

Machine Learning Nanodegree

Capstone Project: Stock Price Predictor.

Author: Oscar Martinez Rico

Definition

[Project Overview](#)

[Figure 2.0](#)

[Problem Statement](#)

[Metrics](#)

Analysis

[Data Exploration](#)

[Exploration Visualization](#)

[Algorithms and Techniques](#)

[Benchmark](#)

Methodology

[Data Preprocessing](#)

[Implementation](#)

[Supervised Learning Algorithms](#)

[Refinement](#)

[Linear Regression Fine Tune](#)

[Lasso Regression Fine Tune](#)

Results

[Model Evaluation and Validation](#)

[Justification](#)

Conclusion

[Free-form Visualization](#)

[Reflection](#)

[Improvement](#)

References

I. Definition

Project Overview

Stock is a type of security that express ownership in a company and represents a claim on the part of the corporation's assets and earnings (Hayes, 2017). Stocks are also known as "shares" or "equity." There are two main types of stocks:

There are two main types of stocks include:

1. A common stock usually gives the owner to vote at shareholders' meetings and to receive dividends.
2. A preferred stock does not have voting rights but has a higher claim on assets and earnings than the common shares.

A holder of stock, also known as a shareholder, has a claim to a part of the company's assets and earnings. In other words, a shareholder is an owner of a company. The ownership of a company is calculated by the number of shares a person has relative to the number of outstanding shares. Below is an example (Figure 1.0) of a stock information sheet.

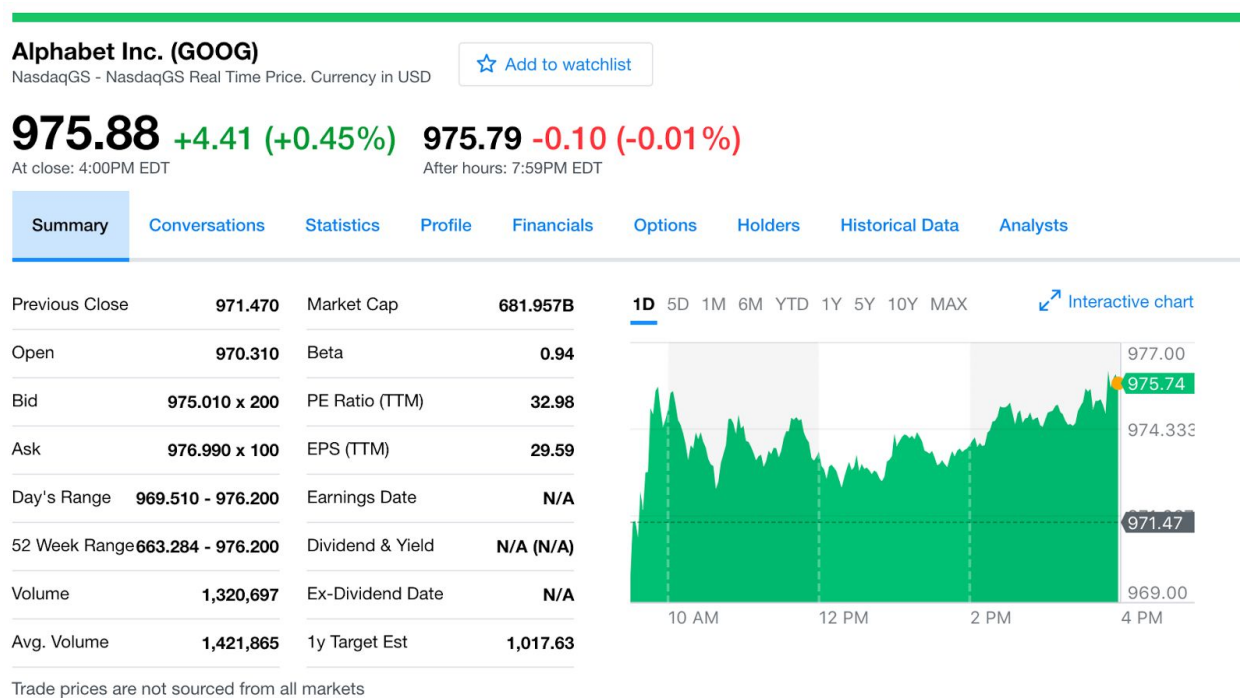


Figure 1.0

A stock quote provides information that includes: the current bid and offer prices, and the last price the stock was traded at. The highest price that an individual is willing to pay at a given time is the bid price. If an individual is interested in buying stocks, they have to make a bid. When the price of a bid and offer coincide, a trade is affected.

There is more information about a stock for example, trading volume, this is number of shares traded. Most of the time stock information is obtained online. To obtain the stocks information, we can search by their symbol. A stock symbol is formed using between one to four capital letters, which corresponds to the company name. Below are a picture(**Figure 2.0**) of some stocks symbols and their prices.

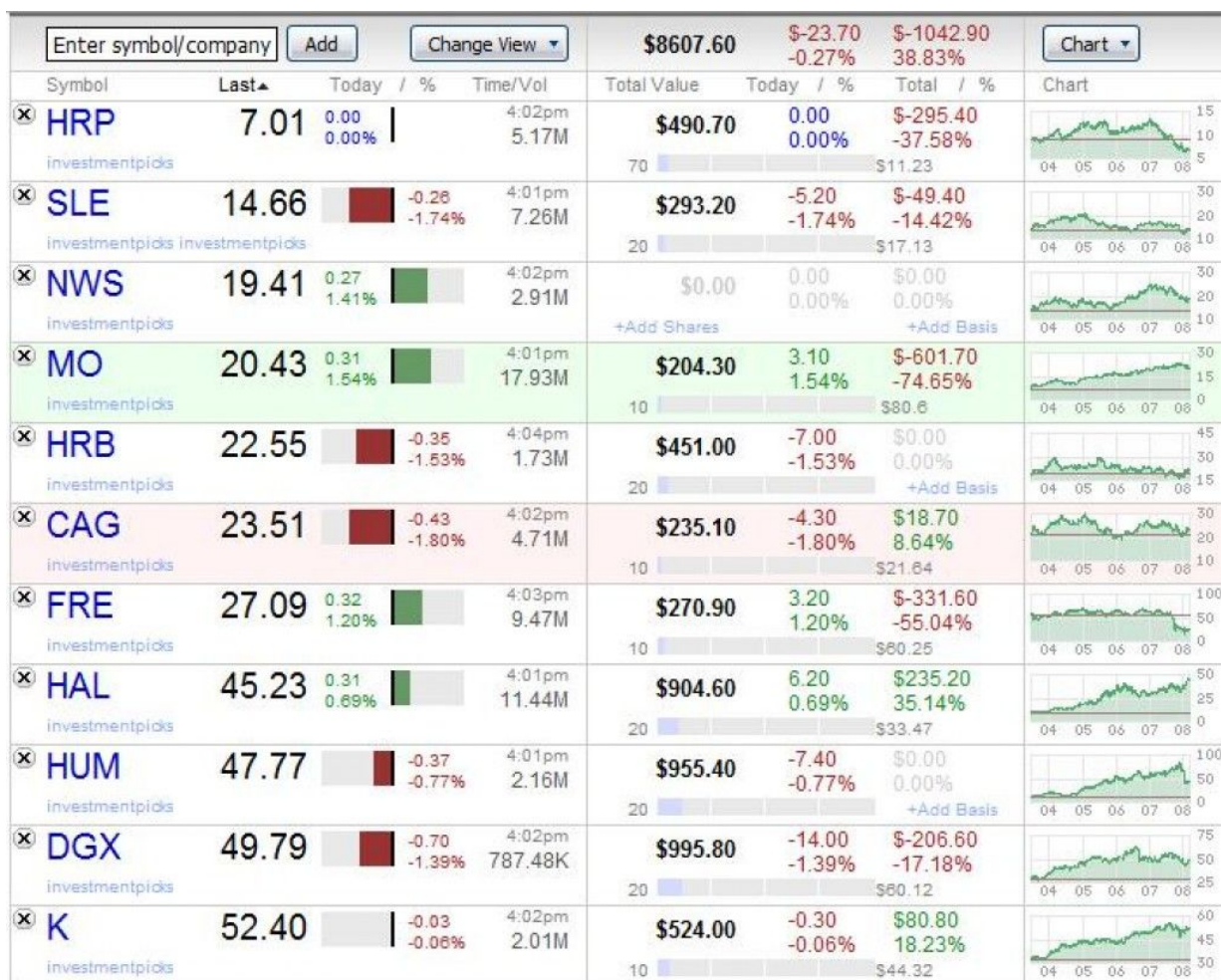


Figure 2.0

Problem Statement

There is no doubt that investing in the stock market can be one of the most exciting ways to invest money, especially when you can see your money multiply and grow.

Building a stock price predictor takes daily trading data over a certain date range as input, and outputs projected estimates for given query dates. By looking at the historical data of a given stock as an input, the stock predictor application will train the model to predict the Adjusted Close value for any given stock in the future. The model will be created using regression algorithms as we are attempting to predict the stock price. Having such a software system will benefit individuals and companies to make more educated decisions managing their stock portfolio.

Metrics

To evaluate the trained, supervised learning model R^2 will be used. The adequacy of the regression specification is often assessed by reference to the coefficient of (multiple) determination, commonly denoted R^2 (**Figure 3.0**), a statistic, whose value is bounded by zero and unity, which measures the proportion of the variation in the response variable explained by the regression. As R^2 declines from its maximum attainable value, it is perforce the case that the variance of the predicted, or fitted, values declines relative to that of the response variable itself (Dancer, 2005).

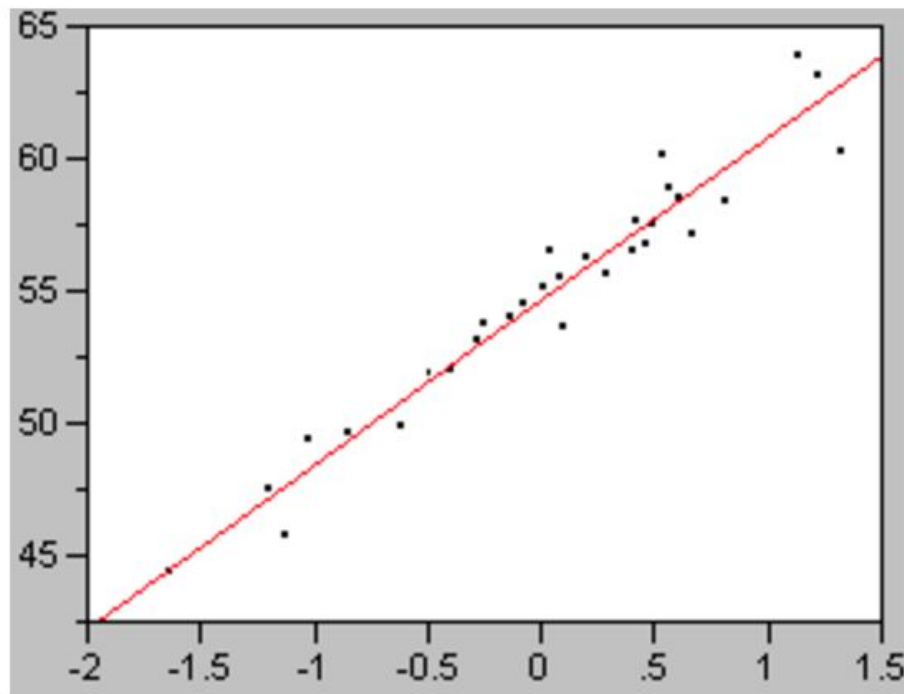


Figure 3.0

R² is a statistical measure of how close the data are to the fitted regression line. It is also known as the coefficient of determination.

The definition of R-squared is fairly straight-forward; it is the percentage of the response variable variation that is explained by a linear model. Or:

$$R^2 = \frac{\text{Explained variation}}{\text{Total variation}}$$

R² is always between 0.0 and 1.0:

- 0 indicates that the model explains none of the variability of the response data around its mean.
- 1.0 indicates that the model explains all the variability of the response data around its mean.

In general, the higher the R-squared, the better the model fits your data. However, there are important conditions for this guideline that I'll talk about both in this post and my next post.

The best possible score is 1.0. If the **R²** score is negative, it shows that the model can be arbitrarily worse. A constant model that always predicts the expected value of y, disregarding the input features, would get an **R²** score of 0.0.

II. Analysis

Data Exploration

Yahoo finance was used to obtain the historical dataset for this project. The dataset contains six features as it can be seen in the following picture(**Figure 4.0**):

- *Open*: The opening price is the price at which a security first trades upon the opening of an exchange on a given trading day.
- *High*: the highest price at which a stock traded during the course of the day.
- *Low*: the lowest price at which a stock traded during the course of the day.
- *Close*: generally refers to the last price at which a stock trades during a regular trading session.
- *Volume*: the number of shares or contracts traded in a security or an entire market during a given period of time.
- *Adjusted Close*: is a stock's closing price on any given day of trading that has been amended to include any distributions and corporate actions that occurred at any time prior to the next day's open.

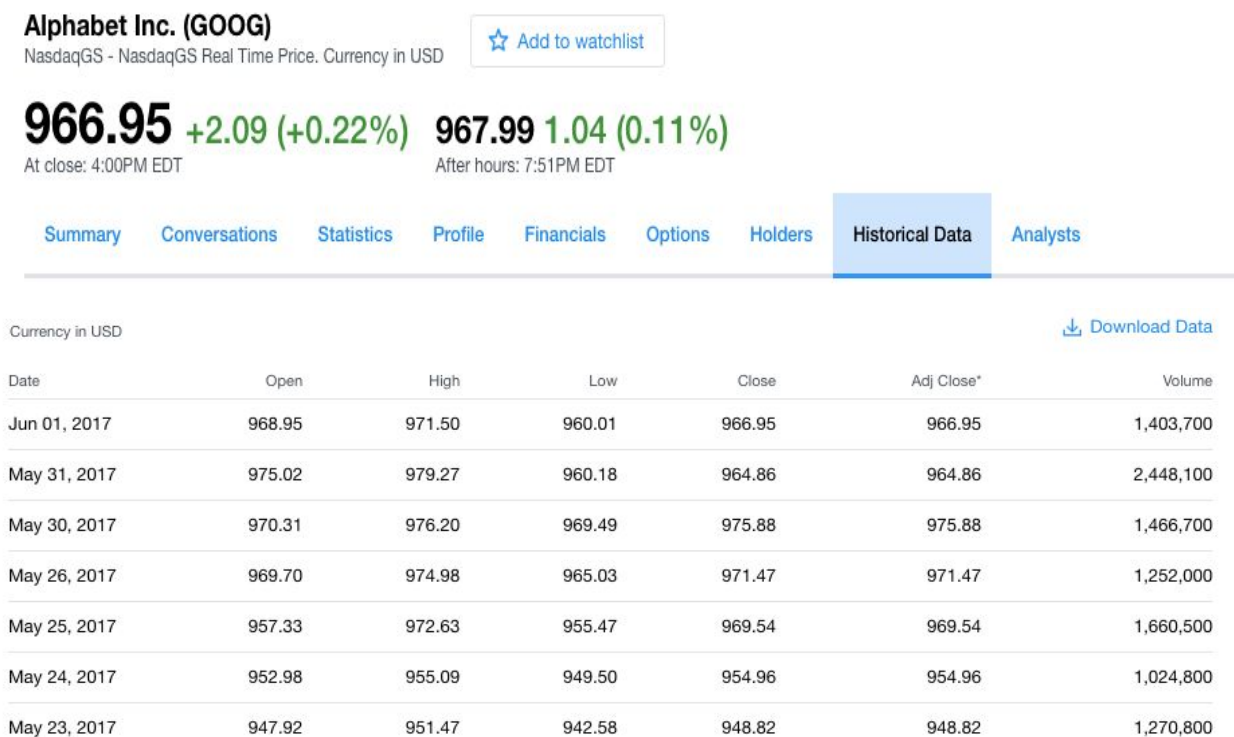


Figure 4.0

The Yahoo finance website was used to extract the historical datasets for the stocks symbols listed below from 2012-05-22 to 2017-05-22:

- GOOG - Google
- AAPL - Apple
- AMZN - Amazon
- MSFT - Microsoft

Each dataset includes 1,258 entries, and this shows each day that the market was open. The features needed to train the model are: 'Open', 'High', 'Low', 'Close', and 'Volume'. 'Adj Close' is the target variable, this is the value that the model is trying to predict. For the same reason, this feature(column) is dropped from the features dataset.

The original dataset is split into training and testing datasets. The size of the training dataset is 20% of the size of the original dataset.

- Size of the training dataset has 1,006 samples.
- Size of the testing dataset has 252 samples.

Exploration Visualization

To better understand the dataset, a scatter matrix is created (**Figure 5.0**). The graph shows each of the features present in the data. For the features 'Open', 'High', 'Low', and 'Close', the graph shows the prices for the stock against at each trading date. For the 'Volume' feature we can see the number of shares traded each date. Although it is hard to see at this point, the graph below slightly shows that there might be a correlation between "Low" and "High" features. It is also important to note the data distribution of the graph below. The graph shows the features that are presented are left-skewed and we can see this at the diagonal of the scatter matrix plot; this indicates the skewness of the data.

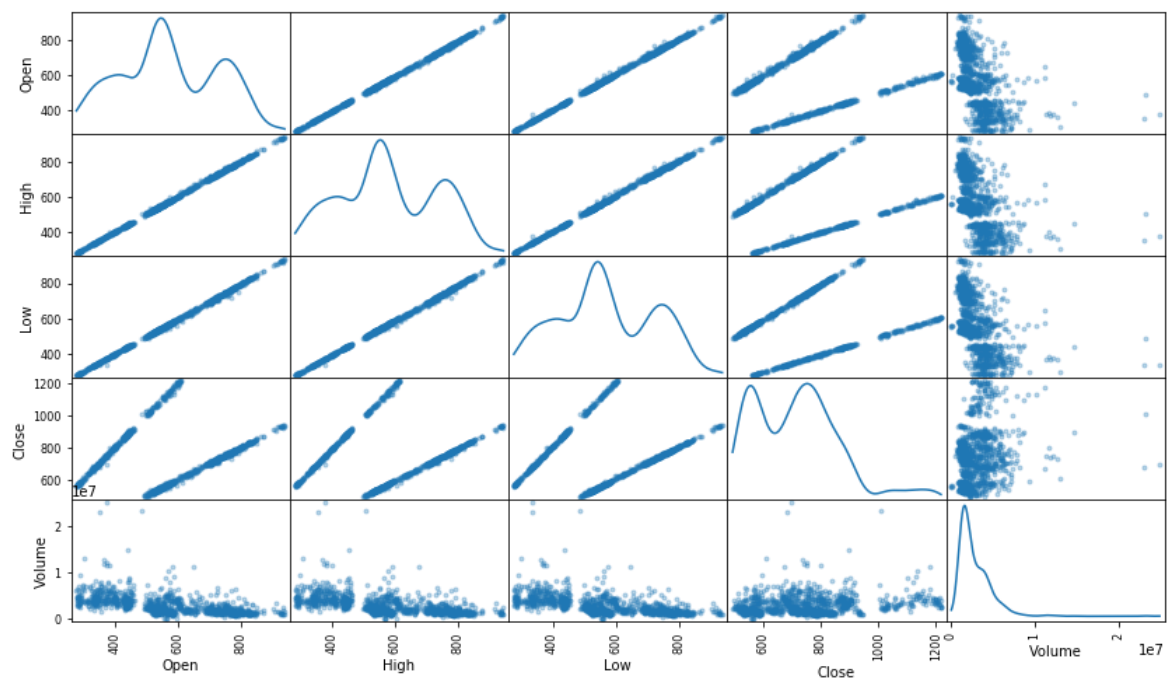


Figure 5.0

Algorithms and Techniques

Supervised machine learning algorithms will be used to create different models. In this problem the supervised machine learning model will try to predict stock prices, this is a good candidate for using a regression algorithm (Ericson, 2017). Supervised learning algorithms make predictions based on a set of examples. The historical stock prices will be used to predict future stock prices. A subset of data from the dataset will be used for training. Each entry is labeled with the value of 'Adj Close' stock price. The supervised learning algorithm will look for patterns in those value labels. After the algorithm has found the best pattern it can, it uses that pattern to make predictions for unlabeled testing data, for instance, future prices (Ericson, 2017).

Benchmark

From sklearn, we will use an out-of-the-box [DummyRegressor](#) for the benchmark. We will compare the R^2 scores between the DummyRegressor and the scores obtained from the selected algorithms (Linear Regression, KNeighborsRegressor, Epsilon-Support Vector Regression, and Lasso). The R^2 scores obtained from models created need to show that significantly outperform the DummyRegressor results so that we can assume that the model can be useful. Below is the implementation of the benchmark for this project:

```
from sklearn.dummy import DummyRegressor

def benchmark(X_train, X_test, y_train, y_test):

    print "*****DummyRegressor Model*****"

    model = DummyRegressor()

    model.fit(X_train, y_train)

    print '{}'.format(model.score(X_test, y_test))

    return model

benchmark(X_train, X_test, y_train, y_test)

*****DummyRegressor Model*****
-0.00434476331597
DummyRegressor(constant=None, quantile=None, strategy='mean')
```

III. Methodology

Data Preprocessing

Once the dataset was obtained from the Yahoo Finance website the following steps were performed to prepare the data before creating the models:

- In the dataset, we used the 'Date' feature as the index column for the dataset.
- Extracted the target variable ('Adj Close') from the dataset.
- Split the data into two different sets, training and testing datasets. The testing dataset consists of 20% of the size of the original dataset.
- Feature scaling of the data. A natural logarithm was applied. **Figure 6.0** shows the data after its preprocessing.

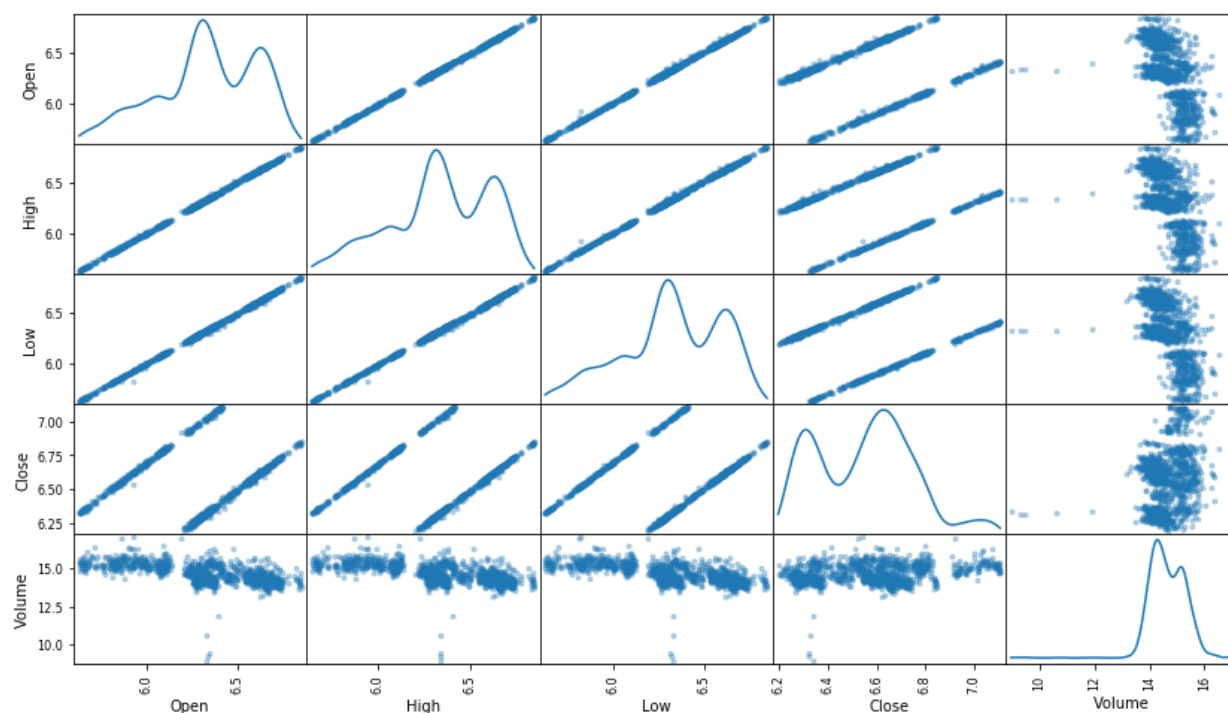


Figure 6.0

Implementation

The implementation of the stock predictor was performed executing the following steps:

1. The application receives the stock symbol of the stock as an input parameter.
The supported symbols correspond to the downloaded datasets(GOOG, AAPL, AMZN, and MSFT).
2. Split historical data collected into two datasets, training data, and test data.
3. Training the model using the historical data,
4. Apply different supervised learning algorithms.
5. Measure the score and the performance of each algorithm.
6. Tweak the model, if necessary, to get better scores results.
7. Select the algorithm that provides a better score.

Supervised Learning Algorithms

Different algorithms were applied to select the algorithm that provides a better score. The following supervised learning models that were used are available in [scikit-learn](https://scikit-learn.org/):

- Linear Regression.
- KNeighborsRegressor.
- Epsilon-Support Vector Regression.
- Lasso.

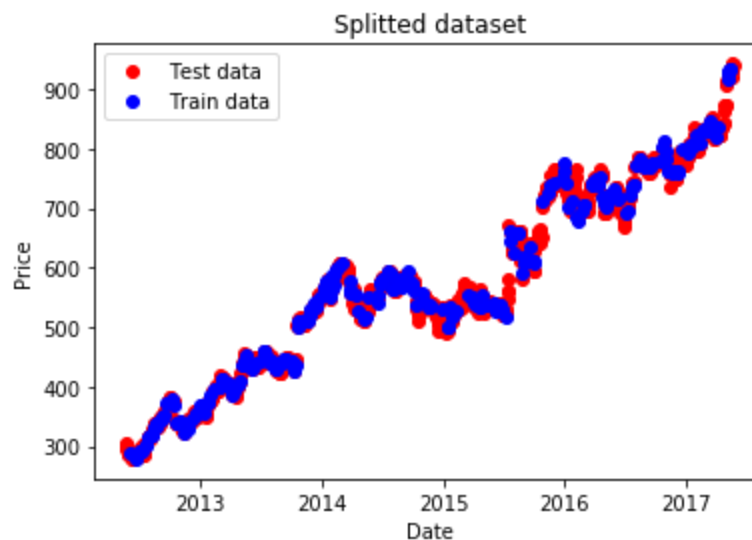


Figure 7.0

Figure 7.0 shows the data after it was split using [train_test_split](#). This function splits

arrays or matrices into random train and test subsets, it useful to prevent overfitting of the supervised model. The dataset was split in the following way:

- Training set has 1006 samples, which corresponds to approximately 80% of the size of the dataset.
- Testing set has 252 samples, this is approximately 20% of the size of the original dataset.

Each of the algorithms were trained with the same training dataset and also were evaluated against the same testing set. R^2 was used to measure the coefficient of determination. In R^2 the best possible score is 1.0 and it can be negative. The following table(**Figure 8.0**) describes the performance of the regressors and their corresponding scores.

Algorithm	R^2 Score	Trained model time.
Linear Regression.	0.999719108411	0.0016 seconds.
KNeighborsRegressor.	0.322343744915	0.0007 seconds.
SVR Regression.	-0.00388351020513	0.0659 seconds.
Lasso.	0.999310878941	0.0058 seconds.

Figure 8.0

It is easy to see that KNeighborsRegressor and SVR Regressor performed poorly as their R^2 is close to 0.0 or produces a negative value. A model that produces a negative R^2 indicates that the model is arbitrarily worse than one that always predicts the mean of the target variable (Scikit-learn, 2016). R^2 provides a measure of how well future samples are likely to be predicted by the model, a score closer to 1.0 is always preferable. Linear Regression and Lasso produced a better score. Therefore, these algorithms are the best candidates. This is described in detail in the refinement section.

Refinement

After the default model of Linear Regression and Lasso regressors obtained an R^2 score of **0.999719604742** and **0.999711366962** respectively. Both models were tuned using different parameters values of each regressor using [grid search](#). The grid search process of each regressor is discussed in detail in the subsequent sections.

Linear Regression Fine Tune

For the Linear Regression model, also known as ordinary least squares, a few parameters were considered to improve the **R² score** of the model (Scikit-learn, 2016). The following table (**Figure 9.0**) describe the parameters and the values that grid search identified.

Parameter	Description	Value
fit_intercept	Whether to calculate the intercept for this model. If set to false, no intercept will be used in calculations .	False.
normalize	If True, the regressors X will be normalized before regression. This parameter is ignored when fit_intercept is set to False. When the regressors are normalized, note that this makes the hyperparameters learned more robust and almost independent of the number of samples.	True.

Figure 9.0

After applying grid to the Linear Regression model the initial **R² score** was slightly improved from **0.999719108411** to **0.999719604742**.

Lasso Regression Fine Tune

Lasso is a linear model that estimates sparse coefficients (Scikit-learn, 2016). For this model, some parameters were considered to improve the **R² score** of the model. The following table (**Figure 10**) describe the parameters and the values that were applied to perform the grid search.

Parameter	Description	Value
fit_intercept	Whether to calculate the intercept for this model. If set to false, no intercept will be used in calculations.	True.
normalize	If True, the regressors X will be normalized before regression. This parameter is ignored when fit_intercept is set to False. When the regressors are normalized, note that this makes the hyperparameters learned more robust and almost independent of the number of samples.	True.
max_iter	The maximum number of iterations.	10000
alpha	Constant that multiplies the L1 term. Defaults to 1.0. alpha = 0 is equivalent to an ordinary least square, solved by the LinearRegression object.	0.1
selection	If set to 'random', a random coefficient is updated every iteration rather than looping over features sequentially by default.	random

Figure 10

After applying the grid search to the Lasso Regression model, the initial **R² score** was slightly improved from **0.999310878941** to **0.999723264224**.

IV. Results

Model Evaluation and Validation

This project explored further Linear Regression and Lasso Regression models because these models provided a far better **R² score** than the other models explored, KNeighborsRegressor and SVR Regression, even before that the tuning process was applied.

A grid search was used to tune Linear Regression and Lasso Regression models and obtained the better estimators; we can see a summary in the following table(**Figure 11**).

Model	Parameters	Final R ² Score	Time to Train
Linear Regression.	{'copy_X': True, 'normalize': True, 'n_jobs': 1, 'fit_intercept': False}	0.999719604742	0.9174 seconds
Lasso Regression.	{'normalize': False, 'warm_start': False, 'selection': 'random', 'fit_intercept': True, 'positive': False, 'max_iter': 10000, 'precompute': False, 'random_state': 42, 'tol': 0.0001, 'copy_X': True, 'alpha': 0.1}	0.999723264224.	7.9846 seconds.

Figure 11

The difference between the two different models is minimal. The Lasso Regression model **R² score** is slightly above the Linear Regressor model by **~0.000003%**. As the purpose of this project is to predict the stock price I believe that having a more precise model is a better idea. In this case, I decided to use the Lasso regression because it produces a **R² score** that it closer to 1.0. This model was tested using different stocks:

- GOOG - Google
- AAPL - Apple
- AMZN - Amazon
- MSFT - Microsoft

I believe that this model needs to be improved to be used to make real life trading decisions. For instance, similar stocks could be grouped together to create a model that would discover new features that can impact the stock prices for similar companies. Stock prices predictions are so complex that we need to build a more robust dataset to be able to build a model that it can predict the prices more accurate. Having a machine learning model like this is a good starting point, however, it needs improvement.

Justification

After tuning both models, Linear Regression and Lasso Regression, the **R² score** obtained significantly outperform the **R² score** obtained by the DummyRegressor.

Model	Parameters	R ² Score
Dummy Regressor	constant=None, quantile=None, strategy='mean'	-0.00434476331597
Linear Regression.	{'copy_X': True, 'normalize': True, 'n_jobs': 1, 'fit_intercept': False}	0.999719604742
Lasso Regression.	{'normalize': False, 'warm_start': False, 'selection': 'random', 'fit_intercept': True, 'positive': False, 'max_iter': 10000, 'precompute': False, 'random_state': 42, 'tol': 0.0001, 'copy_X': True, 'alpha': 0.1}	0.999723264224.

Figure 12

As we can see in the table(**Figure 12**) above, both models were able to outperform the DummyRegressor used for the benchmark for this project. However, while building this project, it was noted that to create a model that predicts the stock prices more accurately, it is necessary to obtain a more robust dataset to captures other important features that help to predict a stock price.

V. Conclusion

Free-form Visualization

After tuning both models, Linear Regression and Lasso Regression models, the **R^2 score** obtained was close to 1.0, which the possible score for the **R^2 score**.



Figure 13

We can see in the **Figure 13** that after approximately 100 training points the **R^2 score** were close to 1.0 . However, this high **R^2 score** does not mean that the model has a good fit. After obtaining these scores for both models, I believe that scores values were too high, in the Reflection section I will explore some possibilities that can be causing these high scores.

Reflection

The steps taken to create the model are the following:

- Explored different Supervised Machine Learning algorithms.
- Obtain the datasets from Yahoo finance for each of the following symbols:
 - a. [Google](#).
 - b. [Amazon](#).
 - c. [Microsoft](#).
 - d. [Apple](#).
- Explore the dataset using different plots such as scatter matrix plot.
- Dataset preprocessing:
 - a. Extract the target variable from the original dataset.
 - b. Perform feature scaling using natural logarithm.
- Split the dataset into training and testing datasets.
- Use and calculate the benchmark for the dataset, DummyRegressor from sklearn.
- Evaluate the following regressor models from sklearn.
 - a. Linear Regression.
 - b. KNeighborsRegressor.
 - c. Epsilon-Support Vector Regression.
 - d. Lasso.
- Calculate the **R^2 scores** and determine the top models.
 - a. Linear Regression.
 - b. Lasso.
- Tune the top models using Grid Search.
- Determine the best model base of the score results.
 - a. Lasso.

The following graph(**Figure 14**) represents the graph for the tuned Lasso Regressor.

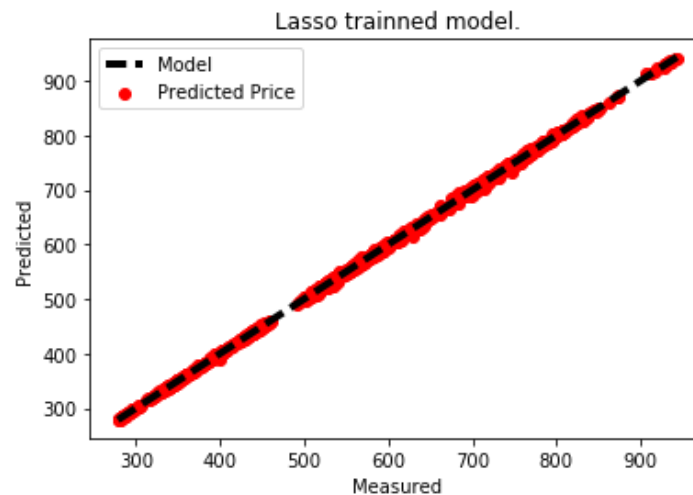


Figure 14

As we discuss earlier high **R² scores** doesn't necessarily mean that the model is more precise. Some factors that could be causing these high **R² scores** are:

- **Overfitted model:** An overfitted model is a model that is too complicated for a data set. Some reasons for these could be including too many terms in your model compared to the number of observations. When this happens, the regression model becomes tailored to fit the quirks and random noise in the sample dataset rather than reflecting the overall population. Overfitting occurs when model describes random error or noise instead of the underlying relationship (Multicollinearity, 2017).
- **Highly correlated data:** Correlation between the features can lead to overfitting. The best regression models are those in which the predictor variables each correlate highly with the dependent (outcome) variable but correlate at most only minimally with each other. Such a model is often called "low noise" and will be statistically robust (Multicollinearity, 2017).

The table below shows the correlation between the all the features and the target variable in the original dataset used to train the models described earlier.

	Open	High	Low	Close	Adj Close	Volume
Open	1.000000	0.999696	0.999512	0.150242	0.999182	-0.520535
High	0.999696	1.000000	0.999474	0.148506	0.999559	-0.516828
Low	0.999512	0.999474	1.000000	0.154026	0.999689	0.527687
Close	0.150242	0.148506	0.154026	1.000000	0.151542	0.229581
Adj Close	0.999182	0.999559	0.999689	0.151542	1.000000	-0.523261
Volume	-0.520535	-0.516828	-0.527687	0.229581	-0.523261	1.000000

Figure 15

By looking at the table(**Figure 15**) above we can see that the features “Open”, “High”, and “Low” are highly correlated. Since multicollinearity causes imprecise estimates of coefficient values in a regressor, the resulting predictions will also be imprecise. These are the reasons why I believe the current dataset needs to be improved.

Improvement

Stock price prediction has become more popular as more individuals have stocks as part of their overall equity. In the tech industry, it is common for companies to offer a stock equity as part of the compensation package for their employees. Finding the right time to sell or buy a certain stock to maximize the capital gains can be a hard problem to solve.

After the creation of the model for this project, it was clear to see that the original dataset needs to be extended to capture more features that can help to create a more realistic model. There can be a significant correlation between the changes in weekly stock prices based on important current events, news, and weather occurring around the world. It is important to review and consider a company's financial health, the value of a company's assets, debts, cash, revenues, expenses, profitability and plans of development.

References

- Dancer, D., & Tremayne, A. (2005). R-Squared and Prediction in Regression with Ordered Quantitative Response. *Journal Of Applied Statistics*, 32(5), 483-493.
- Ericson, G. & Franks, L. (2017). How to choose algorithms for Microsoft Azure Machine Learning. (2017). Retrieved June 3, 2017, from <https://docs.microsoft.com/en-us/azure/machine-learning/machine-learning-algorithm-choice>
- Hayes, Adam. (2017). Stocks Basics: What Are Stocks? Retrieved June 2, 2017 from <http://www.investopedia.com/university/stocks/stocks1.asp>
- Multicollinearity. (2017). In *Wikipedia, The Free Encyclopedia*. Retrieved June 5, 2017, from <https://en.wikipedia.org/w/index.php?title=Multicollinearity&oldid=776794468>
- Scikit-learn developers. (2016). Generalized Linear Models. Retrieved June 4, 2017, from http://scikit-learn.org/stable/modules/linear_model.html#lasso
- Scikit-learn developers. (2016). Model evaluation: quantifying the quality of predictions. Retrieved June 4, 2017, from http://scikit-learn.org/stable/modules/model_evaluation.html#r2-score
- Yahoo Finance. (2017). Alphabet Inc. Retrieved June 1, 2017, from <https://finance.yahoo.com/quote/GOOG/history?ltr=1>