

# MSc in Business Analytics

Academic Year 2019-2020



**Course:** Machine Learning and Content Analytics

**Instructor:** Papageorgiou Haris

**Subject:** Predicting Stock Market from Financial Online News with LSTM

## Description

News media has a tremendous effect on the stock market prices. Financial news articles are read every day by millions, interpreted by analysts, and analyzed by automated stock algorithms at quantitative trading firms. Key news and data about company trades, revenues, market climates, and geopolitical conflicts strongly influence trading and market transactions. In our project, we address the problem of predicting stock price trends using news headlines. We want to determine the extent to which we can predict the movement of stock prices using features extracted from online news.

Related Work Stock market prediction is an area of very active research. Attempts have been made to predict stock prices using technical information such as company technical and price history. Recent work has begun to incorporate natural language text data into prediction – for example, using full-text of news articles, hidden Markov model incorporating news headlines and sentiment in social media. Most of these works have worked with short term price movements (on the order of minutes), while we will focus on end of day stock prices; we find that a longer time horizon is less susceptible to noise, and more beneficial to the common investor.

## Mission

As financial institutions begin to embrace artificial intelligence, machine learning is increasingly utilized to help make trading decisions. Although there is an abundance of stock data for machine learning models to train on, a high noise to signal ratio and the multitude of factors that affect stock prices are among the several reasons that predicting the market is difficult. At the same time, these models don't need to reach high levels of accuracy because even 60% accuracy can deliver solid returns. One method for predicting stock prices is using a long short-term memory neural network (LSTM) for times series forecasting. An approach with the accuracy described above would be a very helpful toolset on hands of investors and especially in trading Companies. Knowing the next day's stock price, the investor/trader would be able to know where to put his/her money. Gradually with an accuracy more than 50% this business workflow would be profitable.

For solving the above business scenario we suggest the following methodology: We have used Neural Networks for predicting Apple and Microsoft Stock Prices, feeding our LSTM Network both with the sentiment of Financial News referring the stock and the previous days actual Open price. We suggest the following models to combat the above project:

- LSTM with Historical Stock Prices and TextBlob Sentiment Score Analysis
- LSTM with Finbert Sentiment Analysis
- LSTM with Finbert Sentence Embeddings
- Finally, we compare the above implementation with ARIMA and make some predictions

## Data Sources

Our data came from 2 data sources. The online news dataset containing titles and headlines about stock news is located publicly on Kaggle(<https://www.kaggle.com/gennadiyr/us-equities-news-data>). The stock prices of the two companies we have analyzed (Apple and Microsoft) came from Yahoo Finance public API. Our Datasets of stock price seems to be imbalanced since stock prices of both companies are growing gradually over the last years.

## Exploratory Data Analysis

	Open	High	Low	Close	Adj Close	Volume
count	2428.000000	2428.000000	2428.000000	2428.000000	2428.000000	2.428000e+03
mean	29.069151	29.328832	28.809544	29.079849	27.325662	2.806688e+08
std	14.768329	14.903127	14.661120	14.797974	15.133855	2.194149e+08
min	6.478929	6.612857	6.453571	6.459286	5.581721	4.544800e+07
25%	18.008660	18.137589	17.842768	17.952144	15.900101	1.194640e+08
50%	26.181249	26.445000	25.988751	26.273750	24.317713	2.046454e+08
75%	38.966249	39.278126	38.585001	38.961249	37.482624	3.833529e+08
max	81.112503	81.962502	80.837502	81.800003	81.432350	1.880998e+09

Figure 1: stats of Apple features

	Open	High	Low	Close	Adj Close	Volume
<b>count</b>	2769.000000	2769.000000	2769.000000	2769.000000	2769.000000	2.769000e+03
<b>mean</b>	63.963868	64.561322	63.333496	63.983063	59.910382	4.037384e+07
<b>std</b>	45.953716	46.474798	45.373649	45.957649	47.325218	2.307689e+07
<b>min</b>	23.090000	23.320000	22.730000	23.010000	18.133703	7.425600e+06
<b>25%</b>	29.660000	29.889999	29.420000	29.709999	24.187122	2.483380e+07
<b>50%</b>	45.930000	46.340000	45.570000	45.939999	41.130180	3.473280e+07
<b>75%</b>	85.430000	86.000000	84.879997	85.500000	82.166878	5.035250e+07
<b>max</b>	229.270004	232.860001	227.350006	231.649994	231.649994	3.193179e+08

Figure 2: stats of Microsoft features

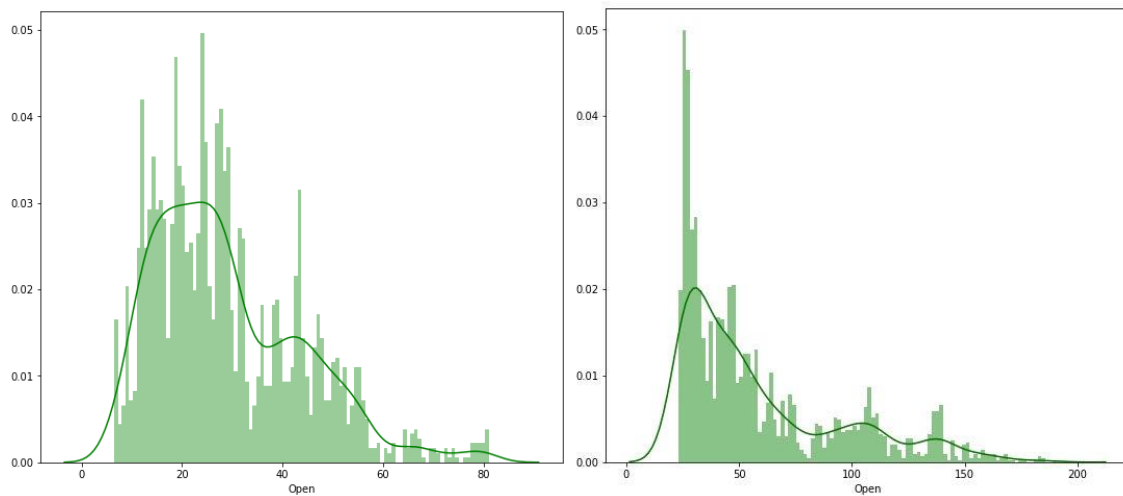


Figure 3,4: distribution of Open Price (left: Apple, right: Microsoft)

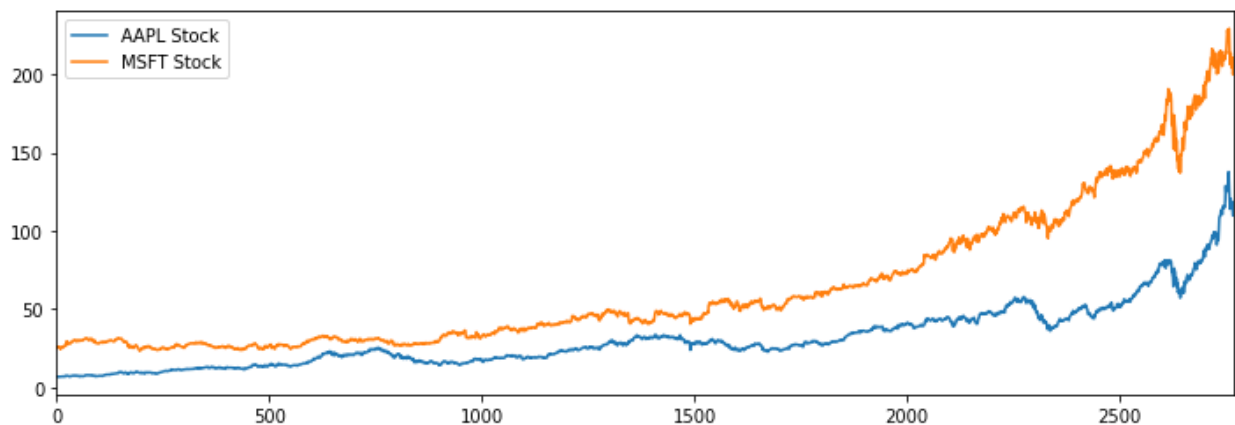


Figure 5: timeseries of both stock Prices

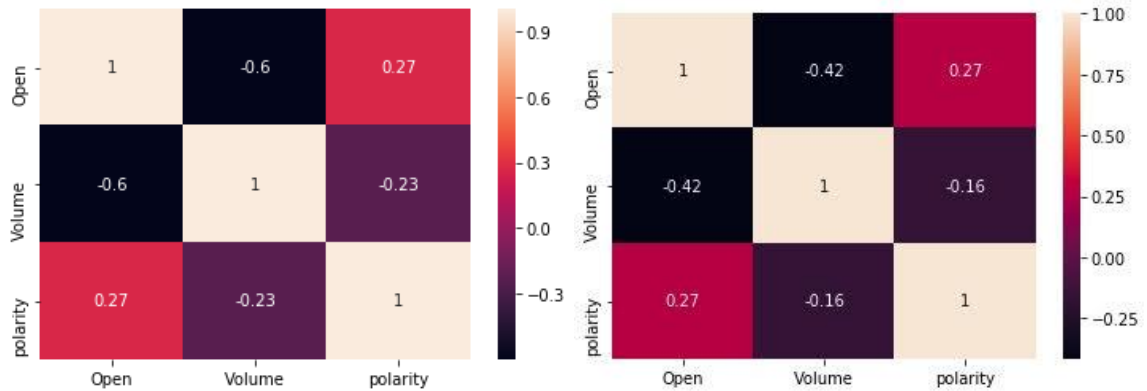


Figure 5,6: Correlation between Open Price and sentiment Polarity(left: Apple, right: Microsoft)

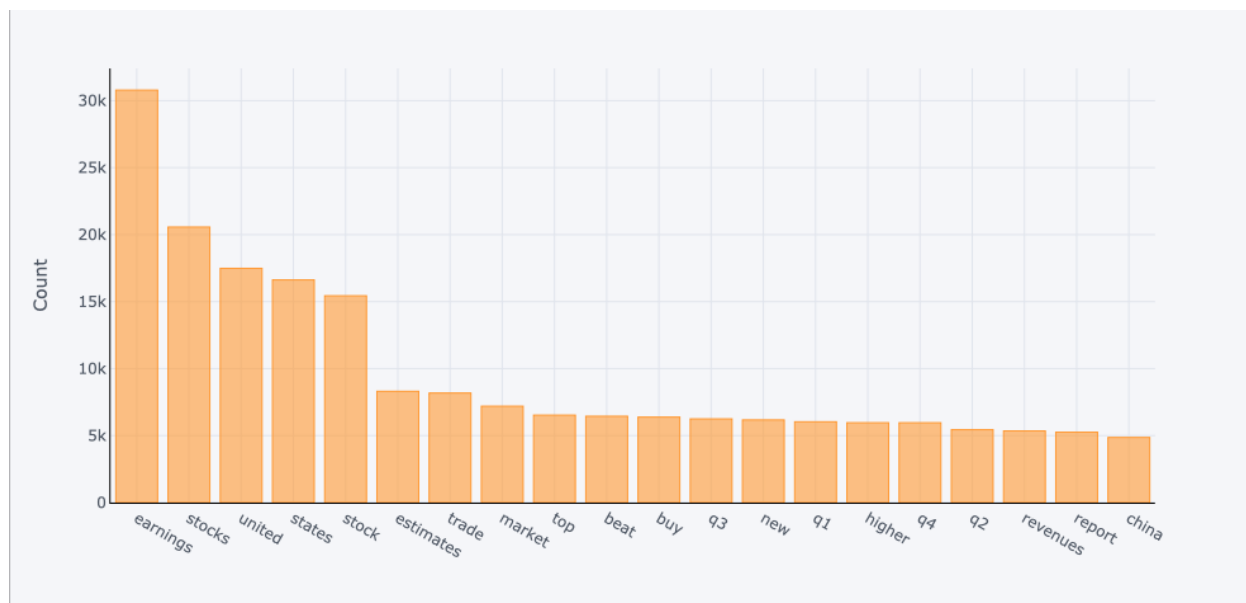


Figure 6: Top 20 Most common words in news titles.

### Data Normalization

Normalization is changing the values of numeric columns in the dataset to a common scale, which helps the performance of our model. To scale the training dataset, we use Scikit-Learn's MinMaxScaler with numbers between zero and one.

## Implementation 1: LSTM with Historical Stock Prices and TextBloB Sentiment Score Analysis.

For this particular implementation we have used both previous years stock prices (all the metrics) and sentiment analysis polarity score comes from the Library Textblob. In order the data to have the same shape we created a dataset contains dates which include both stock prices values and sentiment extracted from Titles News. For each date we join titles of new together, we apply cleaning on text and keeping only english stopwords and finally we apply on them Textblob Sentiment Analysis Library to the titles. The output is two new columns on final dataframe , named “polarity” having the score and “subjectivity”.

release_date	title	Date	Open	High	Low	Close	Adj Close	Volume	polarity	subjectivity
2009-09-21	UPDATE 4 China buys Noble stake seeks commodit...	2009-09-21	6.581786	6.612857	6.486429	6.572143	5.679245	437715600	0.056597	0.380208
2009-09-22	US STOCKS Wall St open commods rise weak dolla...	2009-09-22	6.613929	6.620714	6.530357	6.588572	5.693443	356753600	0.018947	0.319649
2009-09-23	Nikkei climbs U S economy optimism Fed,Nikkei ...	2009-09-23	6.621428	6.746428	6.608214	6.625000	5.724923	593563600	0.042859	0.348008
2009-09-24	Nikkei gains 1 2 pct holiday JAL Aiful sink,Ni...	2009-09-24	6.685714	6.703571	6.527500	6.565000	5.673074	550880400	0.075398	0.305545
2009-09-25	US STOCKS Wall St drops recovery worry RIM res...	2009-09-25	6.500357	6.625000	6.480000	6.513214	5.628324	445239200	0.131481	0.248148

The data were normalized to be in same range and then were given to LSTM input. The feature columns that we keep are ['Open','High','Low','Close','Adj Close','Volume','polarity'] and we want to predict the next days Open price. To achieve that we use a time window of 60 days that LSTM look backs and tries to predict the next day's Open Price.

### Incorporating Timesteps Window into Data

We should input our data in the form of a 3D array to the LSTM model. First, we create data in 60 timesteps before using Numpy to convert it into an array. Finally, we convert the data into a 3D array with X\_train samples, 60 timestamps, and one feature at each step. We are trying to predict the next day price, so the are look up step is 1 day. For this reason we have shifted the open prices to 1 day up in front , assuming that previous days news influence next day's open price.

### Methodology

Our LSTM has been built having 3 layers, can be Bidirectional and we have tried 128 and 256 units/neurons to explore its capabilities in terms of performance. We have used 2 different drop outs values (0.1 ,0.4) and two different batch sizes [32,64] to get the best parameters for our architecture. After each layer, we are adding dropout and for activation function we are using relu. Since Adam optimizer is known as the best optimizer in terms of performance we are keeping adam as default. We are training our model in 100 epochs for each implementation and finally get these parameters that minimize the Mean Absolute Error.

## Learning Curves

To confirm that loss is reducing during epochs we are plotting using tensorboard the epoch and validation loss. Along with Tensorboard callback we are also using Reduce learning rate when the metric has stopped improving.

### Apple

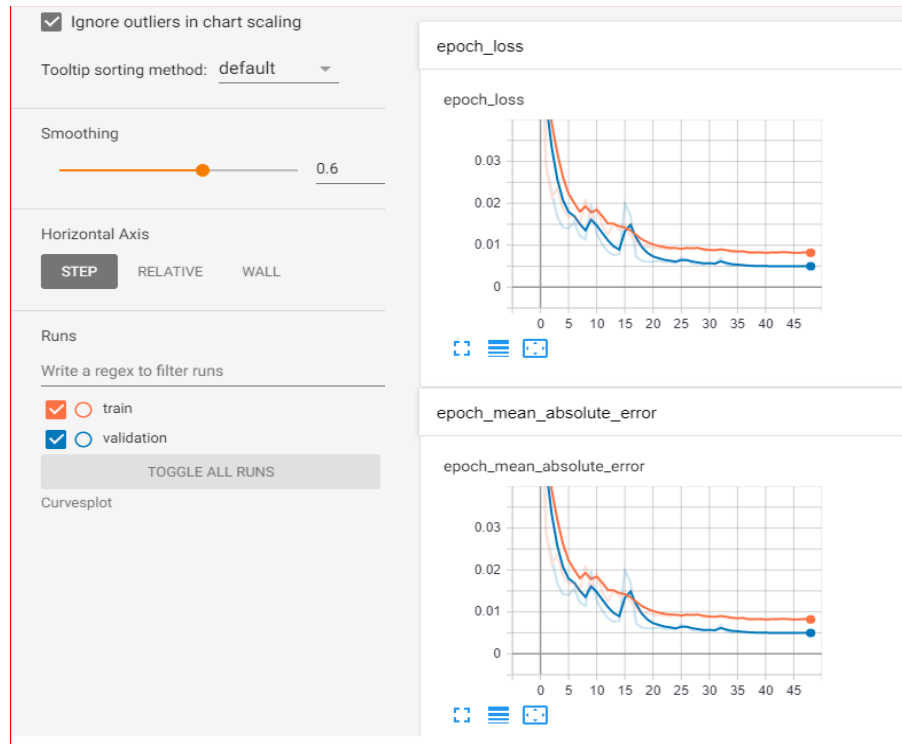


Figure 7: Learning Curves for Stocks &Sentiment model

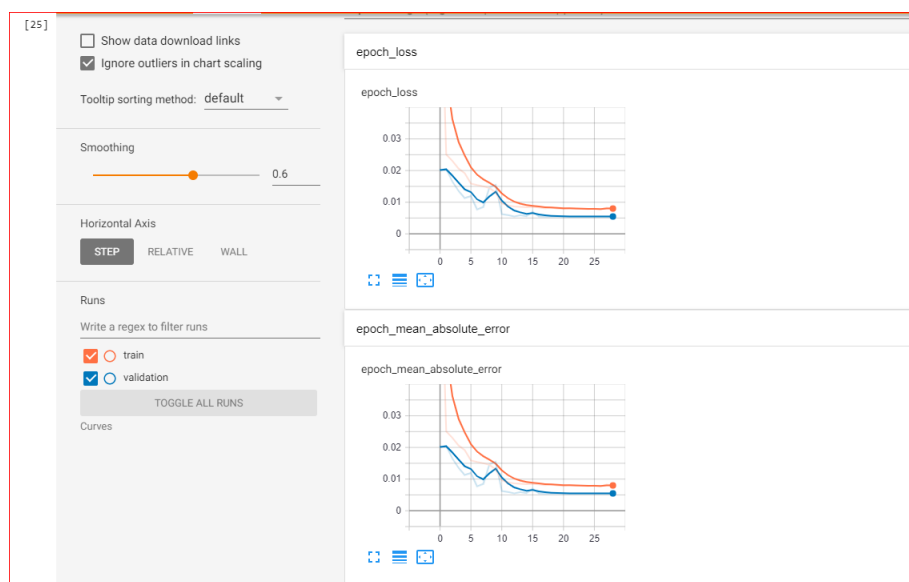


Figure 8: Learning Curves for Stocks only model

## Microsoft

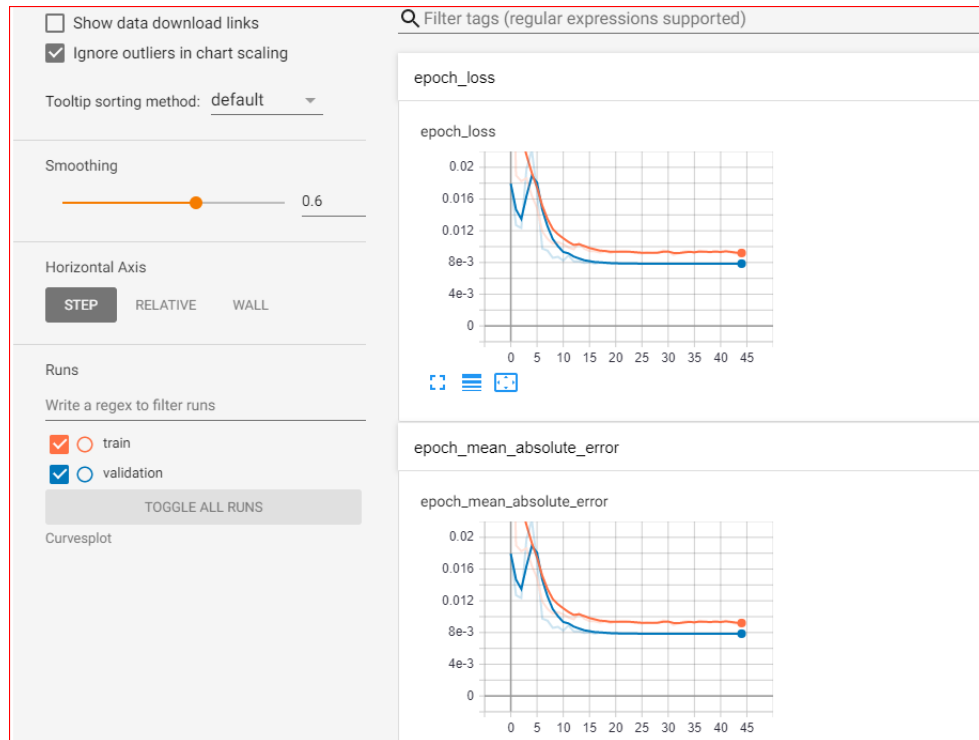


Figure 9: Learning Curves for Stocks & Sentiment model

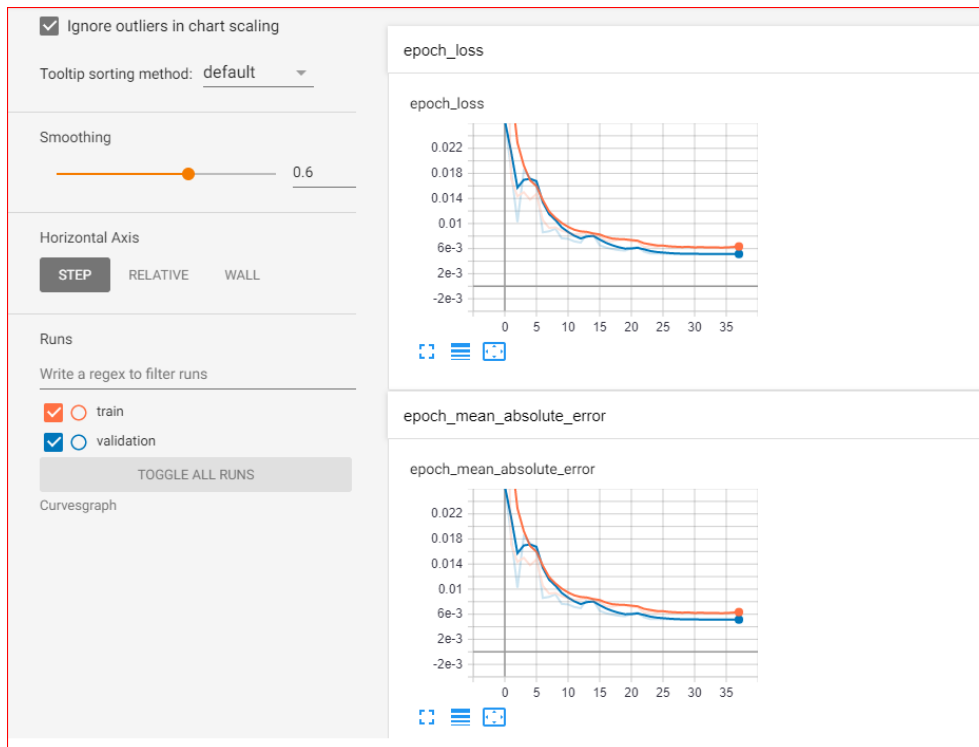


Figure 10 : Learning Curves for Stocks Only model



## Results

In order to interpret the performance of the below implementation we focus on the Mean Absolute Error. The absolute error is the absolute value of the difference between the forecasted value and the actual value. MAE tells us how big of an error we can expect from the forecast on average. One of the main problem was to join two data sources (news and stocks) in a common dataframe having dates with information for both of sources. Cleaning and preparing news title from stopwords, keeping only English language words was also a challenge to have a clearer polarity output. Each of the model needed around 40 mins to be trained and tuned working with Google Colab and using GPU runtime.

batch_size	dropout	units	mae	batch_size	dropout	units	mae
32.0	0.1	128.0	6.987520	32.0	0.1	128.0	7.164780
32.0	0.1	256.0	6.863086	32.0	0.1	256.0	7.098525
32.0	0.4	128.0	7.131041	32.0	0.4	128.0	7.381705
32.0	0.4	256.0	7.088577	32.0	0.4	256.0	7.401638
64.0	0.1	128.0	6.986449	64.0	0.1	128.0	7.409796
64.0	0.1	256.0	6.899435	64.0	0.1	256.0	7.272293
64.0	0.4	128.0	7.244094	64.0	0.4	128.0	7.471470
64.0	0.4	256.0	7.357203	64.0	0.4	256.0	7.352279

Figure 11: Apple only with stocks - Apple with sentiment and stock prices

batch_size	dropout	units	mae	batch_size	dropout	units	mae
32.0	0.1	128.0	24.053155	32.0	0.1	128.0	24.474606
32.0	0.1	256.0	23.796392	32.0	0.1	256.0	24.688490
32.0	0.4	128.0	24.743522	32.0	0.4	128.0	24.954283
32.0	0.4	256.0	24.629776	32.0	0.4	256.0	24.854168
64.0	0.1	128.0	24.503716	64.0	0.1	128.0	24.642374
64.0	0.1	256.0	23.957179	64.0	0.1	256.0	60.699660
64.0	0.4	128.0	24.900420	64.0	0.4	128.0	25.049769
64.0	0.4	256.0	24.477669	64.0	0.4	256.0	24.788361

Figure 12: Microsoft only with stocks - Microsoft with sentiment and stock prices

Choosing the hyperparameters that minimizing the Mean Absolute Error we create our final model. (batch size, dropout,units). Setting these parameters to our final model we end up with a model ready to predict the next day's Open Price. Our last day of having stock Open Price for both companies is 13-02-2020. So, we wanted to predict Open Stock Price for 14-02-2020. The table below presents the final results.

<b>Apple</b> Actual Open Price 14/2/2020	Stock and Sentiment Model	Stock Only Model
81.18	80.63	81.03
<b>Microsoft</b> Actual Open Price 14/2/2020	Stock and Sentiment Model	Stock Only Model
183.25	176.45	177.78

Table 1: Final Results of 4 models

## Implementation 2: LSTM with Finbert Sentence Embeddings

### Data preprocessing

The approach that will be used will try to predict if a stock will rise or fall based on the previous day's news which makes it a typical classification problem.

The first thing that needed to be done was to filter the dataset for the specific companies and then clean the dataset, by removing stopwords or pointless blanks. At this point the news were also grouped based on their release date. This was needed in order to have 1 row for each date and a list of all the headlines in that same row.

After this is done, we examine the days of each date in order to shift Saturdays and Sundays to the date of Friday. This workaround is needed because when we try to predict how the opening price of Monday will behave, we need all the news that were released within the weekend. If we shift the dates of weekends to Friday, then we will always be able to use the previous date's news to predict the next.

Later the stock prices of Apple and Microsoft for the relevant dates were downloaded from Yahoo Finance. In order to integrate these 2 datasets only the dates that were common between the headlines and the stock price datasets were kept.

Then the opening price difference with the previous day's opening price was calculated and eventually the negative difference was replaced with 0 and a positive difference which meant that the stock price went up, was replaced with 1.

### MSFT

	Date	headlines	Open
1	2009-10-01	[US STOCKS SNAPSHOT Nasdaq falls 2 pct selloff...	0
2	2009-10-22	[US STOCKS Futures mixed ahead data earnings, ...	1
3	2009-10-23	[US STOCKS SNAPSHOT Futures rise Microsoft res...	0
4	2009-10-26	[US STOCKS Wall St gains tech energy shares, U...	1
5	2009-10-28	[US STOCKS Wall St falls data lowered Goldman ...	0
6	2011-03-30	[Microsoft co founder Allen blasts Gates book]	0
7	2011-03-31	[US STOCKS Futures suggest slightly lower open...	0
8	2011-04-04	[Google could target FTC antitrust probe report]	1
9	2011-04-05	[US STOCKS S P hovers near resistance volume s...	1
10	2011-04-06	[US STOCKS S P500 index edges higher light vol...	1

## APPLE

	Date	headlines	Open
1	2012-07-19	[Mid Year Update U S And Canadian Stock Market...	1
2	2012-07-23	[Summer Heat Scorches Europe And U S, Apple Ea...	0
3	2012-07-24	[Market Bait And Switch]	0
4	2012-07-27	[Will AAPL Fall From The Tree]	1
5	2012-07-30	[Bulls Snatch Victory From Jaws Defeat]	1
6	2012-07-31	[What Driving China Real Estate Rally Part 3, ...	1
7	2012-08-02	[As Expected Apple Fills Post Earnings Gap, Po...	0
8	2012-08-07	[AAPL With An Eye On The Upcoming iPhone 5, Ch...	1
9	2012-08-10	[Good Knight Public Markets]	1
10	2012-08-14	[VIX Is Under 14 Now What, Largest USA Tech Co...	1

### Finbert Sentence Embeddings on the news Headlines

The implementation that will be used includes creating sentence embeddings for each headline and feed them into an RNN LSTM network to see if the stock price behavior can be predicted using only these headlines. Embeddings are low-dimensional, continuous vector representations of discrete variables that are learned when reading a text. The dimensionality of categorical variables is reduced and that allows us to meaningfully represent categories in the transformed space.

For the sentence embeddings on the headlines of the news the BERT base uncased corpus was used that was trained in a vocabulary of 30,000 words. FinBERT uses the BERT corpus and has been specifically trained on financial data but let us be more specific as to what that means.

BERT stands for Bidirectional Encoder Representations Transformers. It is a state-of-the-art NLP model that trains bidirectionally on text, which has proven to have a deeper sense of language context than common language models. That allows the model to spot the difference e.g between “near the river’s bank” and “the bank near the river”.

The fact that each day has a different number of news headlines was a problem, since the embeddings had to be created in a way that all dates would have an equal number of embeddings and the embeddings should be of equal size. To solve this problem and avoid dropping any information 3 metrics were used to map the different number of embeddings for each day: The mean, minimum and maximum value array. This way we lose less information, and each date now has a list of 3 arrays that represent the mean, minimum and maximum value of all the embeddings from this specific date. On top of that Finbert embeddings use embeddings of 768 length regardless of the sentence’s length.

## Model hyperparameter tuning

For the hyperparameter tuning of the model 3 parameters were tested: a) batch size, b) dropout and c) modes of Cyclic learning rate.

The cyclic learning rate is a callback that allows the learning rate of the model to be adjusted in a cyclical way. This allows the model to escape possible local optimum or minimum that might be encountered during training. Local optimum / minimum traps the model in a specific area and that way there is no way to know if the model reached the true optimal point or if there is a better one.

After implementing a thorough hyper parameter tuning 120 different models were built, all following a wide architecture with no additional LSTM layers. Below the parameters of the best model for Microsoft and Apple can be seen. The accuracy that was obtained from both models is approximately 60%.

Microsoft		Apple	
batch_size	16	batch_size	32
dropout	0.8	dropout	0.8
cyclic_mode	triangular2	cyclic_mode	triangular2
test_accuracy	60.2787	test_accuracy	60.6322
test_loss	0.883939	test_loss	0.802002
Name: 17, dtype: object		Name: 41, dtype: object	

## Model Results

Afterwards a confusion matrix and a classification report were created to see more details regarding the accuracy of the model on each class.

### Microsoft Confusion matrix

				precision	recall	f1-score	support
0				0.62	0.25	0.36	126
1				0.60	0.88	0.71	161
Test accuracy: 60.279 %							
0	1	accuracy				0.60	287
0	32	94	macro avg	0.61	0.56	0.54	287
1	20	141	weighted avg	0.61	0.60	0.56	287

## Apple Confusion matrix

Test accuracy: 60.632 %

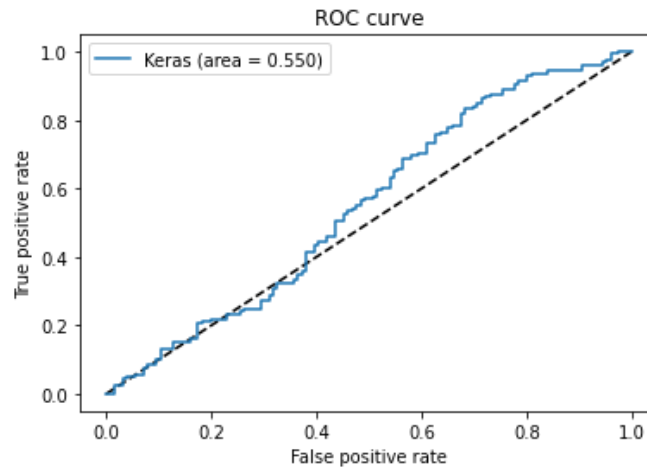
	0	1
0	21	127
1	10	190

## Apple Classification Report

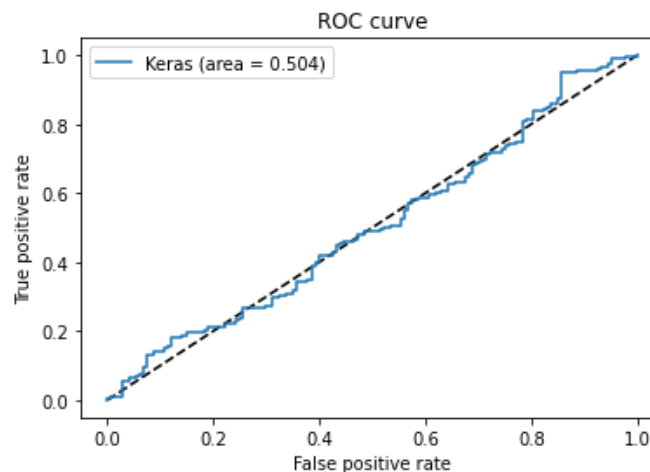
	precision	recall	f1-score	support
0	0.68	0.14	0.23	148
1	0.60	0.95	0.74	200
accuracy			0.61	348
macro avg	0.64	0.55	0.48	348
weighted avg	0.63	0.61	0.52	348

In the plots below we can see the ROC Curves for Apple and Microsoft. The ROC Curve plots the true positive rate versus the false positive. It is considered normal for a low accuracy of 60% to not have an exact curve. Nonetheless for Microsoft it seems that the true positive rate is better than Apple.

### Microsoft ROC Curve



### Apple ROC Curve

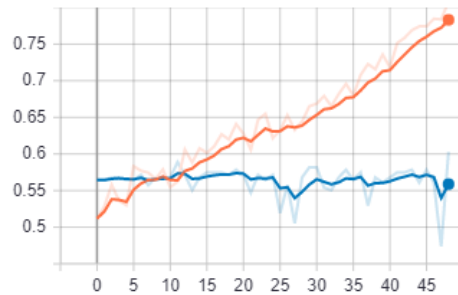


One more useful plot to inspect would be the tensorboard plots of the loss and the accuracy of the model as it was trained. The purple line refers to the train dataset and the blue to the validation dataset.

## Microsoft

epoch\_accuracy

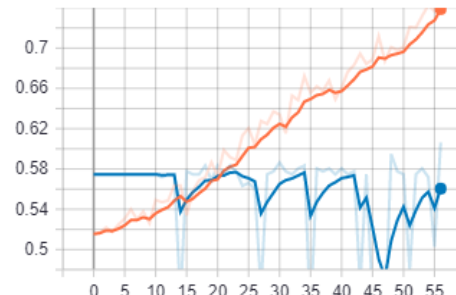
epoch\_accuracy



## Apple

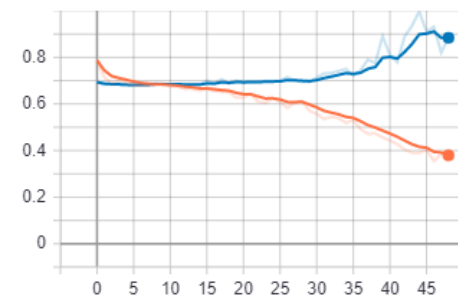
epoch\_accuracy

epoch\_accuracy



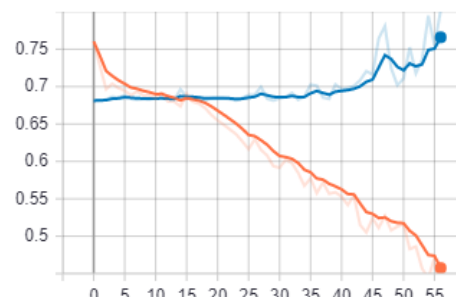
epoch\_loss

epoch\_loss



epoch\_loss

epoch\_loss



From the above analysis we can conclude that the Microsoft model seems to produce higher quality predictions despite the fact that both models have the same accuracy percentage. To make this clearer imagine that a similar dataset with 0 and 1 as target variables comprises of 1s by 57% and we have 2 models with the same accuracy of 57%. The one model predicts everything to be 1 and the other model predicts both classes of the dataset. The model we would choose would be the model that is able to predict both classes more accurately.

### Implementation 3: LSTM with Finbert Sentiment Analysis

The purpose of this side assignment is, initially, to calculate the sentiment of each headline of the dataset, regarding only Apple Inc. (AAPL) and Microsoft Inc. (MSFT). Afterwards, the headline sentiments must be grouped in such a way that each calendar day from the dataset is associated with the average sentiment of the corresponding daily headlines. Finally, defining the daily sentiments as input, a Long Short-Term Memory (LSTM) neural network will be implemented in order to create a classification model, that given a set of past daily sentiments, it will be able to decide if the next day's opening stock price will increase or decrease.

#### Methodology

In an attempt to create such a model, capable of predicting a stock price's surge or fall having news sentiment as input, we utilized the BERT (Bidirectional Encoder Representations from Transformers) model. BERT is a Natural Language Processing (NLP) model used for pre-training language representations. It leverages an enormous amount of plain text data publicly available (*Wikipedia and Google Books*) on the web and is trained in an unsupervised manner [Sun, C., Huang, L., & Qiu, X. (2019). *Utilizing BERT for aspect-based sentiment analysis via constructing auxiliary sentence*]. Because the news' headlines were in finance context and the meaning of a word or a sentence might differ from the meaning BERT would normally assign, we employed a pre-trained FinBERT model for the stock-related headline sentiment analysis. FinBERT is itself a model, able to analyze sentiment of financial text. It is built by further training the BERT language model in the finance domain, using a large financial corpus and thereby fine-tuning it for financial sentiment classification. The FinBERT implementation we applied, relies on HuggingFace's "pytorch\_pretrained\_bert" library [2] and their implementation of BERT for sequence classification tasks. The resulting sentiment labels were added next to its corresponding headline, in a new column in the existing dataset.

In order to obtain the daily stock prices for the stocks in question, we created a Python script and with the assistance of Yahoo Finance API, we created another column in the aforementioned dataset, occupying the Open prices for each stock in the corresponding date.

The resulting dataset will comprise of two (2) columns, the daily average sentiment and the opening price of the stock in question that day.



## Model Implementation

Having loaded the dataset file into a Pandas DataFrame , we retrieved its subset, which contained the stock tickers, headlines and the release date of each headline. Next, we filtered the DataFrame above, so as to retain only the news regarding Apple and Microsoft. The next course of action would be to handle the news published on weekends. We decided to make the assumption that the weekend news would be attributed to the Fridays in order to associate them more efficiently with the opening stock prices, which are affected by the news of the previous days. With the display, which was described above, we used the FinBert model in order to assign the sentiment to each headline, with the possible options being Negative, Neutral and Positive (*Figure 1*). When the sentiment labeling was completed, we transformed the sentiment values, in order to be in a numeric form in order to calculate their average sentiment. Thus, we assigned to the rows having Negative sentiment the number 1, the rows with Neutral sentiment the number 2 and those with Positive sentiment the number 3. Next, we grouped the headline sentiments by each date in such a form that each date was associated with their average and finally utilizing the yahoo-api script we merged the opening stock prices in their corresponding day. After the merge we transformed the prices into a binary categorical variable of 0,1, with 0 being if that day's price decreased from the previous one, and 1 if the price increased.

In summary, the final form of the dataframe had the columns “sentiment” and “Open”, which were the average sentiment of each day and the opening stock price the same day (*Figure 2*).

ticker	title	release_date	sentiment
MSFT	GLOBAL MARKETS U S stocks retreat on recession fears oil falls	2008-10-15	neutral
MSFT	GLOBAL MARKETS U S stocks oil slide on global recession fears	2008-10-15	neutral
MSFT	INTERVIEW UPDATE 2 Nintendo Wii DS sales strong	2008-12-08	neutral
MSFT	Nintendo Dentsu to distribute video on Wii Nikkei	2008-12-24	neutral
MSFT	UPDATE 1 Nintendo Dentsu to offer video on Wii	2008-12-24	neutral
MSFT	COLUMN What Apple loses without Steve Eric Auchard	2009-01-15	negative
MSFT	INTERVIEW Namco Bandai aims to double game profit in 3 yrs	2009-01-19	neutral
MSFT	WRAPUP 3 Qimonda Samsung hit by chip market crash	2009-01-23	neutral
MSFT	UPDATE 2 Microsoft to open own stores take on Apple	2009-02-12	neutral
MSFT	ANALYSIS Video gaming could be bargain buy of UK equities	2009-03-25	neutral

Figure 1

	sentiment	Open
1425	2.200000	1.0
1426	2.000000	0.0
1427	2.500000	1.0
1428	2.300000	0.0
1429	2.000000	1.0
1430	2.000000	0.0
1431	1.700000	1.0
1432	2.000000	1.0
1433	2.274194	1.0
1434	2.500000	0.0

Figure 2

For the LSTM classification model, we implemented, the input will be the sentiments for a specific number of days back and the target variable will be the binary categorical variable. In order to decide which time window of news sentiment provides the best accuracy in the increase / decrease prediction, we introduced multiple time windows, looking 60, 30, 15 and 5 days back. Due to the imbalance of both Apple (Figure 3) and Microsoft datasets (Figure 4), there has been observed a certain level of failure to correctly predict the days that the stock prices will fall, with the exception of relatively correction predictability, when the model had a timestep of only five (5) days. Regarding the aforementioned imbalance, instead of a Receiver Operating Characteristic (ROC) curve, the visualized evaluation method of choice was the Precision-Recall curve of the model compared to a no-skill classifier. A no-skill classifier is one, that cannot discriminate between the classes and would predict a random class or a constant class in all cases. The no-skill line changes based on the distribution of the positive to negative classes. It is a horizontal line with the value of the ratio of positive cases in the dataset. For a balanced dataset, this is 0.5. In conclusion, of the models implemented, the most reliable in terms of both accuracy and low error rate is the one with the 5-day timestep, which produces ~56% accuracy, and the type-2 error is at its lowest value.

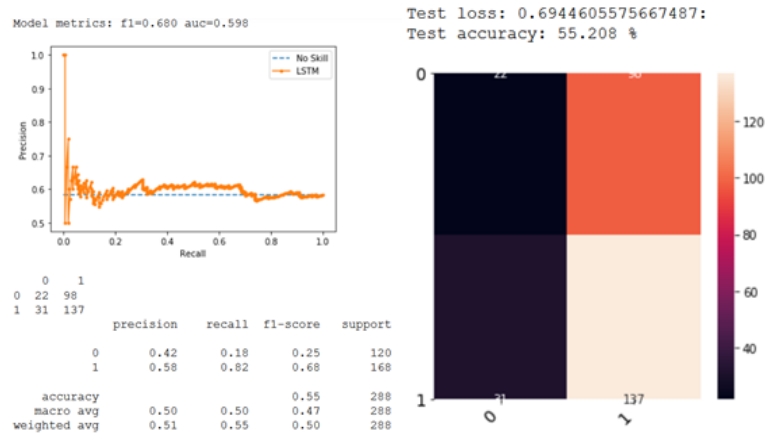
Count		Count	
Open		Open	
0	797	0.0	641
1	942	1.0	793

Figure 3

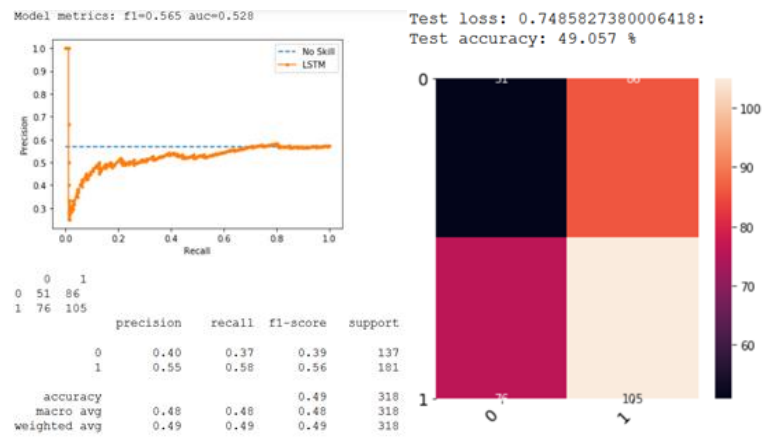
Figure 4

## Apple Dataset Metrics:

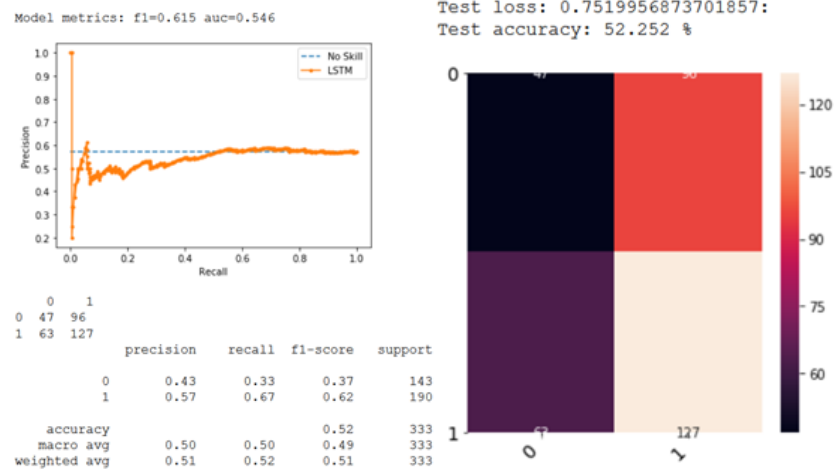
- 60-day timestep



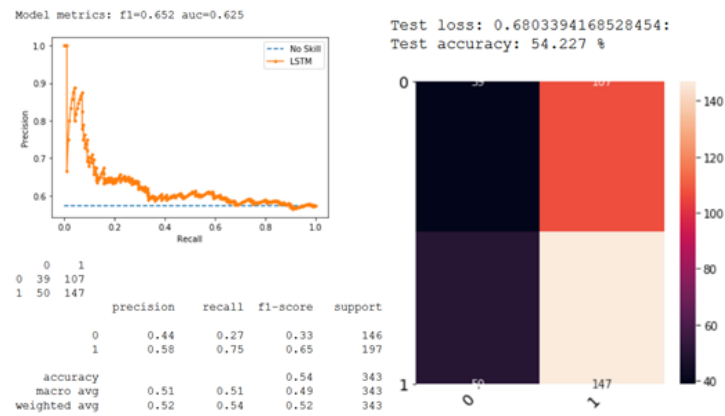
- 30-day timestep



- 15-day timestep

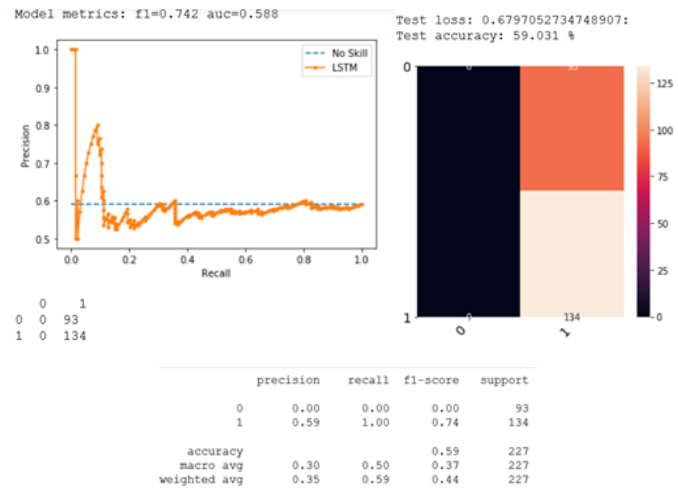


- 5-day timestep

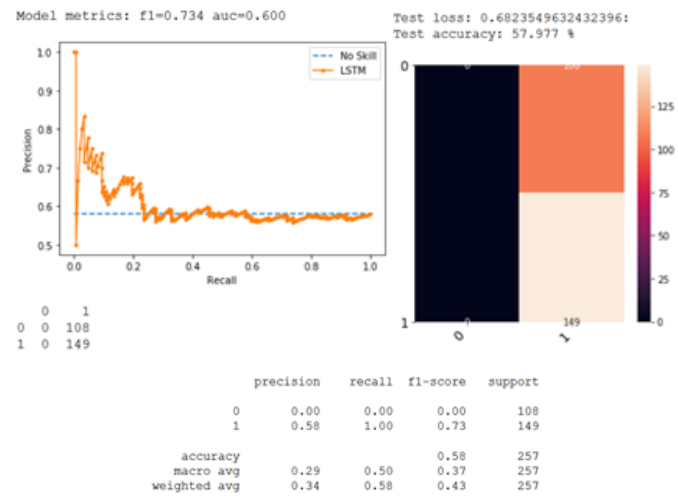


## Microsoft Dataset Metrics

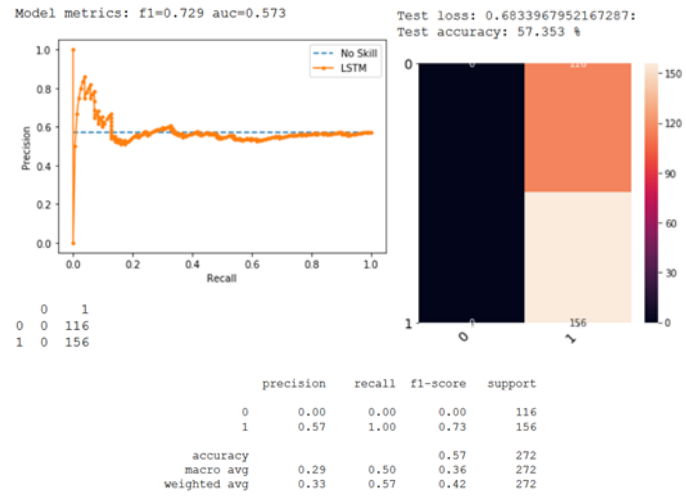
- 60-day timestep



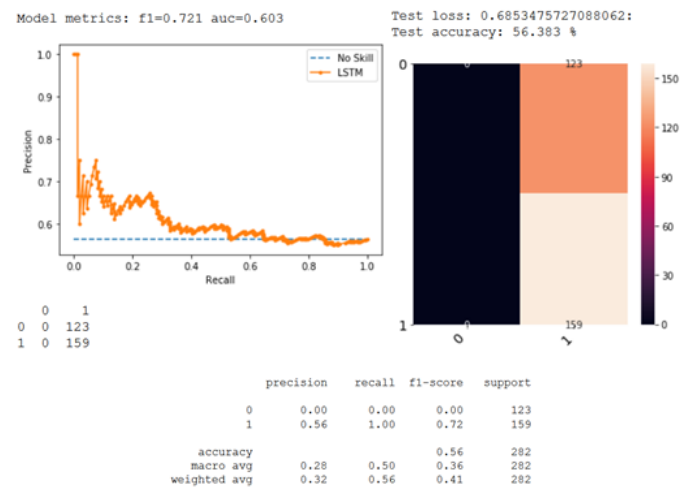
- 30-day timestep



- 15-day timestep



- 5-day timestep



## Implementation 4: Arima Model - Forecasting Using Time Series Analysis

Since stock prices are not randomly generated values, they can be treated as a discrete-time series model which is based on a set of well-defined numerical data items collected at successive points at regular intervals of time. In this implementation we are going to use Autoregressive Integrated Moving Average Model (ARIMA) instead of simple forecasting, because it is essential to analyze trends of stock prices in order to identify a model.

### Autoregressive Integrated Moving Average Model

An ARIMA model is a class of statistical models for analyzing and forecasting time series data. It explicitly caters to a suite of standard structures in time series data, and as such provides a simple yet powerful method for making skillful time series forecasts.

ARIMA is an acronym that stands for AutoRegressive Integrated Moving Average. It is a generalization of the simpler AutoRegressive Moving Average and adds the notion of integration. This acronym is descriptive, capturing the key aspects of the model itself. Briefly, they are:

- AR: Autoregression. A model that uses the dependent relationship between an observation and some number of lagged observations.
- I: Integrated. The use of differencing of raw observations (e.g. subtracting an observation from an observation at the previous time step) in order to make the time series stationary.
- MA: Moving Average. A model that uses the dependency between an observation and a residual error from a moving average model applied to lagged observations.

Each of these components are explicitly specified in the model as a parameter. A standard notation is used of ARIMA(p,d,q) where the parameters are substituted with integer values to quickly indicate the specific ARIMA model being used.

The parameters of the ARIMA model are defined as follows:

- p stands for the number of autoregressive terms, i.e. the number of observations from past time values used to forecast future values. e.g. if the value of p is 2, then this means that two previous time observations in the series are being used to forecast the future trend.
- d denotes the number of differences needed to make the time series stationary (i.e. one with a constant mean, variance, and autocorrelation). For instance, if d = 1, then it means that a first-difference of the series must be obtained to transform it into a stationary one.
- q represents the moving average of the previous forecast errors in our model, or the lagged values of the error term. As an example, if q has a value of 1, then this means that we have one lagged value of the error term in the model.

A linear regression model is constructed including the specified number and type of terms, and the data is prepared by a degree of differencing in order to make it stationary, i.e. to remove trend and seasonal structures that negatively affect the regression model.

Before we can build a model, we must ensure that the time series is stationary. Stationarity is a very important concept in Time Series Analysis. In order to apply a time series model, it is important for the Time series to be stationary; in other words all its statistical properties (mean, variance) remain constant over time. This is done basically because if we consider a certain behavior over time, it is important that this behavior is same in the future in order to forecast the series. In practice we can assume the series to be stationary if it has constant statistical properties over time and these properties can be:

- constant mean
- constant variance
- an auto co-variance that does not depend on time.

There are two primary way to determine whether a given time series is stationary.

- Rolling Statistics: Plot the rolling mean and rolling standard deviation. The time series is stationary if they remain constant with time (with the naked eye look to see if the lines are straight and parallel to the x-axis).
- Augmented Dickey-Fuller Test: The time series is considered stationary if the p-value is low (according to the null hypothesis) and the critical values at 1%, 5%, 10% confidence intervals are as close as possible to the ADF Statistics

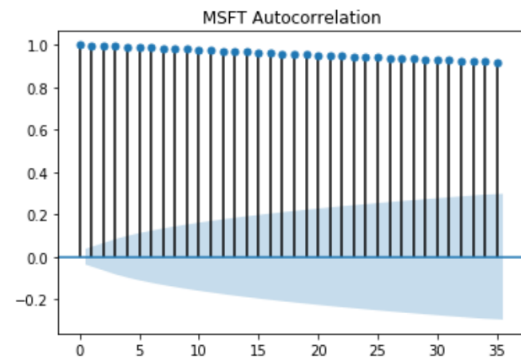
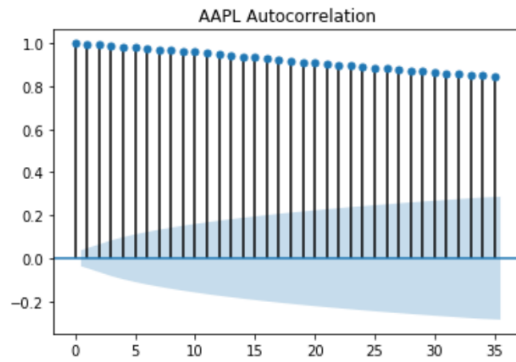
## Methodology

To the begin with, data was transposed in a pandas Series with date as index. We were interested in forecasting opening stock values, so from the original datasets only Open values and dates were kept.

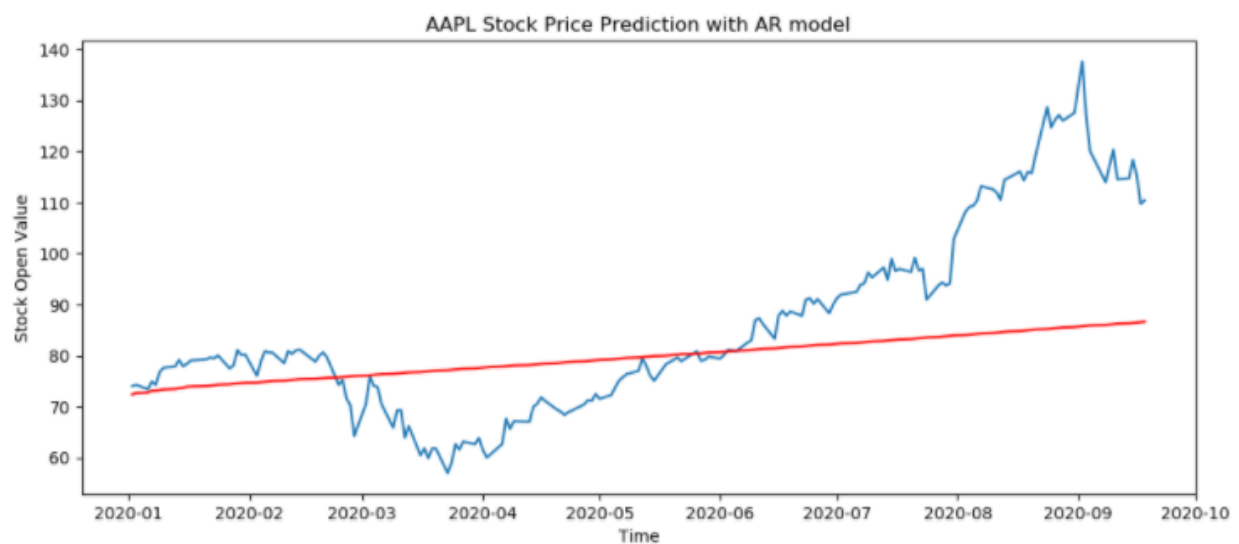
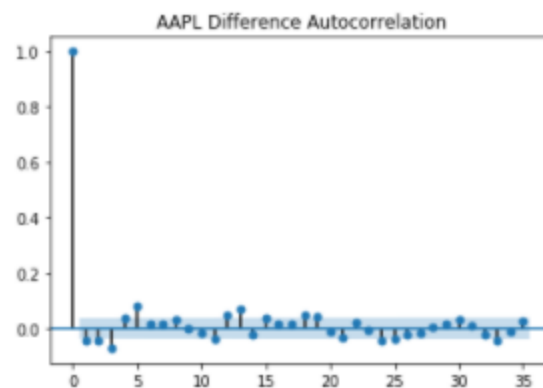
Open		Open	
Date		Date	
2009-09-21	6.581786	2009-09-21	25.110001
2009-09-22	6.613929	2009-09-22	25.400000
2009-09-23	6.621428	2009-09-23	25.920000
2009-09-24	6.685714	2009-09-24	25.920000
2009-09-25	6.500357	2009-09-25	25.690001

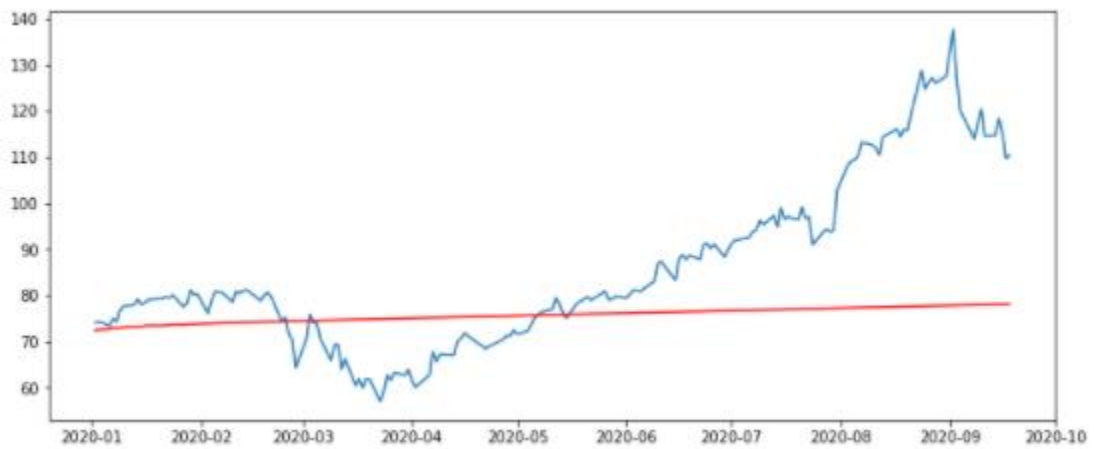
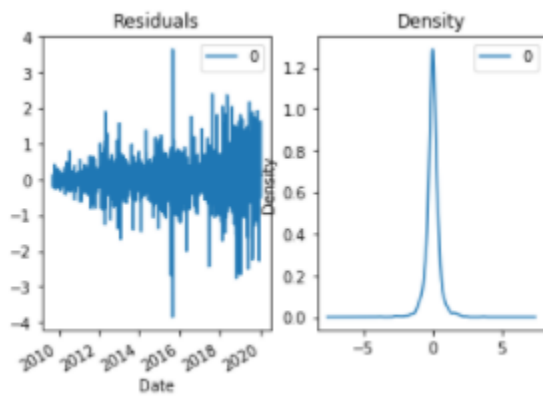
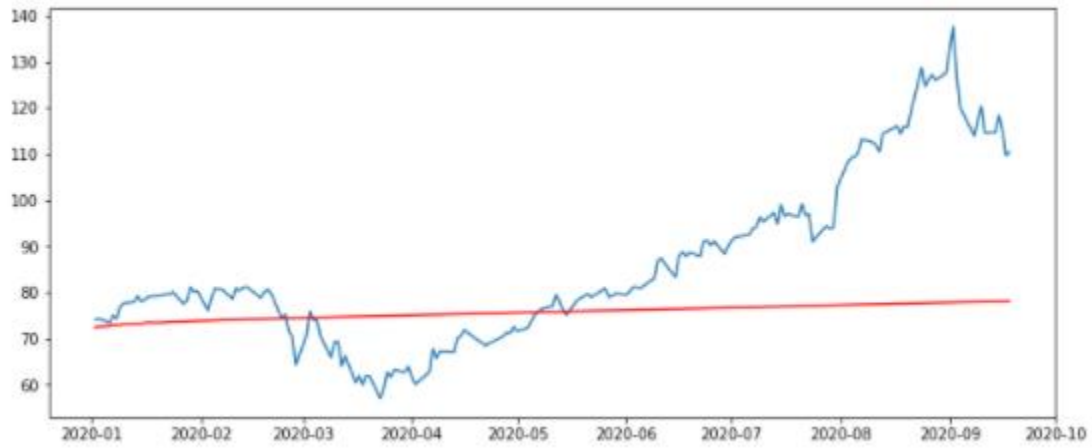
From the Figure 1 it's obvious that data are not stationary. We can also assume that from the autocorrelation plots.

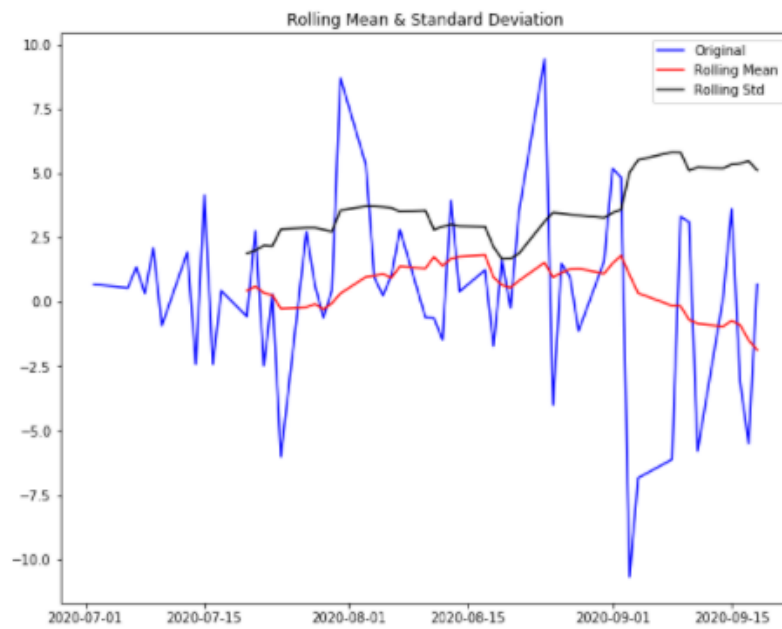
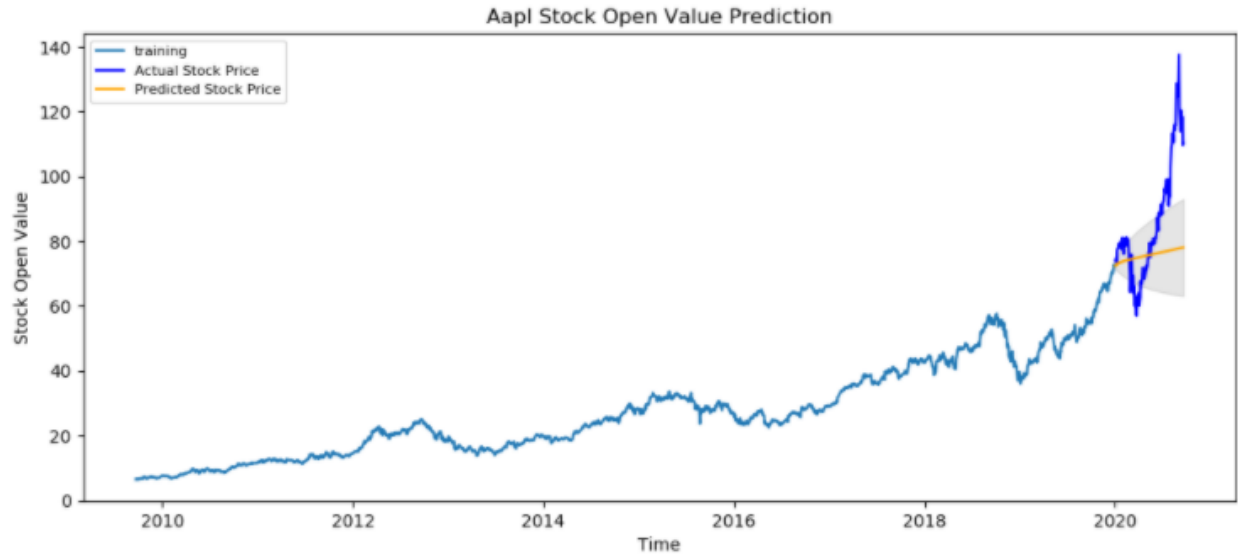




Since autocorrelation is a function of time (decreasing or increasing smoothly following a pattern) data is not stationary. Data has to be transformed in order to become stationary and proceed with our analysis. The most common transformation used in time series is to calculate the difference with previous value.



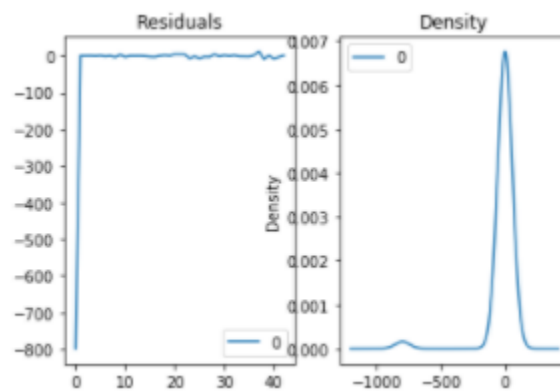
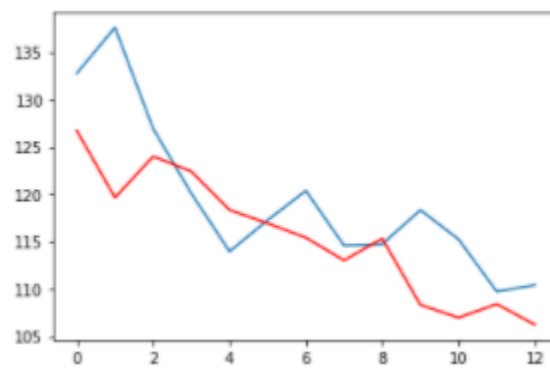
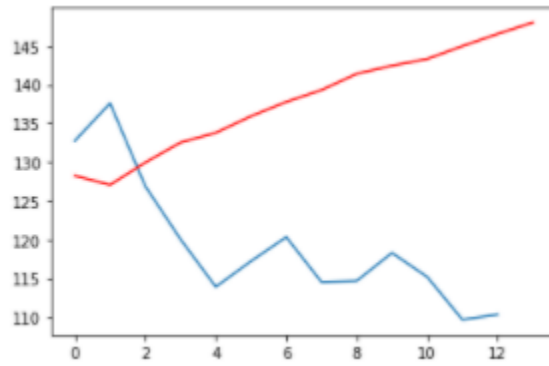


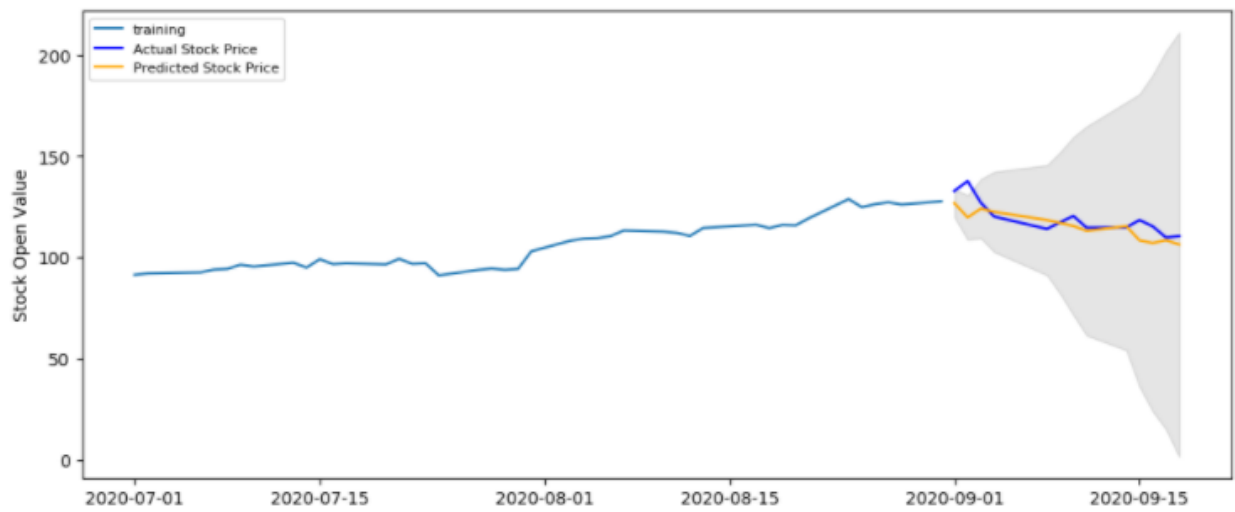
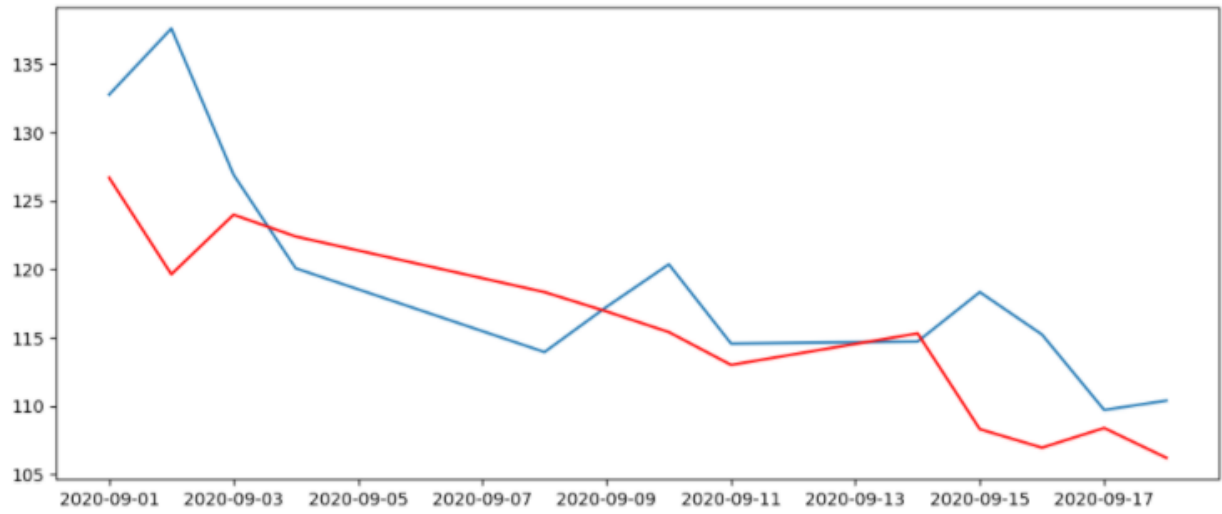


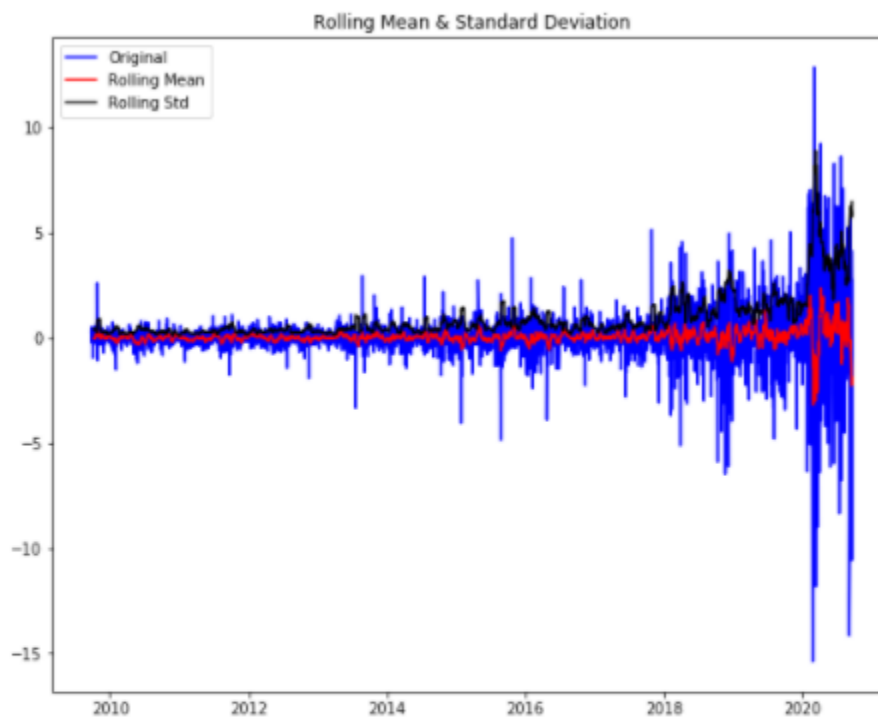
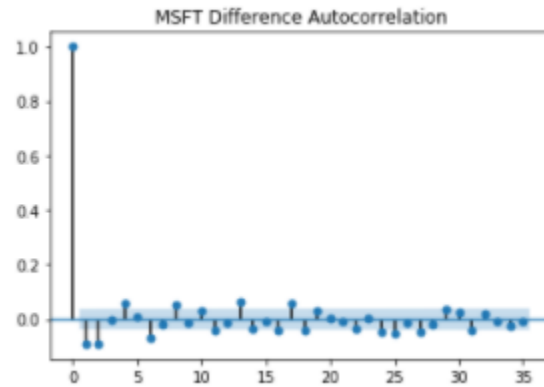
Results of Dickey-Fuller Test:

Test Statistic	-6.755884e+00
p-value	2.872652e-09
#Lags Used	0.000000e+00
Number of Observations Used	5.400000e+01
Critical Value (1%)	-3.557709e+00
Critical Value (5%)	-2.916770e+00
Critical Value (10%)	-2.596222e+00

dtype: float64



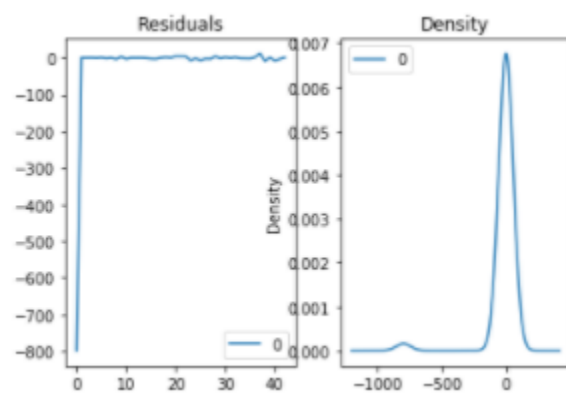
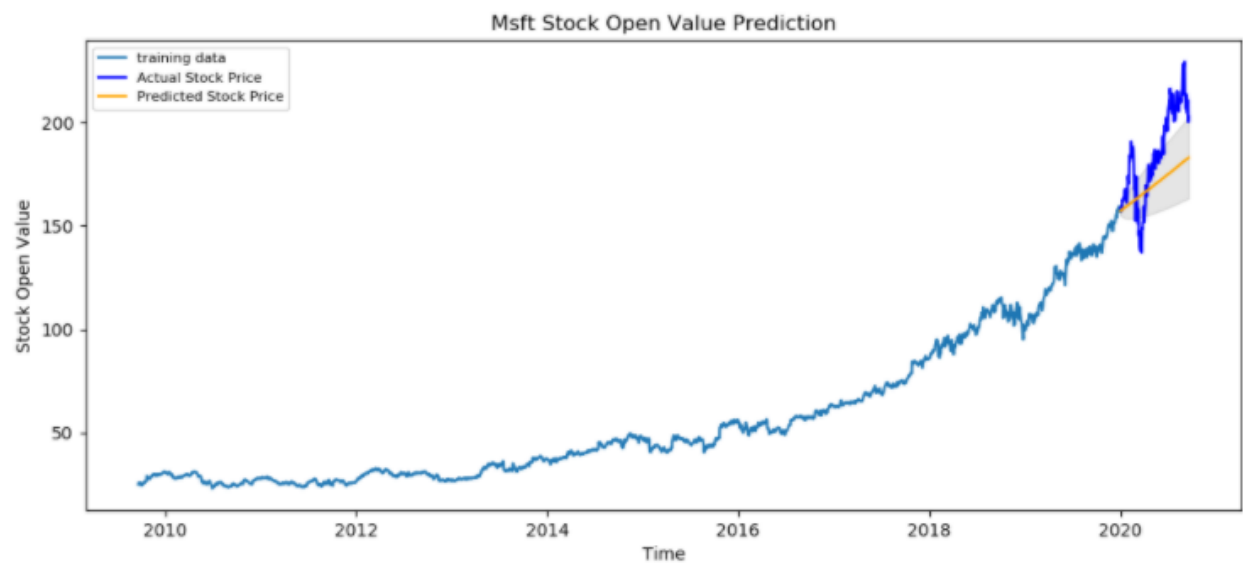
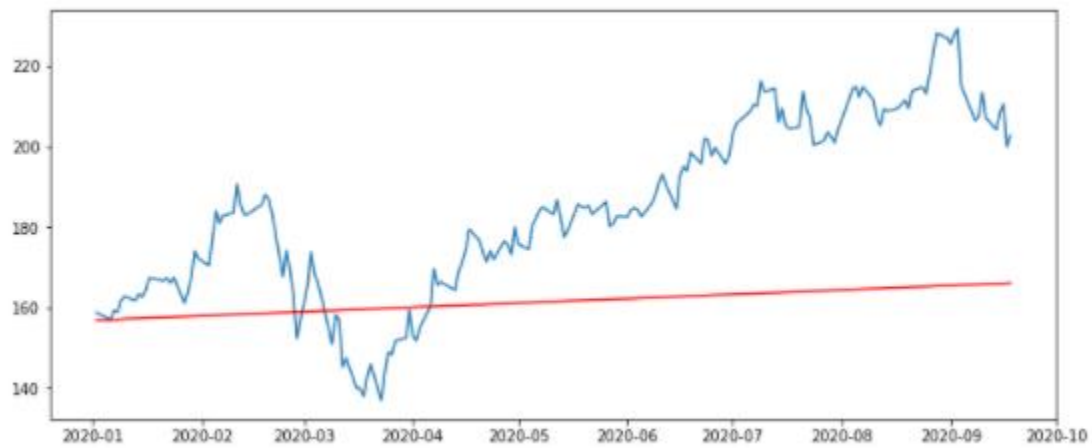


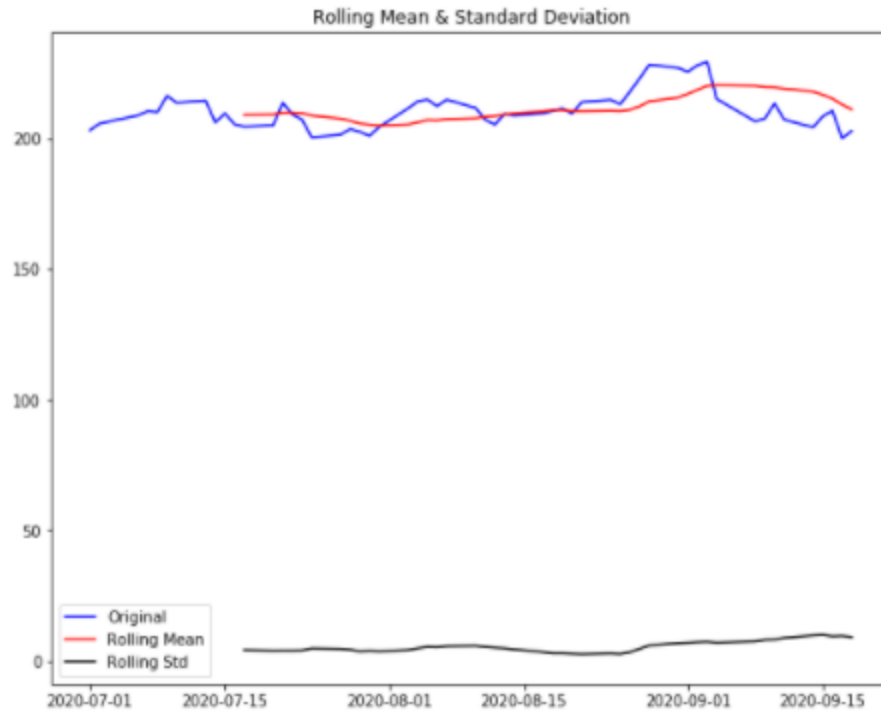


Results of Dickey-Fuller Test:

Test Statistic	-1.272080e+01
p-value	9.837120e-24
#Lags Used	2.700000e+01
Number of Observations Used	2.740000e+03
Critical Value (1%)	-3.432739e+00
Critical Value (5%)	-2.862595e+00
Critical Value (10%)	-2.567332e+00

dtype: float64



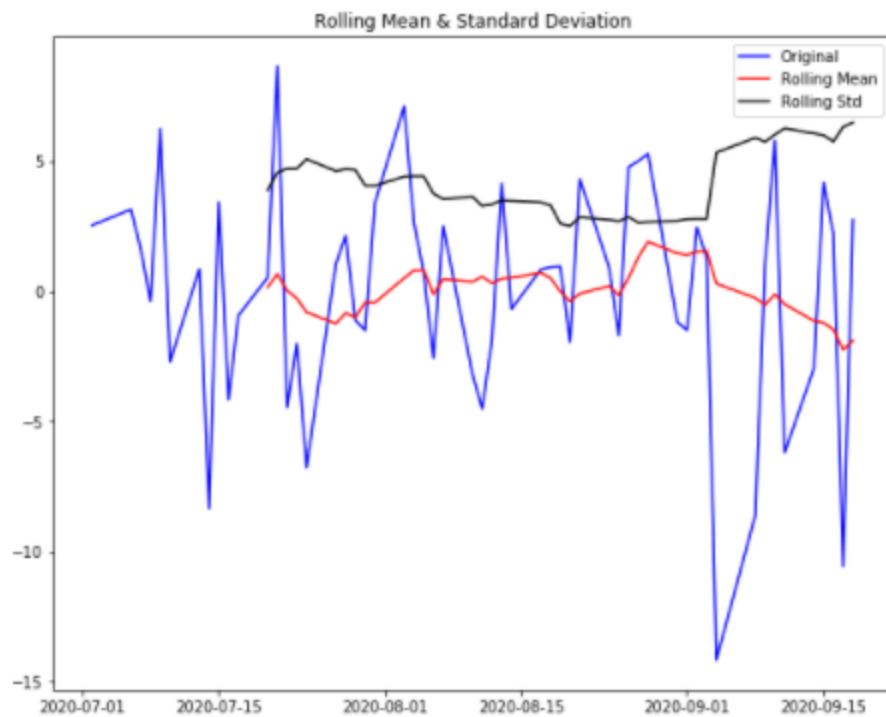


Results of Dickey-Fuller Test:

Test Statistic	-1.987332
p-value	0.292155
#Lags Used	5.000000
Number of Observations Used	50.000000
Critical Value (1%)	-3.568486
Critical Value (5%)	-2.921360
Critical Value (10%)	-2.598662

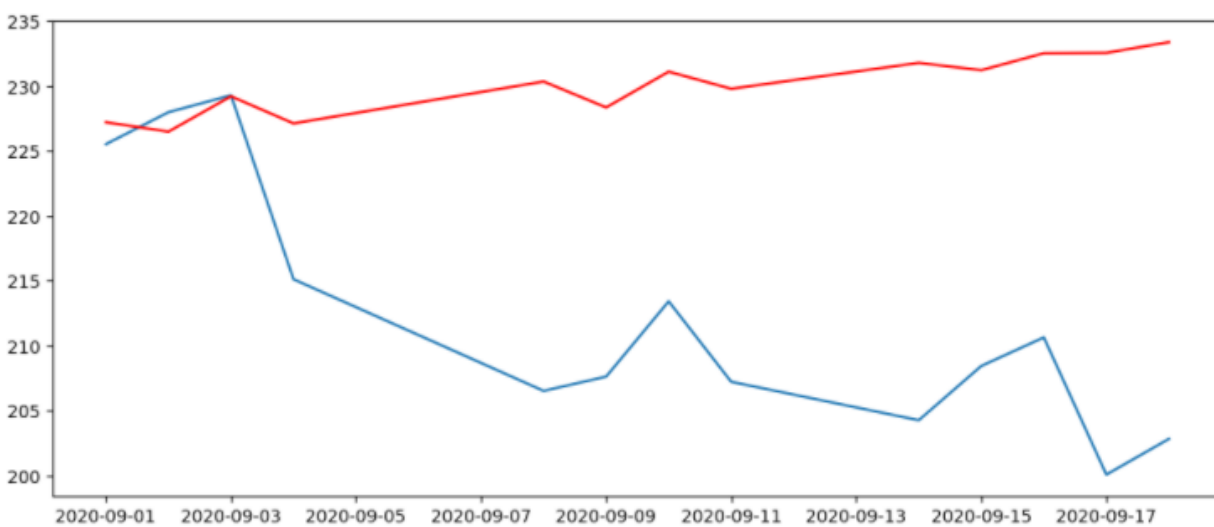
dtype: float64

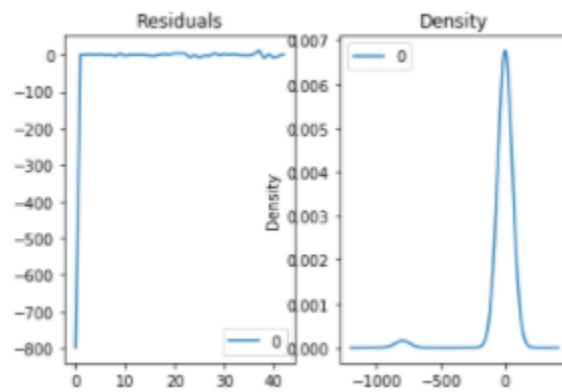
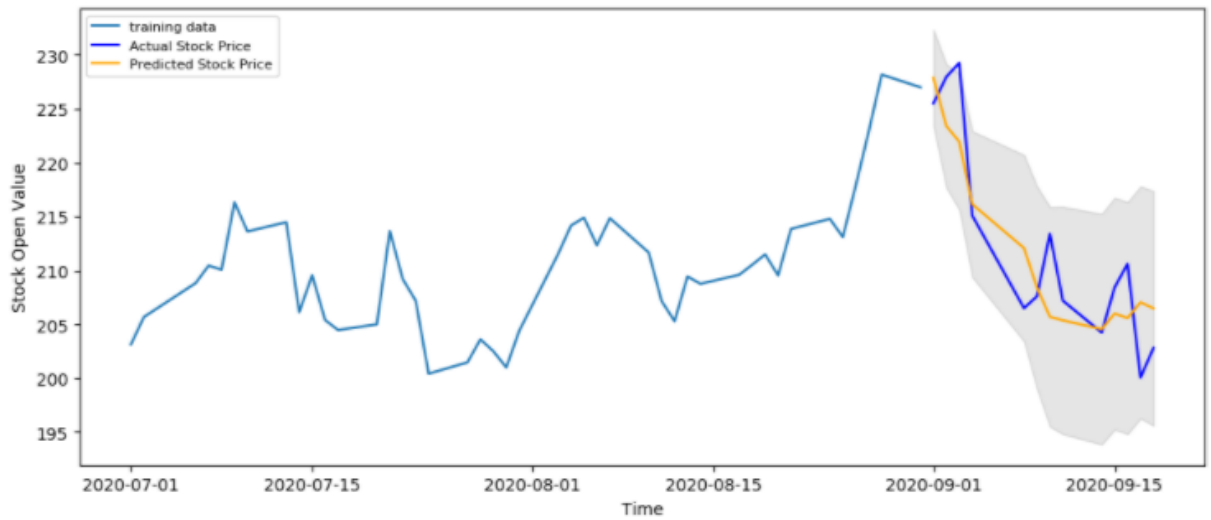
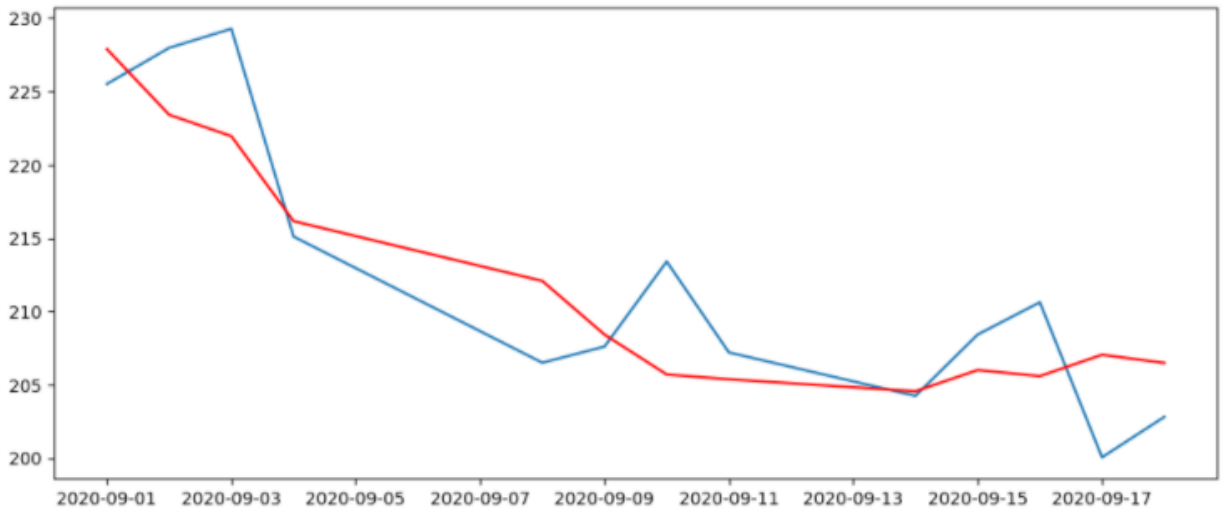




Results of Dickey-Fuller Test:

Test Statistic	-3.115479
p-value	0.025429
#Lags Used	11.000000
Number of Observations Used	43.000000
Critical Value (1%)	-3.592504
Critical Value (5%)	-2.931550
Critical Value (10%)	-2.604066
dtype:	float64





## Timeplan

Start	End	Task
8/5/2020	8/12/2020	Research for Subject. Subject chosen predict stocks movements from twitter sentiment analysis
9/1/2020	9/5/2020	Subject was transformed. We have started collecting previous years tweets with Library GetOldTweets but Twitter has aborted our connected endpoint
9/6/2020	9/8/2020	Then, we have found a new data source to extract sentiment, an online news Dataset on Kaggle. Started explore articles and papers about the subject
9/11/2020	9/18/2020	We have found the different implementations to be done, and we have separate tasks for every team member. Additionally, data cleaning and transformation was started
9/19/2020	9/26/2020	Everybody works on his implementation/tasks
9/27/2020	10/4/2020	Get in touch to aggregate functionalities /implementations, solve bugs, see each other code
10/5/2020	10/11/2020	Run our models and get results. Get everything together in a unified report

### Members/Roles

Name	Role	Task
Mallios Charalampos	ML Engineer Data Analyst Business Analyst	<ul style="list-style-type: none"> <li>Implementation 1: LSTM with Historical Stock Prices and TextBlob Sentiment Score</li> <li>Data Analysis.</li> </ul>
Patrikalos Vangelis	ML Engineer Data Analyst	<ul style="list-style-type: none"> <li>Implementation 2: LSTM with Finbert sentence embeddings</li> </ul>
Protopapas Nikolaos	ML Engineer Data Analyst	<ul style="list-style-type: none"> <li>Implementation 3: LSTM with Finbert sentiment analysis</li> </ul>
Ralis Dionysios	ML Engineer Data Analyst	<ul style="list-style-type: none"> <li>Arima Model - Forecasting Using Time Series Analysis</li> </ul>

### Contact Person (phone & email)

Name	Phone	Mail
Mallios Charalampos	6944137907	Mallioscharis@gmail.com

## Bibliography

Michael Grogan, (August 28, 2019), SARIMA: Forecasting Seasonal Data with Python and R (Online), <https://medium.com/analytics-vidhya/sarima-forecasting-seasonal-data-with-python-and-r-2e7472dfad83>.

Cory Maklin, (May 25, 2019), ARIMA Model Python Example — Time Series Forecasting (Online), <https://towardsdatascience.com/machine-learning-part-19-time-series-and-autoregressive-integrated-moving-average-model-arima-c1005347b0d7>.

Jason Brownlee, (January 9, 2017), How to Create an ARIMA Model for Time Series Forecasting in Python (Online), <https://machinelearningmastery.com/arima-for-time-series-forecasting-with-python/>.

Jason Brownlee, (March 24, 2017), How to Make Out-of-Sample Forecasts with ARIMA in Python (Online), <https://machinelearningmastery.com/make-sample-forecasts-arima-python/>.

[www.hackdeploy.com](http://www.hackdeploy.com), (November 20, 2018), Augmented Dickey-Fuller Test in Python.

Jacob Stallone, (November 9, 2017), Time Series Forecast : A basic introduction using Python (Online), <https://medium.com/@stallonejacob/time-series-forecast-a-basic-introduction-using-python-414fcb963000>.

Selva Prabhakaran, (February 18, 2019), ARIMA Model – Complete Guide to Time Series Forecasting in Python (Online), <https://www.machinelearningplus.com/time-series/arima-model-time-series-forecasting-python/>.

datamites.com, ARIMA in Python - Time Series Forecasting Part 1 & 2, <https://www.youtube.com/watch?v=D9y6dcy0xK8>.

Nagesh Singh Chauhan, (January 20, 2020), Stock Market Forecasting Using Time Series Analysis, <https://www.kdnuggets.com/2020/01/stock-market-forecasting-time-series-analysis.html>

<https://medium.com/@Currie32/predicting-the-stock-market-with-the-news-and-deep-learning-7fc8f5f639bc>

<https://www.thepythoncode.com/article/stock-price-prediction-in-python-using-tensorflow-2-and-keras>

<https://towardsdatascience.com/getting-rich-quick-with-machine-learning-and-stock-market-predictions-696802da94fe>

<http://davidanastasiu.net/pdf/papers/2019-MohanMSVA-BDS-stock.pdf?fbclid=IwAR3fTE1IO wYSUWdeIBKyEXdan0CauOJ5aw7ULq2pnVAXJMpjfyQZl1DjwM>

<http://davidanastasiu.net/pdf/papers/2019-MohanMSVA-BDS-stock.pdf?fbclid=IwAR3fTE1IO wYSUWdeIBKyEXdan0CauOJ5aw7ULq2pnVAXJMpjfyQZl1DjwM>

<http://github.com/huggingface/transformers>

<http://github.com/yya518/FinBERT>

<https://towardsdatascience.com/adaptive-and-cyclical-learning-rates-using-pytorch-2bf904d18dee>

<https://github.com/bckenstler/CLR>

[https://github.com/Currie32/Predicting-the-Dow-Jones-with-Headlines/blob/master/Predict\\_Dow\\_with\\_News.ipynb](https://github.com/Currie32/Predicting-the-Dow-Jones-with-Headlines/blob/master/Predict_Dow_with_News.ipynb)

<https://pypi.org/project/finbert-embedding/>

<https://medium.com/prosus-ai-tech-blog/finbert-financial-sentiment-analysis-with-bert-b277a3607101>

<https://towardsdatascience.com/bert-explained-state-of-the-art-language-model-for-nlp-f8b21a9b6270>

<https://huggingface.co/bert-base-uncased>

<https://medium.com/analytics-vidhya/sentiment-classification-with-bert-embeddings-56607b45daad>

