth. Interestingly, the percentage of tweets classified as negative had the highest positive effect out of the 3 sentiment variables.

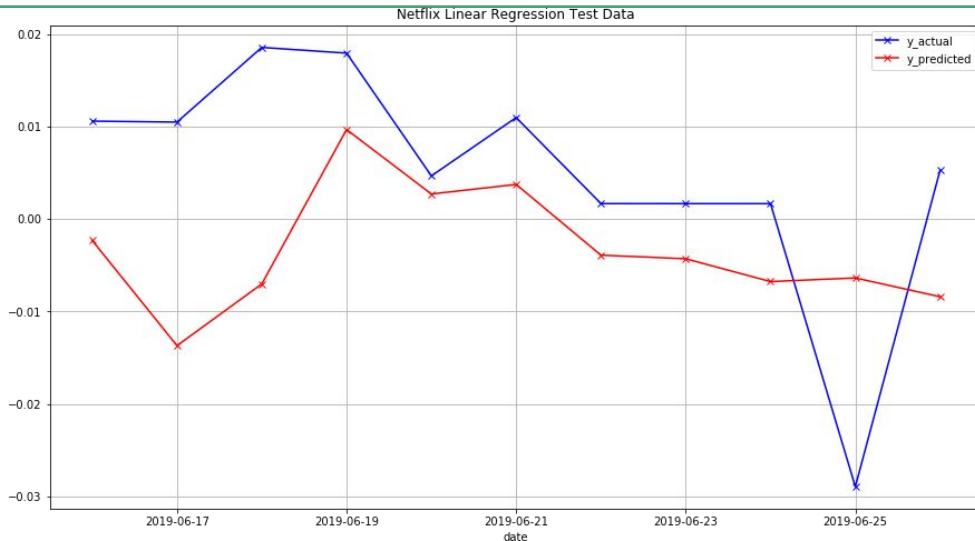
NFLX

Metrics

- RMSE: 0.014
- MAPE: 1.096 %
- Direction: 18.18 %

Coefficients

- sp_growth_l1: 0.716
- percent_neg: 0.063
- percent_pos: 0.025
- w_avg_sent: 0.008
- volume_l1: 0.003
- stock_growth_l1: -0.305

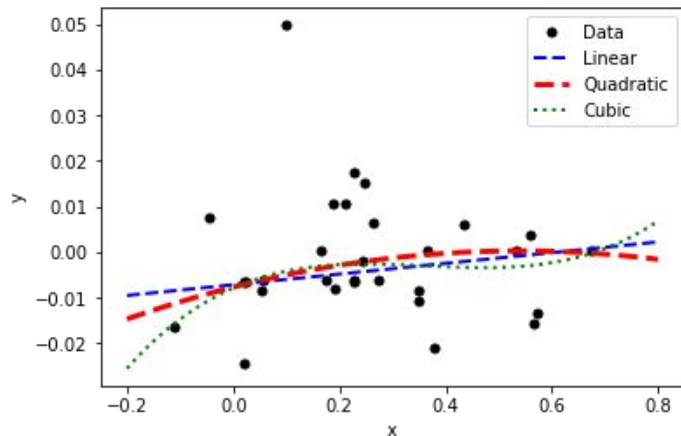
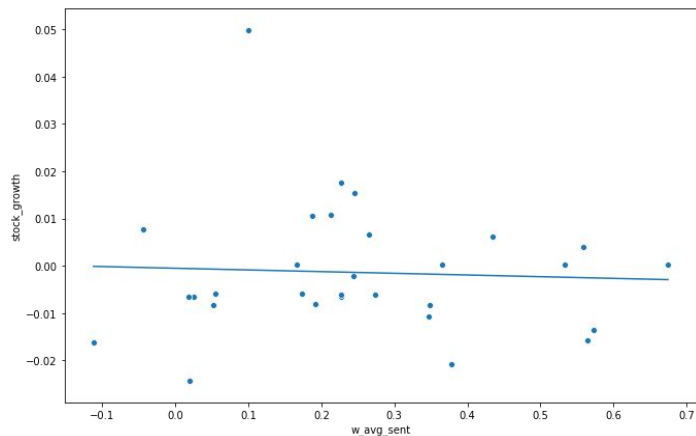


Machine Learning: Building Predictive Model

Polynomial Regression

Linear regression captures the patterns in the data better when the relation between the dependent variable and the independent variable is linear.

Graph of Weighted Average Sentiment vs Stock Growth



However, the relationship between the target variable and predictor variables is not linear. There is more variance towards the lower values in sentiment value than higher values. The cubic function seems to capture the relationship better.

Machine Learning: Building Predictive Model

Polynomial Linear Regression (Degree of 3) All Variables

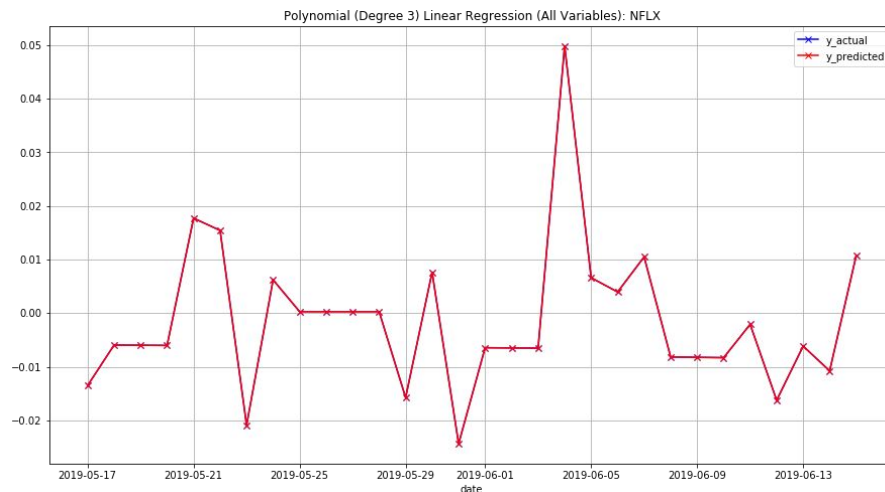
Predictor Variables: Stock Percentage Change (1 day lag), Standardized Volume (1 day lag), S&P 500 Percentage Change (1 day lag), Weighted Compound Score (2 day lag), Proportion of Negative Tweet (2 day lag), Proportion of Positive Tweet (2 day lag) + Polynomial of all variables + Interaction between all variables

Target Variable: Stock Growth

NFLX

- RMSE: 0.11
- MAPE: 79.8 %
- Direction: 52.0 %

OVERFITTING!!

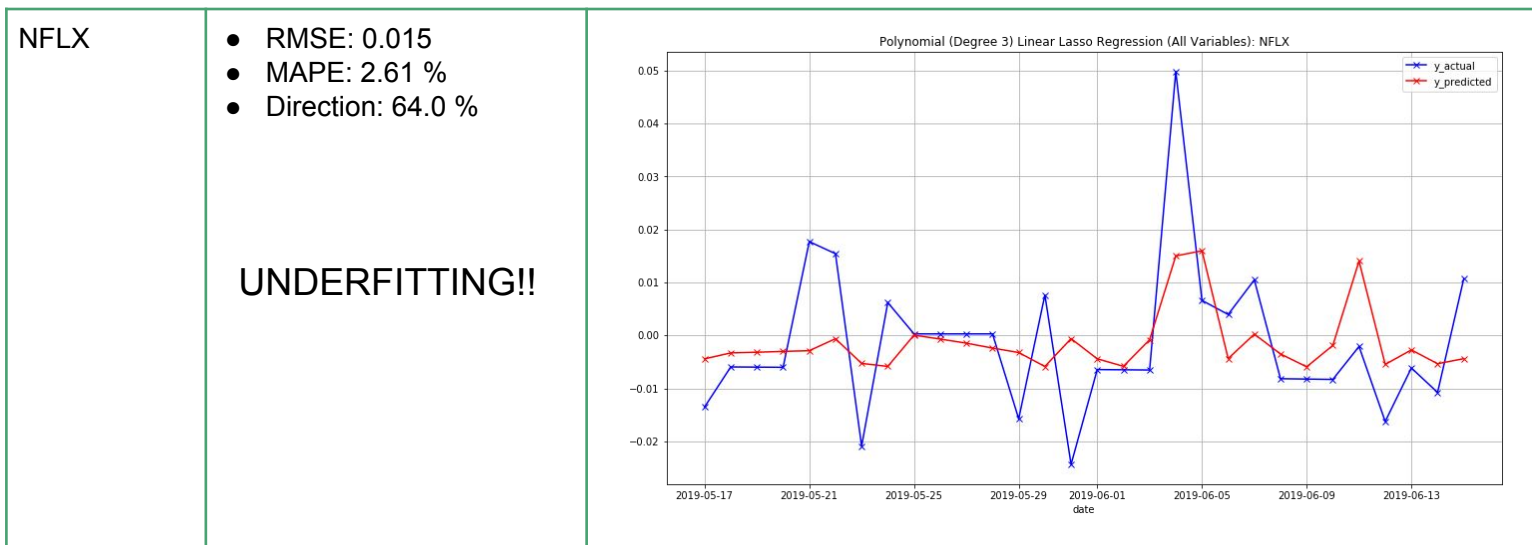


Machine Learning: Building Predictive Model

Polynomial Lasso Regression (Degree of 3) All Variables (“least absolute shrinkage and selection operator”)

Lasso regression helps in reducing overfitting and **feature selection**.

Cost Function = sum of squared prediction error + absolute value of the magnitude of the coefficients



Machine Learning: Building Predictive Model

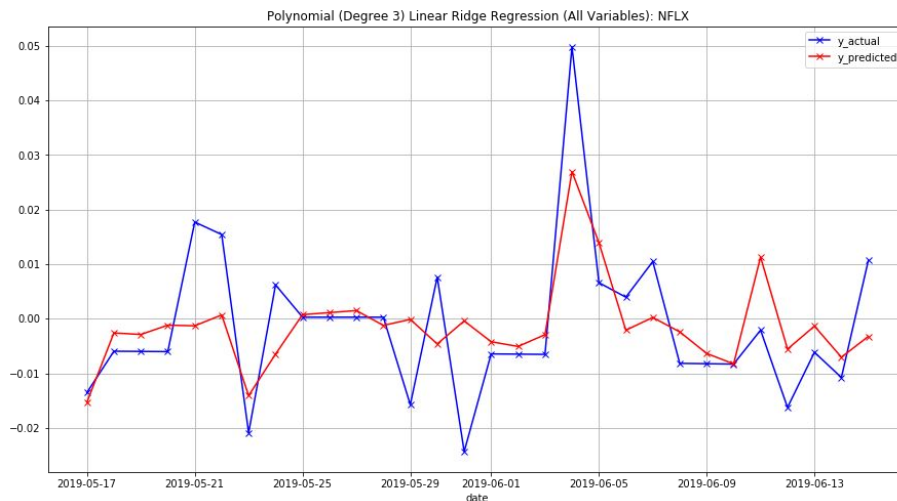
Polynomial Ridge Regression (Degree of 3) All Variables

Ridge regression only helps in reducing overfitting.

Cost Function = sum of squared prediction error + square of the magnitude of the coefficients

NFLX

- RMSE: 0.021
- MAPE: 5.5 %
- Direction: 60.0 %



Machine Learning: Building Predictive Model

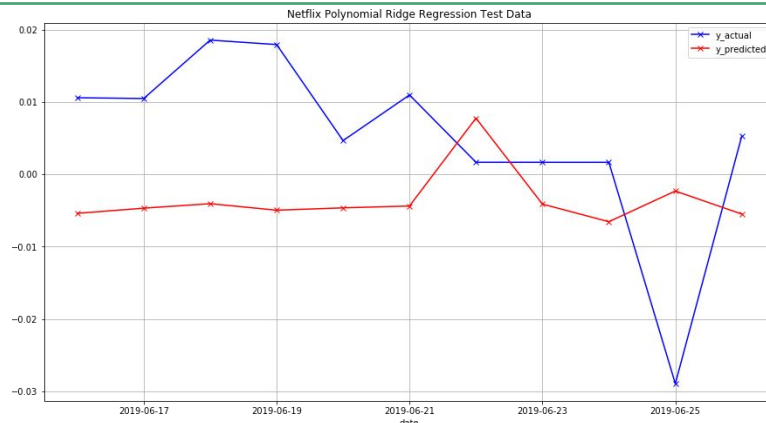
Lasso performed the best in terms of all 3 metrics, but the model does not capture large fluctuations well. → UNDERFITS MODEL

Ridge regression seems to be a good middle ground.

Performance on test data similar to linear regression without polynomial regression. The RMSE is only a little worse (0.001 difference) and the direction accuracy remained the same. Both models did not anticipate the drop in growth on 6/25/2019.

NFLX

- RMSE: 0.015
- MAPE: 2.19 %
- Direction: 18.18 %



Machine Learning: Building Predictive Model

Weighted Moving Average (Stock Growth and Compound Sentiment Value)

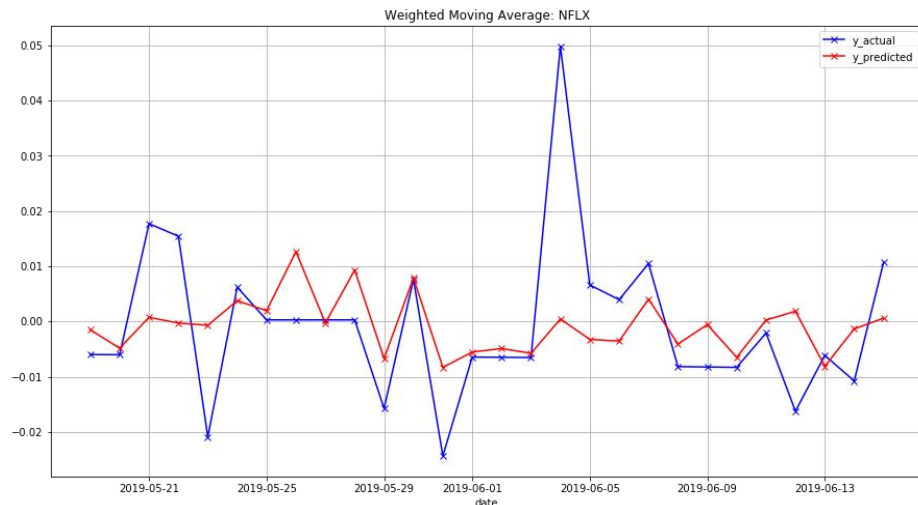
Method: Linear regression

Predictor Variables: stock growth for the past 3 days and sentiment value data for the past 3 days

Target Variable: Stock Growth

NFLX

- RMSE: 0.018
- MAPE: 9.91 %
- Direction: 50 %



Results

Company	Algorithm	RMSE	MAPE	Direction
AMZN	Polynomial Ridge Regression (Degree of 3)	0.012996	3.865511	36.0
AMZN	Linear Regression (Only Stock)	0.014983	3.236855	36.0
AMZN	Polynomial Lasso Regression (Degree of 3, Alpha...	0.016009	2.809316	32.0
AMZN	Linear Regression (All Variables)	0.016407	3.628019	36.0
AMZN	Weighted Moving Average	0.023974	6.787748	55.0
AMZN	Polynomial Regression (Degree of 3)	1.228785	66.054658	44.0
DIS	Polynomial Lasso Regression (Degree of 3, Alpha...	0.008717	2.145738	48.0
DIS	Polynomial Ridge Regression (Degree of 3)	0.009714	3.906355	48.0
DIS	Linear Regression (Only Stock)	0.013188	9.748320	44.0
DIS	Weighted Moving Average	0.013329	3.962521	25.0
DIS	Linear Regression (All Variables)	0.015153	11.118734	36.0
DIS	Polynomial Regression (Degree of 3)	0.236995	202.881103	40.0
GOOGL	Polynomial Ridge Regression (Degree of 3)	0.010487	4.297911	44.0
GOOGL	Polynomial Lasso Regression (Degree of 3, Alpha...	0.010577	4.680170	44.0
GOOGL	Weighted Moving Average	0.018437	15.827871	30.0
GOOGL	Linear Regression (All Variables)	0.023313	10.140069	44.0
GOOGL	Linear Regression (Only Stock)	0.025380	30.702864	28.0
GOOGL	Polynomial Regression (Degree of 3)	2.895230	164.115478	36.0
NFLX	Polynomial Lasso Regression (Degree of 3, Alpha...	0.015140	2.619799	64.0
NFLX	Weighted Moving Average	0.017684	9.912616	50.0
NFLX	Polynomial Ridge Regression (Degree of 3)	0.021401	5.495615	60.0
NFLX	Linear Regression (All Variables)	0.040265	32.931784	32.0
NFLX	Polynomial Regression (Degree of 3)	0.106495	79.800306	52.0
NFLX	Linear Regression (Only Stock)	0.258143	54.941599	44.0

Overall, it seems as if the polynomial regression with some sort of regularization modeled the data best.

Polynomial regression with lasso regularization seems to perform the best in all 3 metrics. BUT it underfits the model with it's extreme regularization and does not capture the fluctuations that investor would typically want to know. Investors do not care if the stock grows by 0.01%; they care more if the stock would grow by 5%. That's when an investor would make the most money.

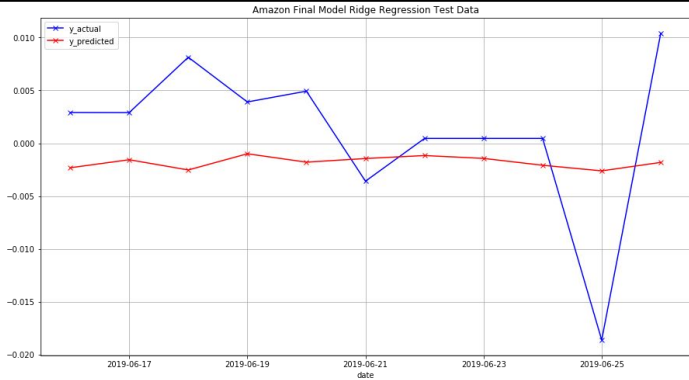
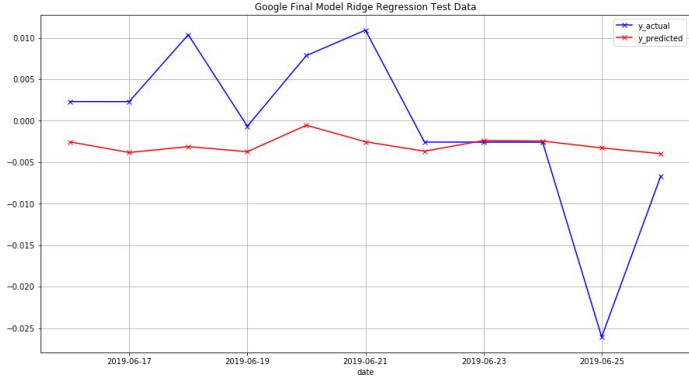
→ FINAL MODEL = polynomial regression with ridge regularization.

	RMSE	MAPE	Direction
Algorithm			
Polynomial Lasso Regression (Degree of 3, Alpha 0.001)	0.012611	3.063756	47.0
Polynomial Ridge Regression (Degree of 3)	0.013650	4.391348	47.0
Weighted Moving Average	0.018356	9.122689	40.0
Linear Regression (All Variables)	0.023785	14.454652	37.0
Linear Regression (Only Stock)	0.077923	24.657409	38.0
Polynomial Regression (Degree of 3)	1.116876	128.212886	43.0

Results

Company	Metrics	Visualization																																	
Netflix	<ul style="list-style-type: none">● RMSE: 0.014● MAPE: 0.844 %● Direction: 45.45 %	<p>Netflix Final Model Ridge Regression Test Data</p> <table border="1"><thead><tr><th>date</th><th>y_actual</th><th>y_predicted</th></tr></thead><tbody><tr><td>2019-06-17</td><td>0.010</td><td>-0.001</td></tr><tr><td>2019-06-18</td><td>0.010</td><td>-0.001</td></tr><tr><td>2019-06-19</td><td>0.018</td><td>-0.001</td></tr><tr><td>2019-06-20</td><td>0.018</td><td>-0.001</td></tr><tr><td>2019-06-21</td><td>0.005</td><td>-0.001</td></tr><tr><td>2019-06-22</td><td>0.011</td><td>-0.001</td></tr><tr><td>2019-06-23</td><td>0.002</td><td>0.002</td></tr><tr><td>2019-06-24</td><td>0.002</td><td>0.002</td></tr><tr><td>2019-06-25</td><td>-0.028</td><td>0.000</td></tr><tr><td>2019-06-26</td><td>0.005</td><td>0.003</td></tr></tbody></table>	date	y_actual	y_predicted	2019-06-17	0.010	-0.001	2019-06-18	0.010	-0.001	2019-06-19	0.018	-0.001	2019-06-20	0.018	-0.001	2019-06-21	0.005	-0.001	2019-06-22	0.011	-0.001	2019-06-23	0.002	0.002	2019-06-24	0.002	0.002	2019-06-25	-0.028	0.000	2019-06-26	0.005	0.003
date	y_actual	y_predicted																																	
2019-06-17	0.010	-0.001																																	
2019-06-18	0.010	-0.001																																	
2019-06-19	0.018	-0.001																																	
2019-06-20	0.018	-0.001																																	
2019-06-21	0.005	-0.001																																	
2019-06-22	0.011	-0.001																																	
2019-06-23	0.002	0.002																																	
2019-06-24	0.002	0.002																																	
2019-06-25	-0.028	0.000																																	
2019-06-26	0.005	0.003																																	
Disney	<ul style="list-style-type: none">● RMSE: 0.008● MAPE: 0.841 %● Direction: 45.45 %	<p>Disney Final Model Ridge Regression Test Data</p> <table border="1"><thead><tr><th>date</th><th>y_actual</th><th>y_predicted</th></tr></thead><tbody><tr><td>2019-06-17</td><td>-0.002</td><td>-0.002</td></tr><tr><td>2019-06-18</td><td>-0.002</td><td>0.001</td></tr><tr><td>2019-06-19</td><td>-0.012</td><td>0.000</td></tr><tr><td>2019-06-20</td><td>0.012</td><td>-0.003</td></tr><tr><td>2019-06-21</td><td>0.008</td><td>0.006</td></tr><tr><td>2019-06-22</td><td>-0.012</td><td>0.002</td></tr><tr><td>2019-06-23</td><td>-0.002</td><td>-0.004</td></tr><tr><td>2019-06-24</td><td>-0.002</td><td>-0.002</td></tr><tr><td>2019-06-25</td><td>0.005</td><td>-0.001</td></tr><tr><td>2019-06-26</td><td>0.003</td><td>-0.004</td></tr></tbody></table>	date	y_actual	y_predicted	2019-06-17	-0.002	-0.002	2019-06-18	-0.002	0.001	2019-06-19	-0.012	0.000	2019-06-20	0.012	-0.003	2019-06-21	0.008	0.006	2019-06-22	-0.012	0.002	2019-06-23	-0.002	-0.004	2019-06-24	-0.002	-0.002	2019-06-25	0.005	-0.001	2019-06-26	0.003	-0.004
date	y_actual	y_predicted																																	
2019-06-17	-0.002	-0.002																																	
2019-06-18	-0.002	0.001																																	
2019-06-19	-0.012	0.000																																	
2019-06-20	0.012	-0.003																																	
2019-06-21	0.008	0.006																																	
2019-06-22	-0.012	0.002																																	
2019-06-23	-0.002	-0.004																																	
2019-06-24	-0.002	-0.002																																	
2019-06-25	0.005	-0.001																																	
2019-06-26	0.003	-0.004																																	

Results

Company	Metrics	Visualization																																				
Amazon	<ul style="list-style-type: none">• RMSE: 0.008• MAPE: 2.207%• Direction: 18.18 %	 <p>Amazon Final Model Ridge Regression Test Data</p> <table border="1"><thead><tr><th>date</th><th>y_actual</th><th>y_predicted</th></tr></thead><tbody><tr><td>2019-06-17</td><td>0.003</td><td>-0.002</td></tr><tr><td>2019-06-18</td><td>0.003</td><td>-0.001</td></tr><tr><td>2019-06-19</td><td>0.008</td><td>-0.002</td></tr><tr><td>2019-06-20</td><td>0.004</td><td>-0.001</td></tr><tr><td>2019-06-21</td><td>0.005</td><td>-0.002</td></tr><tr><td>2019-06-22</td><td>-0.003</td><td>-0.001</td></tr><tr><td>2019-06-23</td><td>0.000</td><td>-0.001</td></tr><tr><td>2019-06-24</td><td>0.000</td><td>-0.001</td></tr><tr><td>2019-06-25</td><td>-0.018</td><td>-0.002</td></tr><tr><td>2019-06-26</td><td>0.010</td><td>-0.002</td></tr></tbody></table>	date	y_actual	y_predicted	2019-06-17	0.003	-0.002	2019-06-18	0.003	-0.001	2019-06-19	0.008	-0.002	2019-06-20	0.004	-0.001	2019-06-21	0.005	-0.002	2019-06-22	-0.003	-0.001	2019-06-23	0.000	-0.001	2019-06-24	0.000	-0.001	2019-06-25	-0.018	-0.002	2019-06-26	0.010	-0.002			
date	y_actual	y_predicted																																				
2019-06-17	0.003	-0.002																																				
2019-06-18	0.003	-0.001																																				
2019-06-19	0.008	-0.002																																				
2019-06-20	0.004	-0.001																																				
2019-06-21	0.005	-0.002																																				
2019-06-22	-0.003	-0.001																																				
2019-06-23	0.000	-0.001																																				
2019-06-24	0.000	-0.001																																				
2019-06-25	-0.018	-0.002																																				
2019-06-26	0.010	-0.002																																				
Google	<ul style="list-style-type: none">• RMSE: 0.0097• MAPE: 1.442 %• Direction: 54.54 %	 <p>Google Final Model Ridge Regression Test Data</p> <table border="1"><thead><tr><th>date</th><th>y_actual</th><th>y_predicted</th></tr></thead><tbody><tr><td>2019-06-17</td><td>0.002</td><td>-0.002</td></tr><tr><td>2019-06-18</td><td>0.002</td><td>-0.003</td></tr><tr><td>2019-06-19</td><td>0.010</td><td>-0.003</td></tr><tr><td>2019-06-20</td><td>-0.001</td><td>-0.003</td></tr><tr><td>2019-06-21</td><td>0.008</td><td>-0.001</td></tr><tr><td>2019-06-22</td><td>0.011</td><td>-0.002</td></tr><tr><td>2019-06-23</td><td>-0.002</td><td>-0.003</td></tr><tr><td>2019-06-24</td><td>-0.002</td><td>-0.002</td></tr><tr><td>2019-06-25</td><td>-0.002</td><td>-0.003</td></tr><tr><td>2019-06-26</td><td>-0.025</td><td>-0.003</td></tr><tr><td>2019-06-27</td><td>-0.012</td><td>-0.003</td></tr></tbody></table>	date	y_actual	y_predicted	2019-06-17	0.002	-0.002	2019-06-18	0.002	-0.003	2019-06-19	0.010	-0.003	2019-06-20	-0.001	-0.003	2019-06-21	0.008	-0.001	2019-06-22	0.011	-0.002	2019-06-23	-0.002	-0.003	2019-06-24	-0.002	-0.002	2019-06-25	-0.002	-0.003	2019-06-26	-0.025	-0.003	2019-06-27	-0.012	-0.003
date	y_actual	y_predicted																																				
2019-06-17	0.002	-0.002																																				
2019-06-18	0.002	-0.003																																				
2019-06-19	0.010	-0.003																																				
2019-06-20	-0.001	-0.003																																				
2019-06-21	0.008	-0.001																																				
2019-06-22	0.011	-0.002																																				
2019-06-23	-0.002	-0.003																																				
2019-06-24	-0.002	-0.002																																				
2019-06-25	-0.002	-0.003																																				
2019-06-26	-0.025	-0.003																																				
2019-06-27	-0.012	-0.003																																				

Conclusion

When predicting stock prices, adding variables related to public sentiment is essential since it explains some of the variance.

The model was able to predict values closer to the actual values and had better accuracy in predicting if the stock price was going to increase or decrease. The model was ultimately improved when more features were engineered via polynomial transformation, which added features up until cubic power and interaction terms between. 81 features were created from this process, causing overfitting when trained. Therefore, Ridge regularization was incorporated to suppress the coefficients of the features that were not as important.

Limitations

Although the RMSE scores are low, the models did not capture the drastic fluctuations as I hoped it would (e.g. drop on 6/25/2019). For future work, the following limitations should be considered:

- More training data
- Better sentiment analyzer (VADER misclassified neutral tweets)
- More predictor variables (Unexplained variance in stock growth)
- Explore more regression algorithms

Coefficients of the weighted average model illustrated that there was a need for more training data.

Overtime, coefficients should converge toward 0.

Coefficients	Netflix	Disney	Amazon	Google
Stock_growth_I1	-0.028	0.1464	0.1915	0.2215
Stock_growth_I2	-0.026	-0.1823	0.2210	0.1260
Stock_growth_I3	-0.058	0.0590	0.1622	0.0819
W_avg_sent_I1	0.007	-0.0033	0.0009	-0.0037
W_avg_sent_I2	-0.006	0.00823	-0.0145	0.0204
W_avg_sent_I3	0.026	-0.0111	-0.0339	-0.0081

References

1. <https://arxiv.org/pdf/1010.3003.pdf>
2. <https://machinelearningmastery.com/decompose-time-series-data-trend-seasonality/>
3. <https://blog.projectpiglet.com/2018/01/causality-in-cryptomarkets/>
4. <https://medium.com/cindicator/backtesting-time-series-models-weekend-of-a-data-scientist-92079cc2c540>
5. <https://towardsdatascience.com/ridge-and-lasso-regression-a-complete-guide-with-python-scikit-learn-e20e34bcbf0b>
6. <https://towardsdatascience.com/machine-learning-with-python-easy-and-robust-method-to-fit-non-linear-data-19e8a1ddbd49>