

Predicting Future Stock Prices

Domain Background

Ever since the Stock Market started, people have been trying to make a living by predicting the future price of a stock. With this information, the buyer or seller can create strategies on which to make money so an accurate model can make a huge difference.

The Efficient Market Hypothesis is an investment theory that states it is impossible to “beat the market” because stock market efficiency causes existing share prices to always incorporate and reflect all relevant information. There are multiple forms of this theory which give tiers of efficiency and allow for some deviation but generally it states that we could not do better than holding a diversified portfolio of stocks.

Contrary to this argument lies the evidence of investment funds, banks, and traders who make a living buying and selling stocks according to different financial models. Some of these models are based on technical analysis of the charts while others adhere to fundamental principles of company valuation.

In this project, I will build a stock price predictor which can predict a stock’s future price to profit from trends in the data.

Problem Statement

The goal is to build a stock price predictor that takes daily trading data over a certain date range as input, and outputs projected estimates for given dates. Although there are many stock prices in a day, the model will only need to predict the adjusted closing price.

To generate future stock prices the model will complete the following tasks:

- I. Build a data set or query data from a financial API
- II. Implement a user interface which allows the user to choose which stock to predict and how far into the future to predict
- III. Select relevant features for the estimation model
- IV. Train a prediction model over the data set
- V. Plot the predictions for the user to see and gain information for stock selection
- VI. Give the user information concerning the reliability of the model by scoring the model and comparing it to a benchmark model

Datasets and Inputs

The dataset will be built by querying from the Quandl financial API to select price information. After watching a series of videos from pythonprogramming.net on financial data analysis, I will be able to repurpose an existing querying script to pull in data over a list of stock tickers and date range specified. For this project, I will be pulling historical price data for all stocks currently in the S&P500. The time frame will be from 2000-01-01 to 2016-12-31 which amounts to 4,032 trading days so the data set will

contain 4,032 rows. The columns of the dataset will be the Adj. Close of each stock so that there will be a total of 500 columns.

After going through the Machine Learning for Finance Udacity videos, 3 calculated fields of moving averages will also be added into the data set. These will be the 10, 20, and 30 day moving averages for the stock being predicted.

The [Date] field will be transformed from a date variable into a [Day_N] field which will essentially be an index variable which

The inputs of the script will be the {ticker} which will tell the program which stock to predict, the number of days into the future to predict, and if possible, the date range over which to train the model.

The target variable will be the stock price column designated by the user from the {ticker} variable and will be shifted up by the number of days the user wants to predict.

Solution Statement

Because this course did not cover time-series data forecasting, and because I feel more familiar with standard estimation techniques, this problem will be approached using linear models and other estimators covered in this course such as the SVM Regressor, KNR, and Random Forest Regressor. I will be using so many different types of models so that an ensemble model can be built at the end of the process. Ensemble modeling is very helpful when considering prediction problems because it can help counteract any biases within each individual model. All models will be judged according to their score and against the benchmark model in order to give the user useful information.

Benchmark Model

The benchmark model for this problem will be a univariate linear regression with the current price of the stock as the single variable to predict the future price of the stock. This was chosen because it is the simplest mode I could think of that could also give some use to the user. This would baseline model would be able to find the general trend of the stock's current price with the future price.

Evaluation Metric

The learned model will be judged against this benchmark using the R-squared metric as this is the most common metric to score regression models. This metric is already the base metric in sklearn for scoring the models so the ".score(X,y)" method will be implemented.

Because this program will have the capability to predict any stock price any number of days into the future, the score of the model will greatly depend on the inputs given by the user. For this reason, I cannot give an acceptable desired score.

Project Design

After the dataset is assembled, I will need to clean the data of any null values. This final set will be broken into training, testing, and a future set based on an 85 – 15% split rate. The future set will contain N number of rows depending on what the user defines as his future price to predict. This is due to the shifting of the data to produce the target variable. These last rows will not allow for verification so they will be kept by themselves as a "future" prediction set.

The dataset will then be passed through multiple models to select the best features. The first model chosen for this process will be PCA because it can find new dimensions that a human could not alone find. Second will be Recursive Feature Elimination, because it will be able to cycle through the given variables and find the best field while excluding unnecessary ones. These two datasets will double the number of predictions possible which will hopefully add to our end ensemble model.

The general learning process will go as follows:

1. A function for training a specific model is defined
2. A dataset is passed to the training function
3. TimeSeriesSplit() splits the data into 10 different subsets on which to perform cross-validation
4. The learning algorithm is defined
5. A dictionary of relevant parameters and their values is defined
6. The regressor, parameters, and cross-validation set are passed to a GridSearchCV() function
7. The training and target sets are passed into the grid
8. The best estimator is returned

The above process will be applied to several regression estimators thus creating multiple predictions over which to ensemble of which two different processes will be tried. The first ensemble method will be another linear regression that will take the predictions of the models as inputs and try to predict the target price. The second way will be a simple average of the models' predictions.

Once all models are built and optimized, scores will be taken for each one and compared to the benchmark. These scores will be presented to the user along with graphs of the training, testing, and future predictions.

Below is the planned outcome of the script which will be presented to the user upon completion of all training and testing of the models.

