

## Machine Learning Engineer Nanodegree

### Capstone Project

# Investment and Trading : Build a Stock Price Predictor<sup>1</sup>

Soyoung Park

March 31, 2017

## I. Definition

### Project Overview

Technical stock price analysis can be approached as a typical time series data analysis without understanding of health of a particular security while fundamental analysis requires deep domain knowledge about a security or domain of interest.

In this project, one particular security, Intel stock, was chosen for forecasting the future stock trend using various machine learning techniques without having to have domain knowledges. The input price data is from Yahoo! Finance. The baseline model was simple time series forecasting model (Autoregression Integrated Moving Average, ARIMA). Two different approaches were attempted for forecasting the trend.

First, the price or return of 7,14,28 days were predicted with various regression models and the performance of each model was ranked using root mean square of errors (RMSE)

Second, to understand how classification can be applied for this forecasting the continuous variable, the daily return was transformed into binary information as increase or decrease, and it was predicted with various classification models. Trading decision (Buy/Sell) was made for the days of query upon prediction. At the end, the total capital growth as a result <sup>2</sup> of trading upon decisions was backtested.

### Problem Statement

Many day stock traders use their own various strategies to find patterns in one particular stock price chart and make their bets without having to understand the health of business or industry because they believe all information is already baked in the stock price. <sup>3,4</sup> The biggest

---

<sup>1</sup> The given problem description of the project:

<https://docs.google.com/document/d/1ycGeb1QYKATG6jvz74SAMqxrlek9Ed4RYrzWNhWS-0Q/pub>

<sup>2</sup> Complete code for this project is

at: [https://github.com/parksoy/Udacity\\_nanoDegree\\_MachineLearning/blob/master/capstone/CapstoneProject\\_ML4Trading\\_SoyoungPark.ipynb](https://github.com/parksoy/Udacity_nanoDegree_MachineLearning/blob/master/capstone/CapstoneProject_ML4Trading_SoyoungPark.ipynb)

<sup>3</sup> "How to Become a Day Trader on the Stock Market | Investopedia." 8 Mar. 2017,

<http://www.investopedia.com/trading/how-to-become-day-trader/>.

<sup>4</sup> "8 Reasons Why You Should Never Become A Day Trader - Business ...." 8 Nov. 2010,

<http://www.businessinsider.com/reasons-not-to-be-a-day-trader-2010-11>.

problem with manual searching for trade signals in price chart is not only risky due to lack of fundamental analysis or statistical significance, but also, takes significant training efforts to build experience to make a consistently reasonable performance. If technical analysis relying on price chart only is believed as valuable, machine learning can be very good alternative solution to guide the day traders to make their trading decisions on each signal. In this project, the problem was approached by following: One company's historical chart in the public domain is taken and various machine learning techniques (such as regression or classification) were applied to understand the trend of price changes more systematically along with the accuracy of prediction for the future. <sup>5</sup> Ultimately, the predicted trading decision can be back tested to forecast how the asset grows over time, and verify if there is any risk so day traders can avoid guessing and relying on luck on their investments.

## Metrics

For ARIMA baseline study and regression approach with various models, RMSE between the prediction and the actual value in test set was used as an evaluation metric. Precisely, RMSE between the predicted log return of 14 days ( $y_{i,predict,testset}$ ) and the actual log return of 14 day value ( $y_{i,actual,testset}$ ) of testset was recorded for ARIMA model, and each RMSE of price, return, log return of 7,14,28 days were recorded for various regression models as shown in eq.(1).

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (y_{i,predict,testset} - y_{i,actual,testset})^2}{n}} \quad \text{eq. (1)}$$

$$f1 = \frac{2 \cdot \text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} \quad \text{eq. (2)}$$

$$\text{where } \text{precision} = \frac{Tp}{Tp+FP}, \text{ recall} = \frac{Tp}{Tp+Fn}$$

The reason RMSE is used here because the quality of model or predictability power of the model should be compared between many different models. RMSE more aggressively punishes big errors than small ones, as a model with the least error is preferred for stock price/return prediction purpose. <sup>6</sup>

For classification of increase/decrease of daily return, its f1 score was used as shown in eq.(2). The INTC data was balanced as the number of bullish days and bearish days were equivalent, however, other company's stock or other test period of INTC stock can be imbalanced. To make the predictor applicable in general, f1 score was used to capture any

<sup>5</sup> "Forecasting stock market movement direction with support vector ...."

<https://pdfs.semanticscholar.org/ed03/6a6f69d192c98a750e8b937061eecf1aba50.pdf>.

<sup>6</sup>

<https://www.quora.com/Why-we-use-Root-mean-square-error-RMSE-Mean-absolute-and-mean-absolute-percent-errors-for-forecasting-time-series-models>

imbalanced case by incorporating precision and recall. Simple accuracy or hit rate may not tell performance of predictor in valid manner for imbalanced cases.

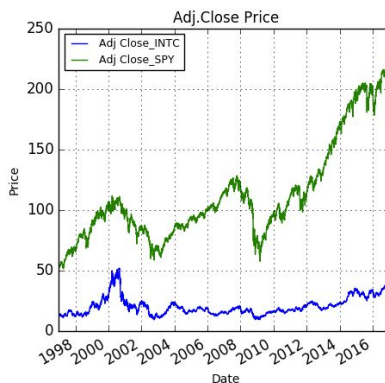
## II. Analysis

### Data Exploration

Input data is Intel (Ticker symbol, INTC) and S&P 500 (SPY) stock price data from January 1, 1997 to March 9, 2017 in csv file from Yahoo! Finance. Total 5080 datapoints for each INTC, SPY were incorporated. Python and its various libraries were used to process the data for this project. User can query the following inputs

- symbol: of interested company
- start\_date = datetime.datetime(1997,1,1)
- end\_date = datetime.datetime(2017,3,9)

A function, get\_data() will download two csv files, one for INTC and other for SPY. Two files are called and combined into a DataFrame and Figure 1 and Table 1 show the first two days of dataset as an example.



**Figure 1.**

	Open_SPY	High_SPY	Low_SPY	Close_SPY	Volume_SPY	Adj Close_SPY	Open_INTC	High_INTC	Low_INTC	Close_INTC	Volume_INTC	Adj Close_INTC
1997-01-02	74.375	74.375	72.750000	74.031197	2031900.0	51.586402	131.75	132.0	127.625	130.375	97639200	11.206980
1997-01-03	74.375	75.125	74.078102	75.093697	2123200.0	52.326773	133.00	138.5	132.625	138.375	95648000	11.894656

**Figure 1. Table 1.** Input dataset that was used for the entire project

Additional features and targets were devised using adjusted close price, 'Adj Close\_INTC':

- 'Rolling\_mean': rolling mean of adjusted close price for 20 days

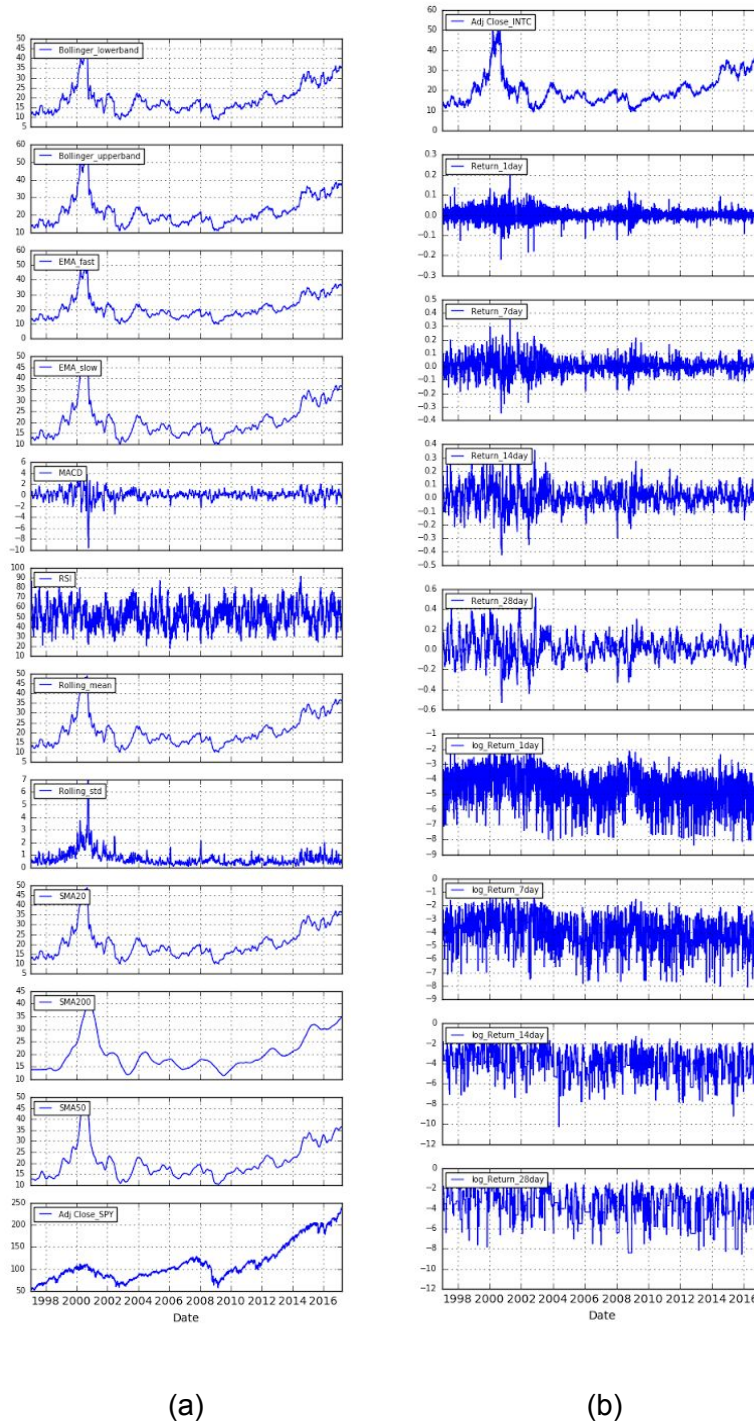
- 'Rolling\_std': rolling standard deviation of 20 days
- 'Bollinger\_lowerband':  $\text{rolling\_mean} - 2 * \text{rolling\_std}$
- 'Bollinger\_upperband':  $\text{rolling\_mean} + 2 * \text{rolling\_std}$
- 'EMA\_fast': exponential moving average of 12 days
- 'EMA\_slow': exponential moving average of 26 days
- 'MACD': Moving average convergence divergence.  $\text{emafast} - \text{emaslow}$
- 'RSI': 14 days relative strength indicator  $100 - 100 / (1 + \text{rs})$  where rs is relative strength, number of price up divided by number of price down
- 'SMA20', 'SMA50', 'SMA200': simple moving average of 20, 50, 200 days
- 'Adj Close\_SPY': Adjusted Close price of S&P 500 market.

Additional target indicators were created on top of simple adjusted close price, 'Adj Close\_INTC':

- 'Return\_1,7,14,28 day'
- 'log\_Return\_1,7,14,28 day'

All input data is continuous numeric variables, there was no abnormalities in terms of missing data or outlier data. However, when the additional features or targets were driven from the Adjusted Close\_INTC data, missing data (NaN) were treated with backfill for the cases that when return or log return was calculated to avoid any unintentional information is interfered by filling with zero or mean. Dropping NaN was avoided to keep a good number of datapoints. Any issue with taking log of negative or zero return was resolved with backfill fill methodology.

## Exploratory Visualization



**Figure 2.** (a) 12 technical indicators (b) 9 different target indicators in time series before normalization.

Figure 2 shows (a) features and (b) targets that were used in this project. X axis is in time and y axis is annotated as in legend box. Price is in US dollar and any ratio, return, or log are unit less. In feature plots, time trend of various moving averages or adjusted close price are

similar over the course of years. Ratio of two different moving average, such as MACD, or RSI metrics are, however, different trend from various moving average itself.

## Algorithms and Techniques

### A.ARIMA

As a baseline of the project, simple forecasting the price out of the sample with its own history was attempted with ARIMA model with guidance of the reference.<sup>7</sup> The following is taken a simple description of ARIMA (an autoregressive integrated moving average) from Wikipedia<sup>8</sup>:

ARIMA is fitted to time series data either to better understand the data or to predict future points in the series (forecasting). ARIMA models are applied in some cases where data show evidence of non-stationarity, where an initial differencing step can be applied one or more times to eliminate the non-stationarity.

- AR : evolving variable of interest is regressed on its own lagged (i.e., prior) values.
- I : data values have been replaced with the difference between their values and the previous values (and this differencing process may have been performed more than once). The purpose of each of these features is to make the model fit the data as well as possible.
- MA : regression error is actually a linear combination of error terms whose values occurred contemporaneously and at various times in the past.

Non-seasonal ARIMA models are generally denoted  $ARIMA(p,d,q)$

- where parameters  $p$ ,  $d$ , and  $q$  are non-negative integers,
- $p$  is the order (number of time lags) of the autoregressive model,
- $d$  is the degree of differencing (the number of times the data have had past values subtracted),
- $q$  is the order of the moving-average model.

The INTC price data itself was non-stationary, transforming the price into log return was necessary. 14 day window were chosen to avoid any possible noise by too frequent sampling. After taking log return of 14 days of adjusted close price, the data became stationary, forecast is made one day ahead (one step) out of history as that point is added to history as progresses.

$ARIMA(\text{history}, \text{order}=(1,1,0))$  model was precisely used on 2046 data history points. This means,

- Number of time lag=1
- Degree of differencing=1
- Order of moving average model=0

As model is fit one step, data is added from test set to history set. Starting testset points were 681 points after history points. Fine hyperparameter tunings were not proceeded.

<sup>7</sup> "How to Create an ARIMA Model for Time Series Forecasting with Python." 9 Jan. 2017, <http://machinelearningmastery.com/arima-for-time-series-forecasting-with-python/>.

<sup>8</sup> [https://en.wikipedia.org/wiki/Autoregressive\\_integrated\\_moving\\_average](https://en.wikipedia.org/wiki/Autoregressive_integrated_moving_average)



## B. Regression

Various regressors were used to model price or return 7,14,28 days, log return of 7,14,28 days. The dataset consisted of 12 features, and 9 different target metrics. The dataset were splitted to trainset (first 2 out of 3 splits in time series) and testset (last one out of 3 splits). Using 'for loop' on 9 different targets, each model was fit on '12 features-one target' pair in train set and 12 features were input to predict the target on testset. RMSE was recorded on each model per each target and ranked with respect to baseline ARIMA model's performance. The following 4 regressors were attempted and its short descriptions are taken from each reference:

- **KNeighborsRegressor(n\_neighbors=3):**
  - Regression based on k-nearest neighbors. The target is predicted by local interpolation of the targets associated of the nearest neighbors in the training set. In this project, n\_neighbors=3 was used.<sup>9</sup>
- **RandomForestRegressor:**
  - A random forest is a meta estimator that fits a number of classifying decision trees on various sub-samples of the dataset and use averaging to improve the predictive accuracy and control over-fitting. The sub-sample size is always the same as the original input sample size but the samples are drawn with replacement if bootstrap=True. In this project, bootstrap was set as True.<sup>10</sup>
- **DecisionTreeRegressor(max\_depth=4):**
  - Decision tree learning uses a decision tree as a predictive model which maps observations about an item (represented in the branches) to conclusions about the item's target value (represented in the leaves). Decision trees where the target variable can take continuous values are called regression trees.<sup>11,12</sup>
- **AdaBoostRegressor (DecisionTreeRegressor(max\_depth=4),n\_estimators=300):**
  - An AdaBoost regressor is a meta-estimator that begins by fitting a regressor on the original dataset and then fits additional copies of the regressor on the same dataset but where the weights of instances are adjusted according to the error of the current prediction. As such, subsequent regressors focus more on difficult cases.<sup>13</sup>
- **SVR:<sup>14</sup>**
  - Support vector machine regression uses kernels to find hyperplane that maximizes margin.<sup>15</sup>
  - Default, C=1, epsilon=0.1, kernel=rbf (radial basis) were used in this project.

<sup>9</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsRegressor.html>

<sup>10</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html>

<sup>11</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeRegressor.html>

<sup>12</sup> [https://en.wikipedia.org/wiki/Decision\\_tree\\_learning](https://en.wikipedia.org/wiki/Decision_tree_learning)

<sup>13</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostRegressor.html>

<sup>14</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVR.html>

<sup>15</sup> <http://kernelsvm.tripod.com/>

Fine hyperparameter tunings were not proceeded.

### C.Classification

For target metric, daily return was converted to 1(bullish or no change) or 0(bearish) for each day, and 12 features were used to fit and predict those binary target metric. Grid search cross validation was used on 3 timely splitted datasets for the following classifiers with various hyperparameter configuration in the nested for loop:

- AdaBoostClassifier:
  - Similar learning algorithm to its regressor that was mentioned previously.<sup>16</sup>
- GradientBoostingClassifier:
  - GradientBoostingClassifier builds an additive model in a forward stage-wise fashion. It allows for the optimization of arbitrary differentiable loss functions. In each stage, `n_classes_` regression trees are fit on the negative gradient of the binomial or multinomial deviance loss function. Binary classification is a special case where only a single regression tree is induced.<sup>17</sup>
- XGBClassifier<sup>18</sup>:
  - Short for Extreme Gradient Boosting.
- LogisticRegression<sup>19</sup>:
  - A regression model where the dependent variable is categorical, "0" and "1".<sup>20</sup>
- RandomForestClassifier<sup>21</sup>:
  - Similar learning to its regressor as mentioned earlier.
- KNeighborsClassifier<sup>22</sup>:
  - Similar learning to its regressor as mentioned earlier.
- SVC<sup>23</sup>:
  - Similar learning to its regressor as mentioned earlier.

F1 score was recorded for all cases, and the best model with the best hyperparameter set was selected. At the end, trade decision was made with the best prediction of everyday bullish, bearish information. Buy signal is made when two consecutive bullish days are predicted. Sell signal is recorded when two consecutive bearish days are predicted. Simple buy-sell-hold exercise based on prediction was backtested to understand how portfolio grows over the course of investment period.

---

<sup>16</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.AdaBoostClassifier.html>

<sup>17</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.GradientBoostingClassifier.html>

<sup>18</sup> <http://xgboost.readthedocs.io/en/latest/model.html>

<sup>19</sup> [http://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LogisticRegression.html](http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LogisticRegression.html)

<sup>20</sup> [https://en.wikipedia.org/wiki/Logistic\\_regression](https://en.wikipedia.org/wiki/Logistic_regression)

<sup>21</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html>

<sup>22</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.neighbors.KNeighborsClassifier.html>

<sup>23</sup> <http://scikit-learn.org/stable/modules/generated/sklearn.svm.SVC.html>



## Benchmark

One step forecasting out of sample using ARIMA model time was a benchmark/baseline for this project. The evaluation metric is RMSE to compare with other regression models.

RMSE of forecast and actual value(log return of 14 days) in testset was 0.86 as shown in Figure 4(a). Figure 8 shows RMSE of all regressors in rank. ARIMA is marked in red and other models performed better than ARIMA. Figure 4(b) shows superimposed price of trainset and testset, especially out of sample log return of 14 days were transformed back to price after the first date in testset that the black arrow indicates.

## III. Methodology

### Data Preprocessing

As two (INTC, SPY) csv files were downloaded from Yahoo! Finance, they were read and merged into one DataFrame. Total 5080 data points in 14 different price metrics were starting points. There was no missing values, however, when adjusted close price was transformed into return of X days by shifting X rows (where X=1,7,14,28 days), the return before X days were missing(NaN) and they were backfilled with the first value of return. When log of return was taken, any zero or negative return became infinite. In this case, 'infinite' numbers were simply dropped in ARIMA modeling since there were enough number of datapoints (initially 2046 points in history + upto additional 681 points) to forecast one step as one data point is added to history set while training. However, for the case of regression task, the author preferred to keep the same number of datapoints for all targets, (price, return of 7,14,28 day, and log return of 7,14,28 day) to understand if any target is better metric than other or not by comparing RMSE in valid manner. So any 'infinite's before X days were backfilled to avoid any further projection of additional information or losing too much of datapoints.

For classification task, the daily return was converted to binary format as 1 for increase/no change, 0 for decrease of price with respect to the previous day.

### Implementation

#### A. ARIMA - Baseline

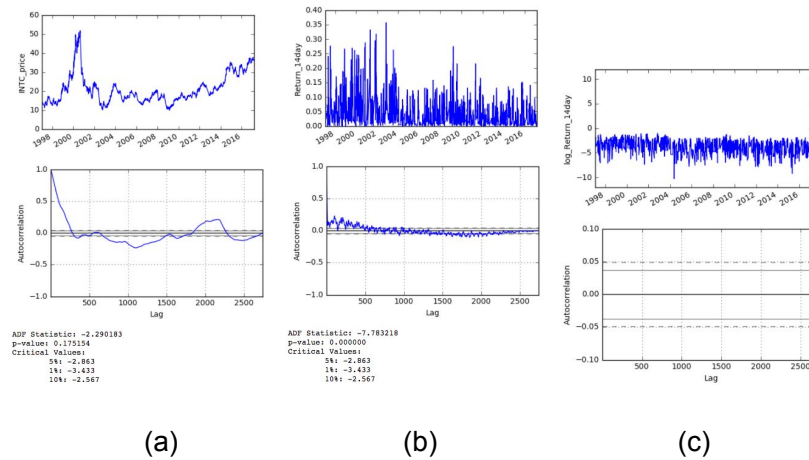
Based on this time series forecasting tutorial,<sup>24,25</sup> stationary test value on price data over the period of time was -2.290183, which is greater than all of the critical values at the 1%, 5%, and 10% confidence levels. This "non-stationary" data needed to be converted to stationary data to make a valid forecasting. As suggested in the tutorial, return of 14 days instead of price was calculated. Stationary test value on return became smaller than any of confidence level as shown in (b) compared to Figure 3(a). Furthermore, log was taken on return of 14 days to

<sup>24</sup> "How to Check if Time Series Data is Stationary with Python - Machine ...." 30 Dec. 2016, <http://machinelearningmastery.com/time-series-data-stationary-python/>. Accessed 31 Mar. 2017.

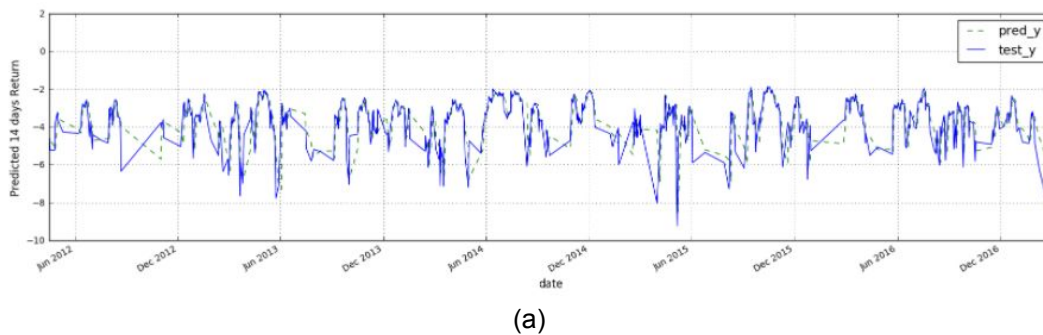
<sup>25</sup> "How to Create an ARIMA Model for Time Series Forecasting with Python." 9 Jan. 2017, <http://machinelearningmastery.com/arima-for-time-series-forecasting-with-python/>. Accessed 31 Mar. 2017.

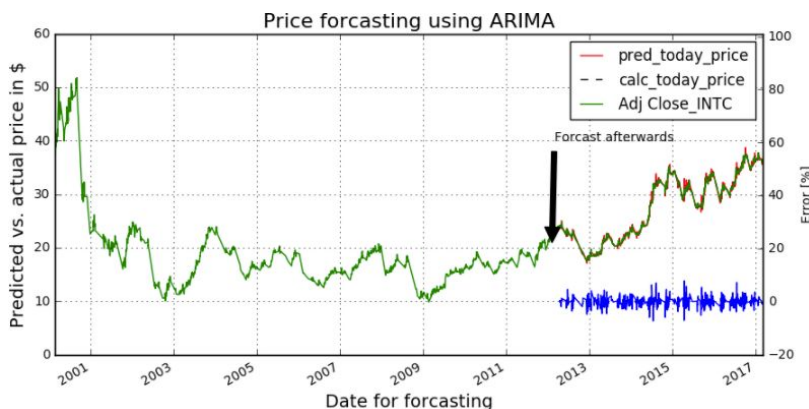
completely remove any autocorrelation. Log of return 14 days in Figure 3(c) suggests the data is stationary and good to use for valid forecasting as shown.

Taking negative or zero return led to infinite value, and they were dropped from the dataset. Final data points were time series splitted and first 2046 points were used as history in trainset, the last 681 points were used for one step forecasting in testset. As one step forecast is made, that datapoint is added to history. RMSE of forecast and actual value(log return of 14 days) in testset was 0.86 as shown in Figure 4(a). Figure 4(b) shows superimposed price of trainset and testset, especially out of sample log return of 14 days were transformed back to price after the first date in testset that the black arrow indicates.



**Figure 3.** (a) non-stationary (b) stationary time series after taking return of 14 days





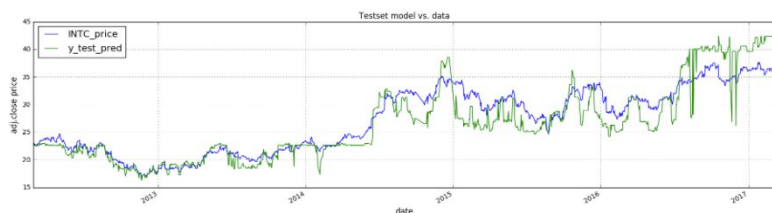
(b)

**Figure 4.** (a) ARIMA forecasting result on return of 14 days (b) Comparison of the price in and out of the sample: the price was transformed back from the forecasted return out of sample.

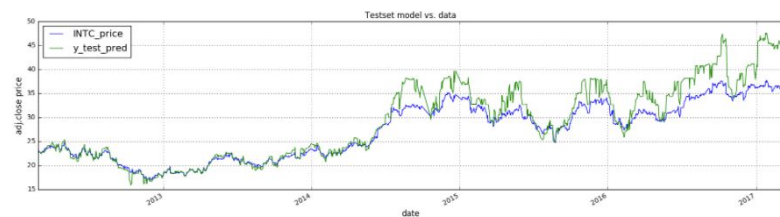
## B. Regression

Various regressors were used to train and predict using the following 12 features and the superimpose price of predicted value and actual values for each regressor is shown in Figure 5(a)-(e).

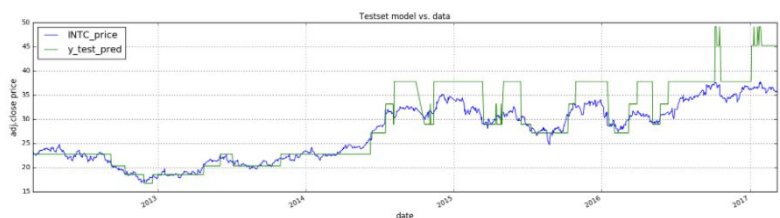
- 'Bollinger\_lowerband',
- 'Bollinger\_upperband',
- 'EMA\_fast',
- 'EMA\_slow',
- 'MACD',
- 'RSI',
- 'Rolling\_mean',
- 'Rolling\_std',
- 'SMA20',
- 'SMA200',
- 'SMA50',
- 'SPY\_price'



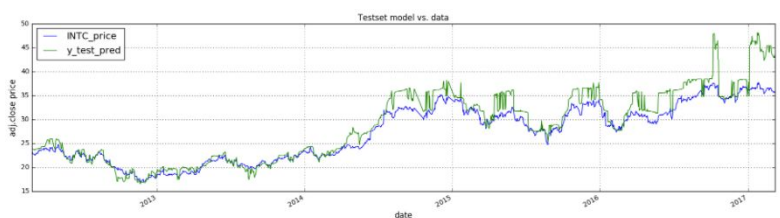
(a) knn



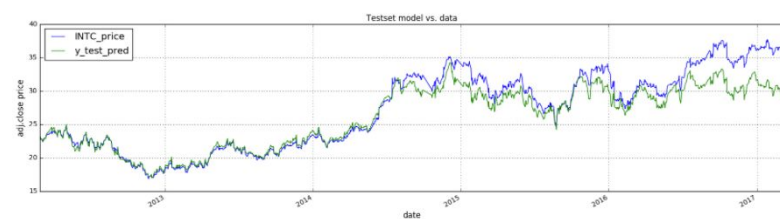
(b) RandomForest



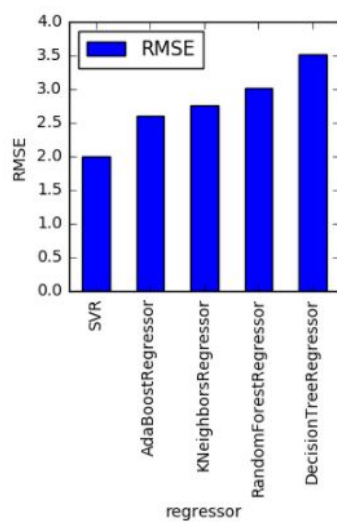
(c) Decision tree



(d) Adaboost with decision tree



(e) SVR regressor



(f) Summary of RMSE in prediction of price in order for various regressors with default args

**Figure 5.** Prediction result of various regressors (a) Knn (b) RandomForest (c) Decision tree (d) Adaboost with decision tree (e) SVR regressor (f) Summary of RMSE in prediction of price in order for various regressors with default args

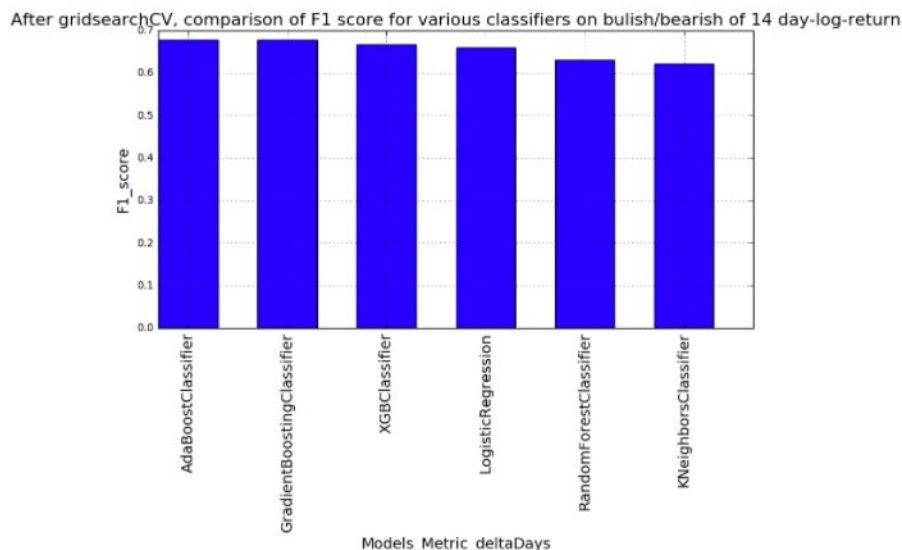
Figure 5(f) is summary of RMSE in order between the predicted value and actual value in testset for various regressors. There was no missing data for target value(price), but for the features, backfill was used to fill any missing data due to shifting X days to make moving average or return.

### C. Classification

Given with the significant RMSE of all regressors, and a concern of time dependence of price data itself or the correlation of moving average characteristics with the target price that may leads to false interpretation, the daily return was taken and converted to binary information. 1(bullish or no change) or 0(bearish) for each day were predicted for each day, and initially 12 features were used to fit and predict those binary target metric. When various hyperparameters were attempted to be tuned, it was very computationally heavy process, therefore, feature importance was taken and sorted in descending order using random forest classifier as a start as following:

1. feature RSI (0.187172)
2. feature MACD (0.105709)
3. feature Rolling\_std (0.089269)
4. feature SPY\_price (0.083011)
5. feature EMA\_fast (0.072642)
6. feature Bollinger\_lowerband (0.070900)
7. feature Bollinger\_upperband (0.070351)
8. feature SMA50 (0.069771)
9. feature SMA200 (0.066733)
10. feature EMA\_slow (0.066463)
11. feature Rolling\_mean (0.061869)
12. feature SMA20 (0.056111)

Therefore, only top 3 meaningful features, 'RSI','MACD', and 'SPY\_price' are selected to predict increase(or no change) or decrease of price. The data was balanced between increase/no change(1.0: 2564 points) and decrease (0.0: 2365).



**Figure 6.** Summary of f1 score in prediction of price increase/decrease in order for various classifiers with many

7 different classifiers were used along with grid search cross validation with different sets of hyper parameters. Cross validation was done on time series splitted 3 chunks of the entire dataset. Figure 6 shows f1 score ranks in descending order for 6 classifiers. SVC was dropped due to significant computational time. The best classifier was AdaBoostClassifier with the highest f1 score, 0.68. Among the sets of hyperparameters, `[{'n_estimators': [5, 50], 'learning_rate': [0.05, 0.5, 1.0]}]`, `{'n_estimators': 50, 'learning_rate': 0.5}` made the highest f1 score.

After prediction is made on the testset with the best model, Adaboost(`'n_estimators': 50, 'learning_rate': 0.5`), the total capital was backtested based on the following rules<sup>26</sup>: when price increase(1's) occurs in two consecutive days, trade decision is 'buy', and 'sell' for two consecutive 0's. Initial capital is \$100,000. Each day, 500 shares were purchased in the morning (at open price) when trade decision is 'buy', and sold at 'sell' signal. Figure 7(c) shows portfolio of initial cap + gain from trades over the course of investment period (testset). Total 54 times of 'sell' in red and 54 times of 'buy' in blue are made over the course of 5 years from February 27, 2012 to March 9, 2017 and it made ~\$45,000 profit. Trading cost for each transaction (e-trade \$4.95 per trade<sup>27</sup>) was not included in this plot.

## Refinement

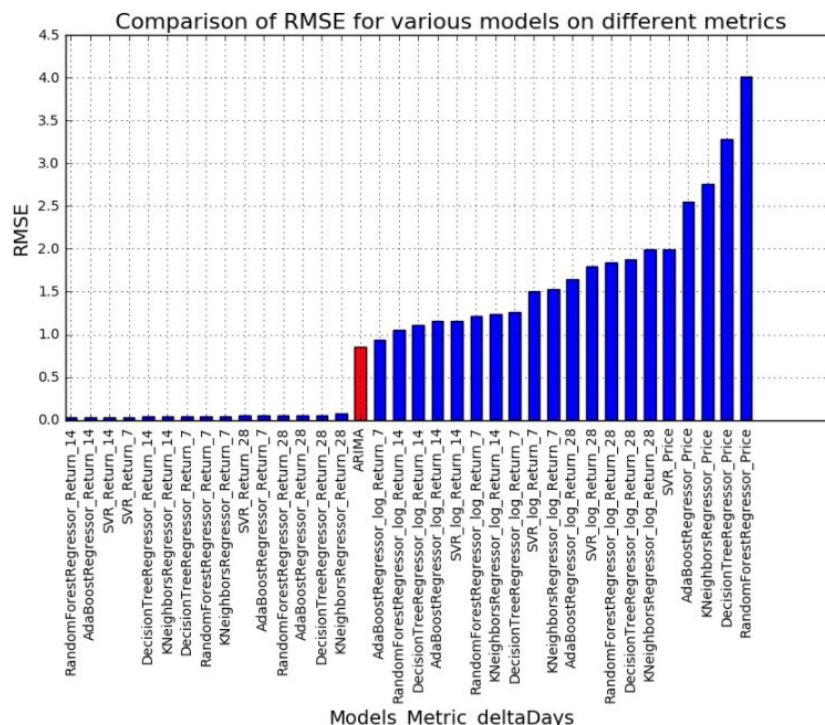
For the regression test, only price was used as target initially as discussed and shown in Figure 5(f). In order to remove a concern on time dependency of the price, return of 7,14,28 days, and further log of them were used as target for regression tasks. Figure 8 shows a significant RMSE reduction in prediction of return as opposed to use price only or ARIMA

<sup>26</sup> <http://francescopochetti.com/stock-market-prediction-part-introduction/>

<sup>27</sup> "E\*TRADE Fees and Rates | Pricing for Investing & Trading | E\*TRADE."  
<https://us.etrade.com/what-we-offer/pricing-and-rates>.



model. RandomForestRegressor predicts the best with the lowest RMSE of 0.027769 on return of 14 days.



**Figure 8.** Summary of RMSE in prediction of price, return of 7,14,28 day, log return of 7,14,28 days in order for various regressors with default args

## IV. Results

### Model Evaluation and Validation

For regression, RandomForestRegressor on return of 14 has the smallest RMSE, 0.028 in a hold out testset. Interestingly enough, any price regression show higher RMSE than return regression.

For classification that were tried to address any time dependency concern in ARIMA or regressor models, Adaboost classifier showed the highest f1 score, 0.68 after gridsearch cross validation as seen in Figure 6, which is much higher than simple knn regressor.

To elaborate how additional validation has performed especially for classification that would be extended to portfolio backtest, total 4929 datapoints were splitted into trainset(1997-01-02 to 2012-02-24,3697 points) and testset (2012-02-27 to 2017-03-09,1232 points). Grid search cross validation was carried out only in training set(3697 points) and the best model in this case was LogisticRegression with RMSE 0.64. RMSE on the hold out testset was 0.70. The

RMSE number seems to be pretty stable in this work flow with multiple classifiers with multiple grid spaces.

## Justification

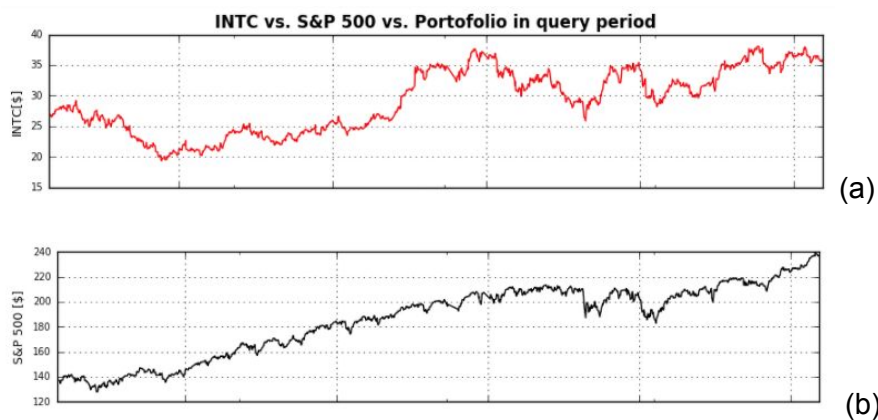
Figure 8 shows RMSE of all regressors in rank. The benchmark model ARIMA is marked in red and other models performed better than ARIMA. RandomForestRegressor on Return of 14 performed much better with the smaller RMSE, 0.027769. ARIMA predicted return, but almost most of regressors that were attempted here showed better performance. This means, the final model, randomforest regressor can predict the future return with significantly smaller error than typical time series ARIMA model can predict.

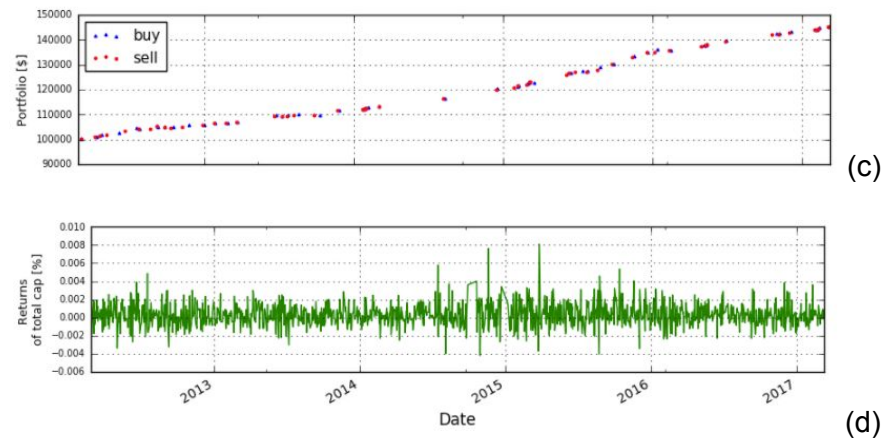
For classification task, the final model, Adaboost classifier is proved to provide better solution with f1 score 0.68 than relying on luck, which is 0.5 chance of flipping a coin.

To comment about net profit of this trading scheme with respect to just holding INTC or SPY, holding INTC stock (3846 shares which is \$100,000 worth on the one time purchase) for 5 years without any transaction will gain \$38k.

On the other hand, this machine learning based day trading scheme will make \$45k minus \$535 trading cost for 'Buy' and 'Sell' 108 times over 5 years, which is slightly more than what holding INTC stock only (\$~6400) with the same amount of money. Most important justification of this scheme is that even if there is down trend in INTC stock, the portfolio in Figure 7(b) shows only uptrend. That is meaningful enough to convince that machine learning based trading scheme is worth to investigate more about strategy.

However, in simply investment point of view, holding S&P 500 stocks (724 shares) will gain \$74k over 5 years of period of time. Even if the tax rate is assumed to be the same for the simplicity of calculation, it may not be worth to pursue buy/sell trading INTC stocks over 5 years compared to just holding S&P500 stock. More volatile stock with the greater range of swinging rather than stable INTC stock might be of next interest to try the model.





**Figure 7.** (a) INTC price per share in testset (b) SPY price per share (c) portfolio growth over the investment period. Initial capital is \$100,000, trading is made based on the prediction. Total value in y axis is initial capital+ cumulated profit based on 54 'buy' and 54 'sell' actions. (d) daily return of portfolio. Daily return itself doesn't grow over time.

## V. Conclusion

### Reflection

One particular security, Intel stock, was chosen for forecasting the future stock trend using various machine learning techniques. Regression and classification approaches were attempted for forecasting the trend. Instead of forecasting log return using ARIMA or price itself by regression, return of 14 is better target metric to predict than any other metrics or return of different number of days. For classification, the daily return was transformed into binary information and predicted with various classification models. Adaboost was the best performer. Trades (Buy/Sell) were simulated for the days of query upon prediction and portfolio backtest suggests that this machine learning based trading scheme on this INTC stock is proven to be slightly more profitable compared to just holding the equivalent amount of its stocks.

Interesting learning I had while researching the application of machine learning application in trading is that day traders are doing the pattern recognition every day to take much higher risk with commercial software that no one understands how it works. This simple, open source based, machine learning techniques can be a good assist for those day traders to take less risk with more confidence.

Difficult part of this project was that whether to decide to use only technical indicators or involves fundamental indicators as well. Finding fundamental indicators to evaluate company's health was much more difficult process than expected since most of critical database service is either provided for big organizations (such as Bloomberg terminal) or, with expensive fee was charged for any critical data (quandl). Eventually, only technical indicators were included for this paper, however, it would be very interesting to know in the future how powerful those fundamental indicators could impact on prediction performance.

### Improvement

Only INTC stock was tested for various regression and classification tasks. However, it would be interesting to know if this scheme is applicable for any trend or volatility of stock. Also once its robustness is proven to other stock, it can be expanded to create a portfolio of multiple stocks to make even greater profits with lesser risk.