

Brief Overview:

Name: Tony Jiang

I have been very interested in the concept of automated stock trading for a long time but had never really had any chance to do any project on it. After learning from the past couple lessons about different models of machine learning, I thought maybe I can apply my knowledge and practice into something I am interested in. Thus, I decided to use the Nasdaq datasets on yahoo finance to test whether we can predict statistical values of a trading date based on the statistical datas of the previous trading date. In the dataset (Nasdaq.csv), each row of the file contains several statistical datas (such as Open, High, Close, Volumn,etc) of a specific trading data. I also manually pasted the same statistical datas from the previous trading data on to the same row. The entire dataset spans for the past last 5 years, and I arbitrarily used the last 25 dates as the testing set and rest of the dataset as the training set.

Independent and Dependent Variables:

There are 6 features (independent variables) that the decision trees to take into consideration.

They are all statistical values of the previous trading date: *continuous value

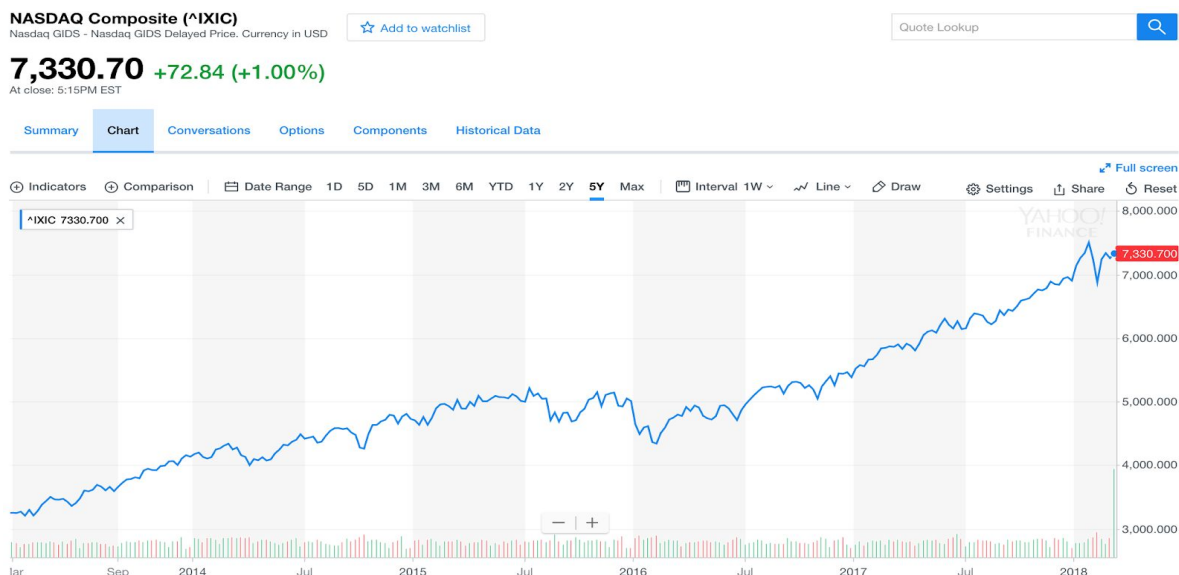
```
feature_names = ['Yesterday_Open', 'Yesterday_High', 'Yesterday_Low',  
'Yesterday_Close', 'Yesterday_AdjClose','Yesterday_Volumne']
```

The dependent variable could be any statistical value of the current trading date. For my tests, I decided to use the High value of the current trading date:

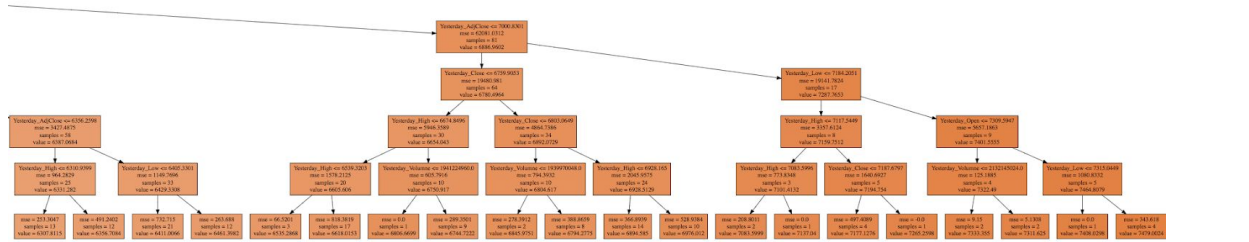
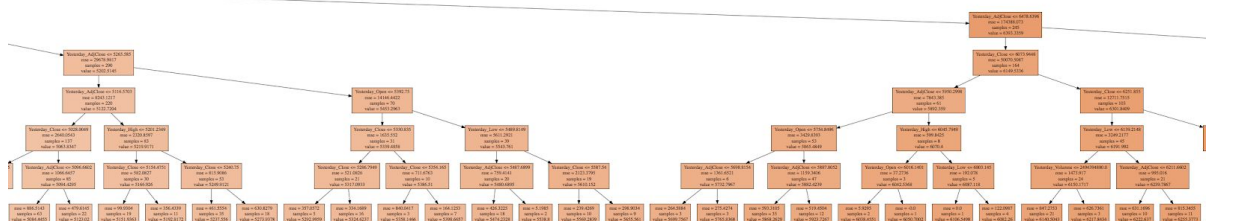
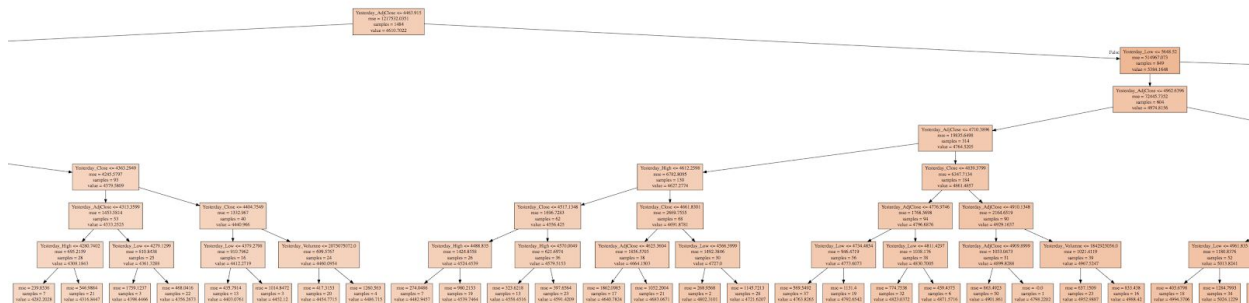
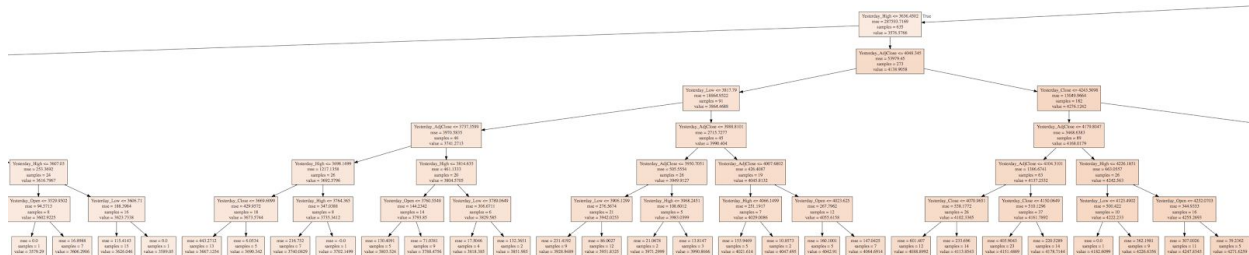
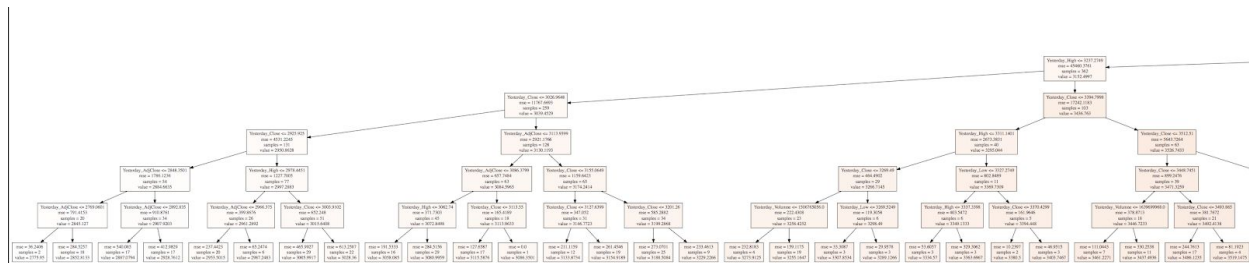
```
target_names = ['Today_High'] *continuous value
```

Generic Data Graphics:

The following graph illustrates the Nasdaq closing value of all trading dates for the past 5 years



Decision Tree Image and Description:



Here are the interpretation of the first 2 layers of the regression tree:

The first layer box states that :

Yesterday_AdjClose <= 4463.915; mse = 1217532.0351; samples = 1484; value = 4610.7022

This means that out of all the 1484 data values for high values, the average is 4610.7022. The tree will further break down into two boxes in the second layers. The first box represents all of the data that have Yesterday AdjClose Value less than or equal to 4463.915. The second box, similarly, represents all of the data that have High values greater than 4463.915.

On the second layer, the left box states that:

Yesterday_High <= 3636.4502 mse = 287593.7169 samples = 635 value = 3576.5766

This means that out of 635 data here, the average High values are 3576.5766 and it would further branch out base on the condition Yesterday High <= 3636.4502.

The right box on the second layer states that:

Yesterday_Low <= 5648.52 mse = 514967.073 samples = 849 value = 5384.1648

Here, it means that out of the 849 data, the average of their high values are 5384.1648, and it would further branch out base on the condition: Yesterday Low value <= 5648.52

DT and RF feature importances:

1. The following is the Decision Tree feature importances of the 6 features:

[4.92583996e-04 9.09353968e-02 1.95303670e-01 8.95871836e-03
7.04285977e-01 2.36542763e-05]

Order: ['Yesterday_Open', 'Yesterday_High', 'Yesterday_Low', 'Yesterday_Close', 'Yesterday_AdjClose', 'Yesterday_Volume']

2. The following is the Random Forest feature importances of the 6 features:

[3.09105237e-02 3.02474340e-01 5.84258076e-02 3.13251016e-01,
2.94897595e-01 4.07169488e-05]

Order: ['Yesterday_Open', 'Yesterday_High', 'Yesterday_Low', 'Yesterday_Close', 'Yesterday_AdjClose', 'Yesterday_Volume']

Final Thoughts:

Overall, I think it's a successful mini project. For both decision tree and randomized forest, I am able to get CV Values in the 90s with the optimal max_depth and number_estimator(Max_depth = 7 for both and n_estimator = 200). Nevertheless, as you can see from the result (printed in the program), though my predicted value is mostly in the range + - 100 from the actual value, the High values fluctuates just as much every day in the dataset. So, I can't really use my predicted value to accurately predict whether the market would increase or decrease based from yesterday's values. Nevertheless, It gives us a range of difference in values from one trading day to another, and validates the strength of machine learning. I am quite pleased with the overall result.