

## **1. Definition:**

### **Project Overview:**

I have chosen the Investment and Trading Capstone project provided as one of the options. I wanted to use my Machine Learning knowledge to solve a real problem and where I can learn the domain too. I have always been interested in investment and trading as an amateur and now I want to use my machine learning knowledge and see if it could help in predicting stock prices in the short-term future as strategy for good investment. I extensively took help of the free courses – “Machine learning for Trading” and “Time Series Forecasting” for this Capstone Project. Based on the instructions provided in “Machine Learning free course”, I performed various financial analysis on the stocks data I downloaded from Yahoo Finance. I built few predictive models for predicting future prices one day, seven day, 14 ahead in future etc. from the last date provided for training:

- LINEAR REGRESSION based model
- KNN (K nearest Neighbors based model)
- ARIMA (Auto regressive integrated moving Average) based model (The free course “Time Series Forecasting” gave me a pretty good background on Time series data analysis and prediction)

### **Datasets:**

I have downloaded the following Datasets from yahoo finance: for e.g. – Apple, Google, Amazon, Facebook, Ge, Microsoft, SPDR Gold Shares (GLD), S&P 500. S&P 500 is the American market index, I used it to see trending behavior of other stocks by the overall market changes tracked by S&P 500.

**Problem Statement:** Build machine learning models that learn from historical stock price attributes and predict the stock price on a future date (any day after the last training date), I am only predicting the Adjusted Closing pricing.

**Strategy to solve the problem:** Building few models such as Linear Regression, KNN Regression which learn from historical stock data, attributes such as: Open, High, Low, Close, Volume, Adjusted Closing and “Adjusted Closing price” n days ahead in future. After the model is built on the historical data, it would predict the Adjusted Closing price for a date (only one date) which is “n” days ahead in future from the last date of training. Please note splitting the datasets has to be done in the ascending order of dates manually, and not using the SKLEARN test train split module which randomly split the data. In the stock price prediction world, we must not train a model on future datasets and try to predict the stock prices for past dates (Source - “Machine learning for Trading” Udacity course)

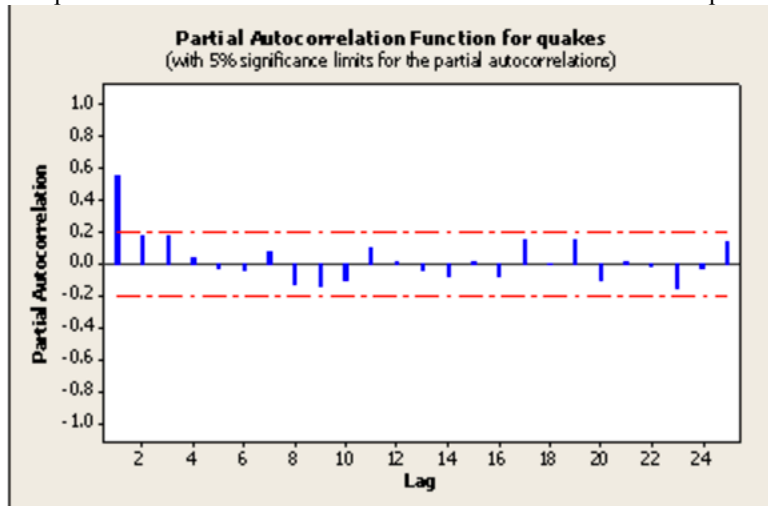
**Non-seasonal ARIMA model:** Non-seasonal ARIMA model is based on three terms: AR – Autoregressive - P, I – differencing - d, MA – Moving Average - q. Parameters p, d, and q are non-negative integers, p is the order (number of time lags) of the autoregressive model, d is the degree of differencing (the number of times the data have had past values subtracted), and q is the order of the moving-average model.

Three items should be considered to determine a first guess at an ARIMA model: a time series plot of the data, the ACF, and the PACF.

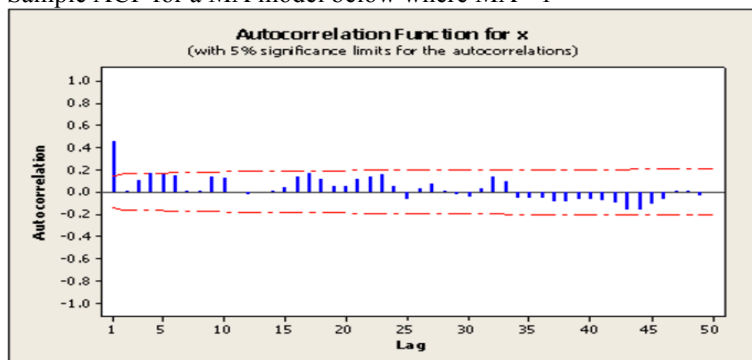
**AR model:** **Identification of an AR model is often best done with the PACF.** the theoretical PACF “shuts off” past the order of the model. The phrase “shuts off” means that in theory the partial autocorrelations are equal to 0 beyond that point

**MA model:** For an MA model, the theoretical PACF does not shut off, but instead tapers toward 0 in some manner. A clearer pattern for an MA model is in the ACF. The ACF will have non-zero autocorrelations only at lags involved in the model.

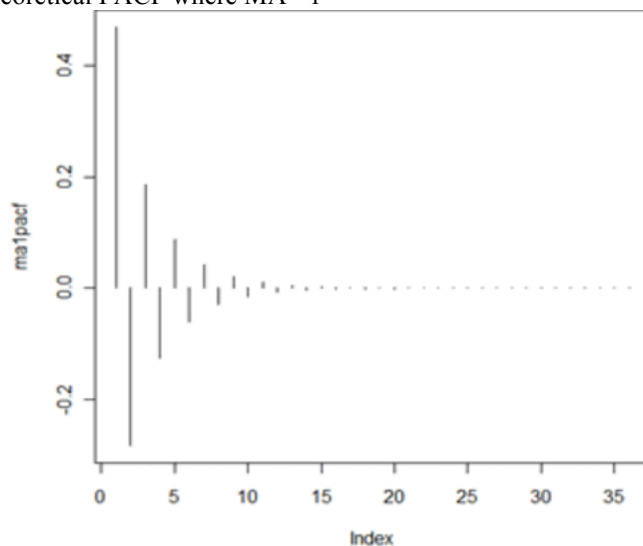
Sample PACF for an AR model below where the PACF shuts off past the 1<sup>st</sup> order so, AR =1 in this case:



Sample ACF for a MA model below where MA =1



Theoretical PACF where MA =1



Sources:

<https://onlinecourses.science.psu.edu/stat510/node/62>

<https://onlinecourses.science.psu.edu/stat510/node/49>

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

<http://www.inertia7.com/projects/time-series-stock-market-python>

**Metrics:** R2 score on the train and test datasets are calculated to measure the performance of the model. Best possible score is 1. Besides R2 score, I am also calculating the variation of the predicted stock price compared to the actual prices in the test dataset. R2 Score is used to test the performance of all three models – Linear Regression, KNN Regression and ARIMA

For e.g.:

Score on training data for GOOG: 0.97

Score on test data for GOOG: 0.88

The prediction on the test dataset 6 days after the training date for GOOG is +/- 2.721% compared to the actual values

Explaining R2 Score:

In statistics, the coefficient of determination, denoted R2 or r2 and pronounced "R squared", is a number that indicates the proportion of the variance in the dependent variable that is predictable from the independent variable. The coefficient R<sup>2</sup> is defined as  $(1 - u/v)$ , where u is the regression sum of squares  $((y_{\text{true}} - y_{\text{pred}})^2).sum()$  and v is the residual sum of squares  $((y_{\text{true}} - y_{\text{true.mean}})^2).sum()$ . Best possible score is 1.0 and it can be negative (because the model can be arbitrarily worse). A constant model that always predicts the expected value of y, disregarding the input features, would get a R<sup>2</sup> score of 0.0.

Source:

[http://scikit-learn.org/stable/modules/generated/sklearn.metrics.r2\\_score.html](http://scikit-learn.org/stable/modules/generated/sklearn.metrics.r2_score.html)

[http://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LinearRegression.html#sklearn.linear\\_model.LinearRegression.score](http://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html#sklearn.linear_model.LinearRegression.score)

[https://en.wikipedia.org/wiki/Coefficient\\_of\\_determination](https://en.wikipedia.org/wiki/Coefficient_of_determination)

## **2. Analysis:**

### **Data Exploration:**

I have downloaded daily historical stocks data from Yahoo Finance between 12-05-2011 and 12-02-2016 for the following stocks along with SNP 500 (American Market index, I want to see how these stock trend with respect to market trends explained by SNP 500):

- 'AAPL', 'GOOG', 'AMZN', 'FACEBOOK', 'GE', 'GLD', 'MICROSOFT'

Each of this stock has the following attributes:

Open, High, Low, Close, Volume, Adj Close and the Date

Sample Google Data and as well as the statistical description of the dataset:

	Date	Open	High	Low	Close	Volume	Adj Close
0	2016-12-02	744.590027	754.000000	743.099976	750.500000	1448700	750.500000
1	2016-12-01	757.440002	759.849976	737.025024	747.919983	2996900	747.919983
2	2016-11-30	770.070007	772.989990	754.830017	758.039978	2365800	758.039978
3	2016-11-29	771.530029	778.500000	768.239990	770.840027	1604500	770.840027
4	2016-11-28	760.000000	779.530029	759.799988	768.239990	2172200	768.239990

	Open	High	Low	Close	Volume	Adj Close
count	1258.000000	1258.000000	1258.000000	1258.000000	1.258000e+03	1258.000000
mean	711.782122	717.308367	705.651885	711.604109	3.211480e+06	527.225467
std	165.433068	166.190613	164.513760	165.309748	2.283091e+06	150.522426
min	494.652237	495.978230	487.562205	492.552239	7.900000e+03	279.246220
25%	578.488910	581.805133	573.535970	578.320983	1.639925e+06	396.808875
50%	689.695581	697.279998	682.376171	691.010010	2.633550e+06	536.272422
75%	779.232499	782.829986	773.508514	778.055327	4.198000e+06	635.259995
max	1226.802152	1228.882066	1218.602083	1220.172036	2.497790e+07	813.109985

All these Data sets are loaded in panda Data Frame with Date as the index, sample below with all rows with no values for 'SNP 500' removed. Meaning, we want to keep SNP 500 as the standard and base of our analysis, if we missing data for 'SNP 500' for certain dates, we want to remove all data for all stocks. Please refer the code snippet:

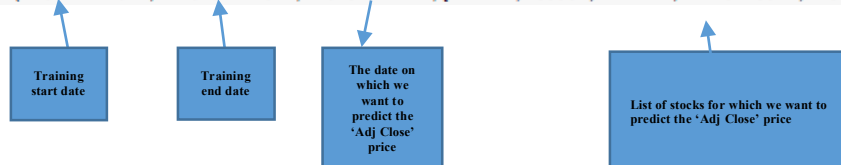
```
def get_data(symbols, dates):
    """Read stock data (adjusted close) for given symbols from CSV files."""
    df = pd.DataFrame(index=dates)
    if 'SNP_500' not in symbols: # add SPY for reference, if absent
        symbols.insert(0, 'SNP_500')

    for symbol in symbols:
        # TODO: Read and join data for each symbol
        df_temp = pd.read_csv("Data/{}.csv".format(symbol), index_col='Date', parse_dates=True, usecols=['Date', 'Adj Close'],
                              na_values=['nan'])
        df_temp = df_temp.rename(columns={'Adj Close': symbol})
        df = df.join(df_temp) #use default how='left'
        if symbol == 'SNP_500':
            df = df.dropna(subset=['SNP_500'])
    return df
```

An addition to the above attributes, I will be appending one more attribute – 'Adj\_Close\_req\_Days\_Later'. The days to be shifted to get the future 'Adj Close' price will be based on the user input – the number of days between training\_end date and the future date for which I need to predict 'Adj Close' price.

This is how I will call the “query\_regressor” function:

```
if __name__ == "__main__":
    query_regressor('2012-12-01', '2015-12-01', '2015-12-11', ['AAPL', 'GOOG', 'AMZN', 'FACEBOOK', 'GE', 'MICROSOFT'])
```



Code snippet below to get the number of days between the last training Date and the Prediction date:

# Capstone Project on Investment and Trading: Build a Stock Price Indicator Machine Learning Engineer Nanodegree

Partha Deka

Feb 12<sup>th</sup>, 2017

```
def train_regressor(start_date, end_date, stock_list, days_shift):  
    #prepare the date range for the DataFrame  
    #I am keeping the the maximum date available on my collected data from Yahoo Finance  
  
    dates = pd.date_range(start_date, end_date) #'2016-12-01'  
  
    #Setting the variables to pick the indices for the training and test datasets  
  
    print "days between the two training dates"  
    date_format = "%Y-%m-%d"  
    a = datetime.strptime(start_date, date_format)  
    b = datetime.strptime(end_date, date_format)  
    delta = b - a  
    num_days = delta.days
```

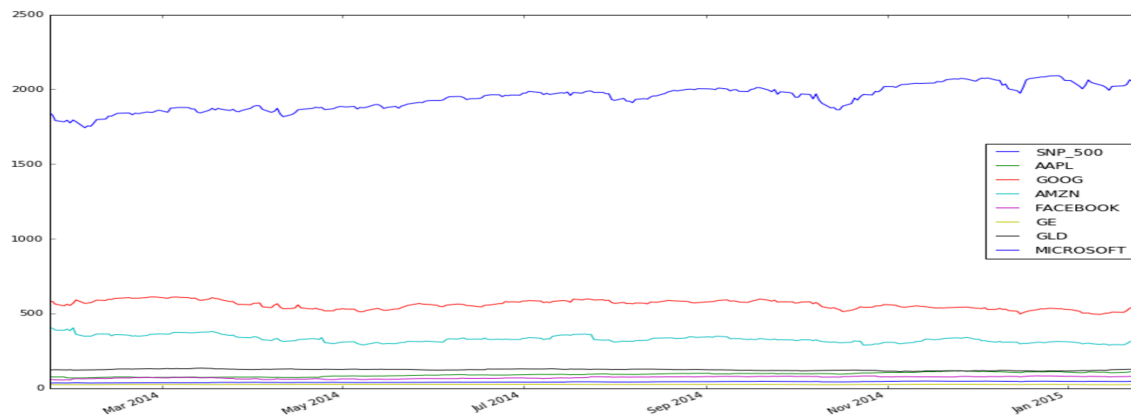
If there are non-trading days, for e.g. Saturday, Sunday or holidays, we want to forward fill first and then backward filling, code snippet below:

```
# Fill empty trade dates with forward filling dates first and then backward filling  
df.fillna(method="ffill", inplace="True")  
df.fillna(method="bfill", inplace="True")
```

## Exploratory Visualization:

Please refer the ipython notebook “All\_Analysis.ipynb” for all exploratory visualization code:

### HEADER 1: “Plotting Adj Close price for all stocks in comparison to the SNP 500 index”:



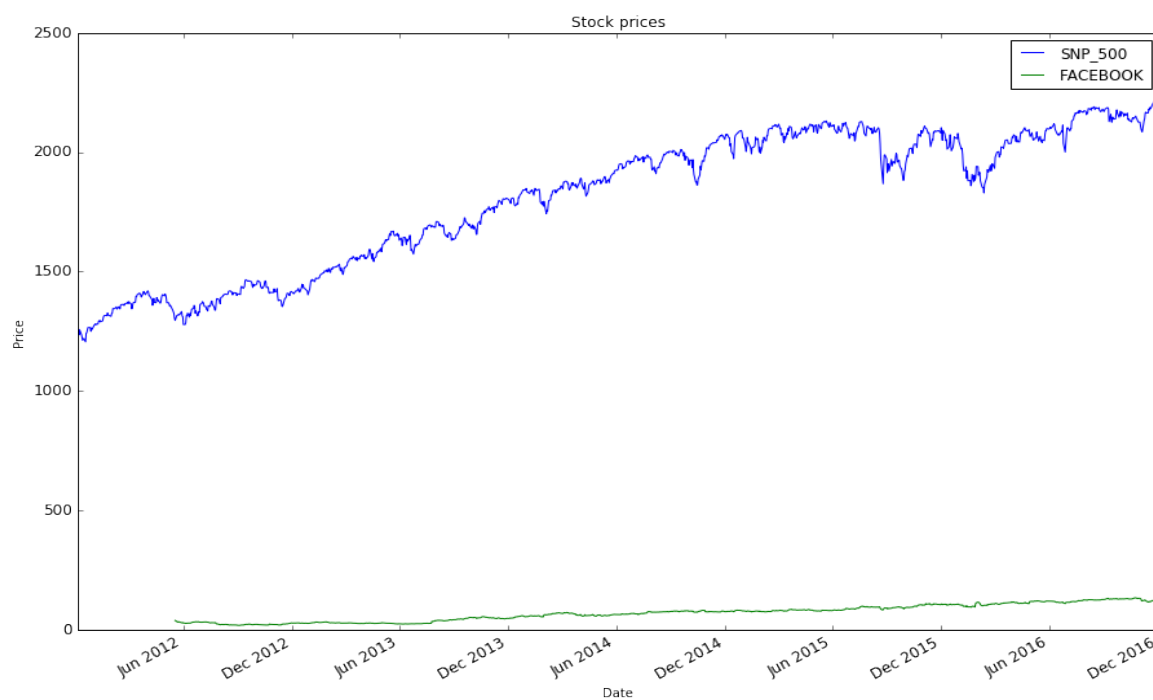
### HEADER 2: “Further Plotting and Analysis”:

Please find below the plot- Comparison SNP\_500 vs FACEBOOK:

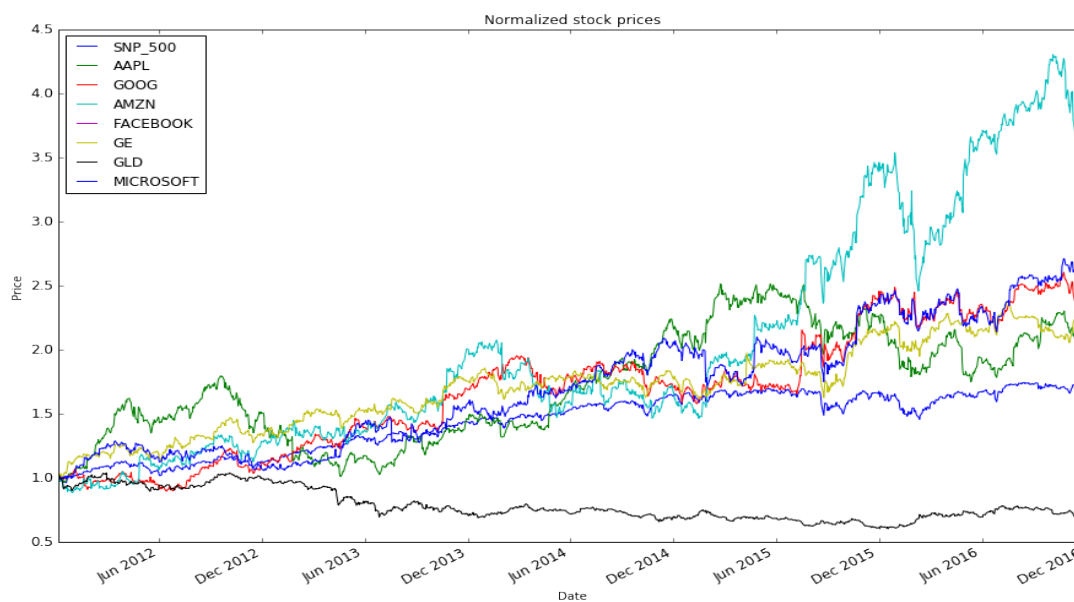
**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**



Please find below the Normalized plot for all stocks for “Adj Close” along with S&P\_500”:



Please find below the Normalized plot GE vs S&P\_500:

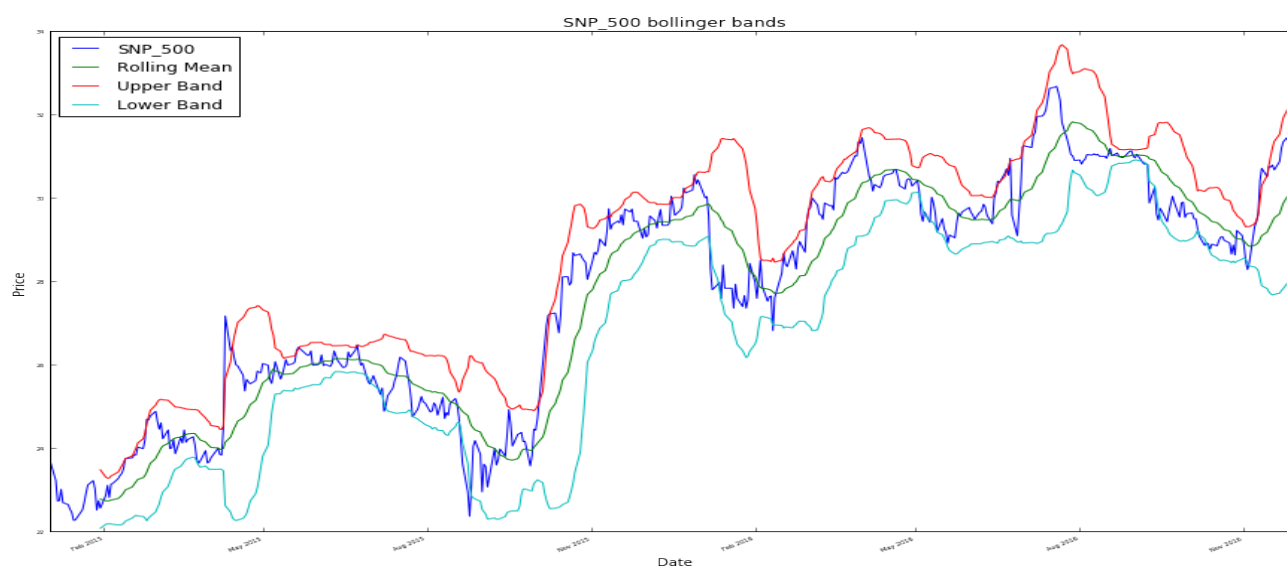
**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**



Bollinger Bands for GE: The popular notion is that when the stock price meets the lower Bollinger band (mean – 2\*Standard deviation), then it's usually a buy signal & when the stock price meets the upper Bollinger band (mean + 2\*Standard deviation) its usually a sell signal:



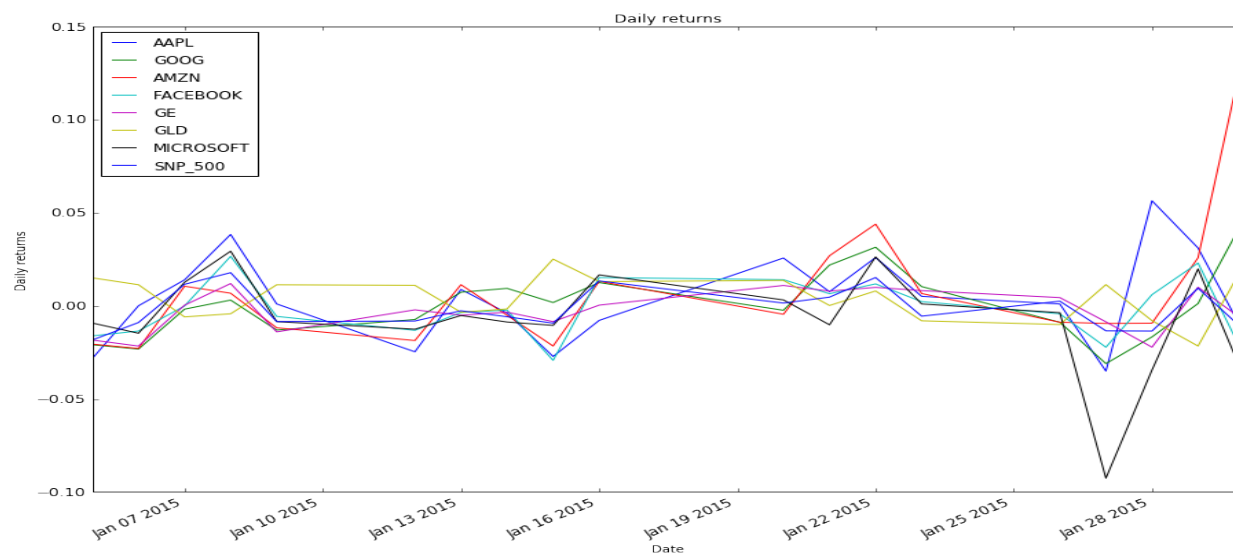
Daily returns:

- Comparing the Daily returns of all stocks for the month of Jan 2015:

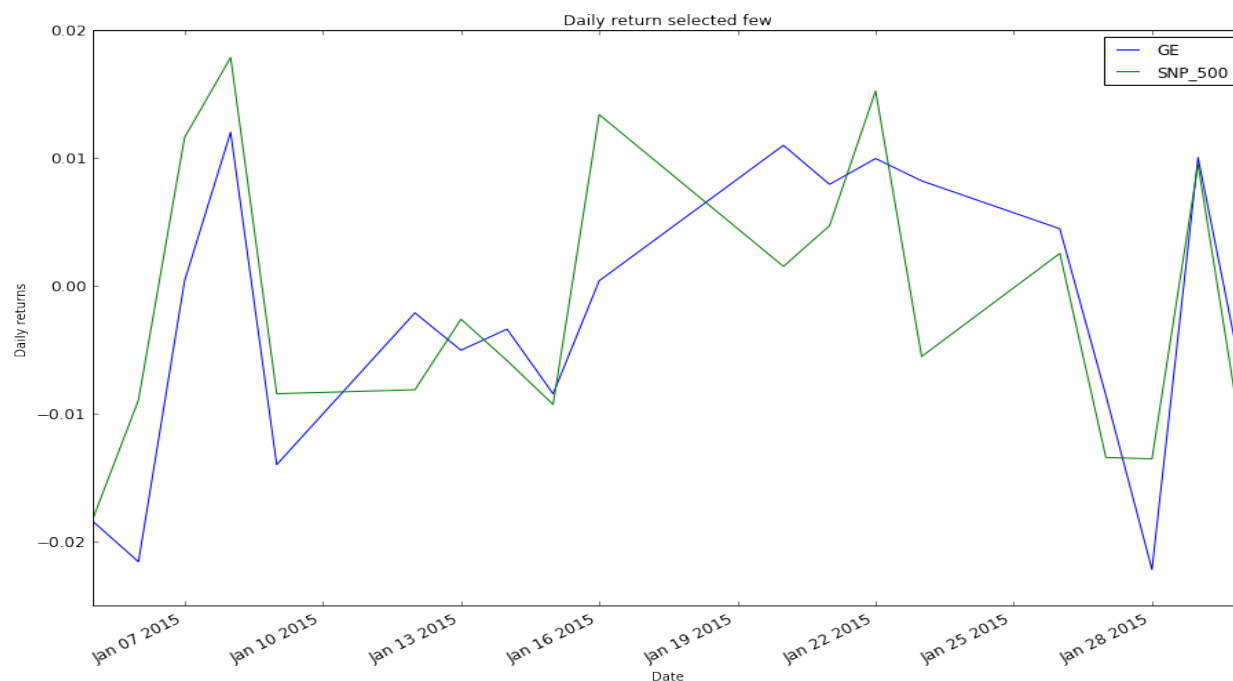
**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**

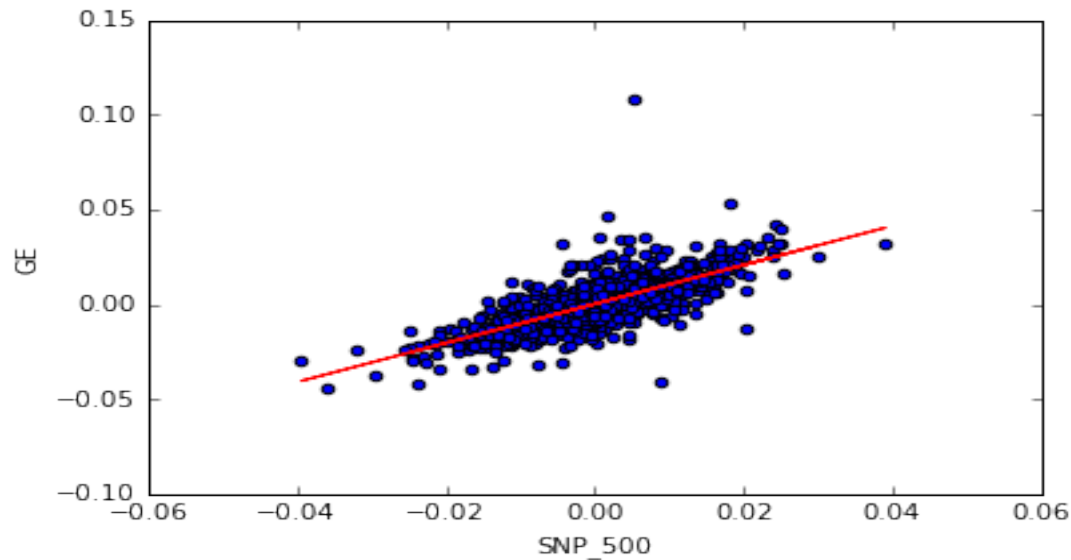


- Comparing daily returns between GE and SNP\_500 for the month of Jan 2015, it seems GE frequently moves in the same direction as SNP\_500 mostly, but in few instances, it also moves farther from SNP\_500:



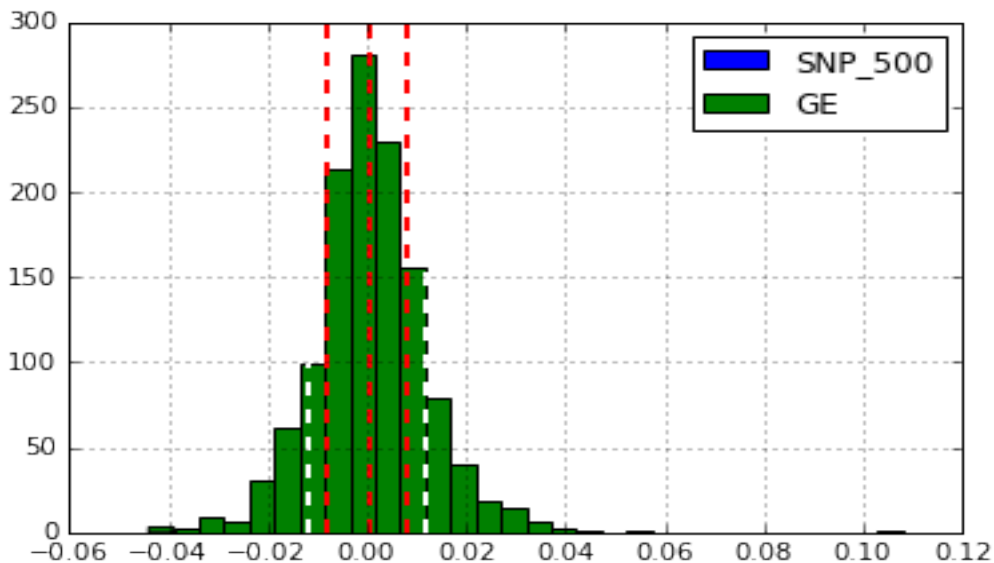
- Scatter plot between GE and SNP\_500 (between dates '2011-01-01', '2016-12-01') and fitting a line & looking at the statistics below. As we can see the slope or Beta is almost 1, this infers GE is very reactive to the market meaning if the market goes up by 1%, GE also goes up by 1 percent. Alpha is positive, which means GE performs little bit better than SNP 500 everyday on an average:





beta\_GE = 1.0291719923  
alpha\_GE = 0.000233231271241

- Daily return distribution of GE compared to SNP\_500 represented by histogram: The 'white' colored dotted line identifies GE's mean and std, while the 'red' colored identifies SNP\_500's mean and std. as we can see the distribution for both of them is very similar, other than GE has a higher std and higher Kurtosis which infers GE's distribution is not exactly Gaussian with flat tails.



Mean, Standard Deviation and Kurtosis for SNP\_500

Mean = SNP\_500 0.000476

std = SNP\_500 0.008204

kurtosis = SNP\_500 1.891211

**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**

Mean, Standard Deviation and Kurtosis for GE

Mean = GE 0.000723

std = GE 0.011632

kurtosis = GE 6.805789

- Pearson co-relation among Daily returns of the stocks, please refer below. We can see that GE is the only stock which has the best co-relation with SNP\_500 as we would infer from explanation above:

	SNP_500	AAPL	GOOG	AMZN	FACEBOOK	GE
SNP_500	1.000000	0.499828	0.561288	0.494226	0.332701	0.725937
AAPL	0.499828	1.000000	0.317679	0.236503	0.206110	0.298751
GOOG	0.561288	0.317679	1.000000	0.486442	0.327207	0.381830
AMZN	0.494226	0.236503	0.486442	1.000000	0.319441	0.321093
FACEBOOK	0.332701	0.206110	0.327207	0.319441	1.000000	0.220835
GE	0.725937	0.298751	0.381830	0.321093	0.220835	1.000000
GLD	-0.001239	0.028924	-0.040297	-0.023005	-0.024789	0.006026
MICROSOFT	0.613608	0.324849	0.428412	0.351674	0.217526	0.402501

	GLD	MICROSOFT
SNP_500	-0.001239	0.613608
AAPL	0.028924	0.324849
GOOG	-0.040297	0.428412
AMZN	-0.023005	0.351674
FACEBOOK	-0.024789	0.217526
GE	0.006026	0.402501
GLD	1.000000	-0.014787
MICROSOFT	-0.014787	1.000000

- Cumulative returns for all the stocks for the period between '2015-01-01' and '2016-12-01':

commulative return for SNP\_500 is:

6.46%

commulative return for AAPL is:

4.12%

commulative return for GOOG is:

42.51%

commulative return for AMZN is:

141.04%

commulative return for FACEBOOK is:

46.72%

commulative return for GE is:

32.64%

commulative return for GLD is:

-2.23%

commulative return for MICROSOFT is:

33.62%

Portfolio Statistics between the dates '2014-12-01', '2016-12-01':

- Code snippet with equal allocation to a portfolio of Six stocks:

# Capstone Project on Investment and Trading: Build a Stock Price Indicator Machine Learning Engineer Nanodegree

Partha Deka

Feb 12<sup>th</sup>, 2017

```
# Define a date range
dates = pd.date_range('2014-12-01', '2016-12-01')

# Choose stock symbols to read
symbols = ['AAPL', 'GOOG', 'AMZN', 'FACEBOOK', 'GE', 'MICROSOFT']

# Get stock data
df = get_data(symbols, dates)

# Fill empty trade dates with forward filling dates first and then backward filling
df.fillna(method="ffill", inplace=True)
df.fillna(method="bfill", inplace=True)

df_port = df.drop('SNP_500', axis=1)

# Normalize stock prices
df_normalized = normalize_data(df_port)
print "\n"
print "normalized: "
print df_normalized.head()
print "\n"

# Portfolio allocation for each stock
allocation = [0.166, 0.166, 0.166, 0.166, 0.166, 0.166]
df_allocation = df_normalized * allocation
print "allocated: "
print df_allocation.head()
print "\n"

# starting value for each stock -- dividing 1 million equally among all stocks
starting_value = [1000000]
df_values = df_allocation * starting_value
print "Portfolio Stock value : "
print df_values.head()
print "\n"

# Portfolio value by day
portfolio_values = pd.DataFrame()
portfolio_values['portfolio_value'] = df_values.sum(axis=1)
```

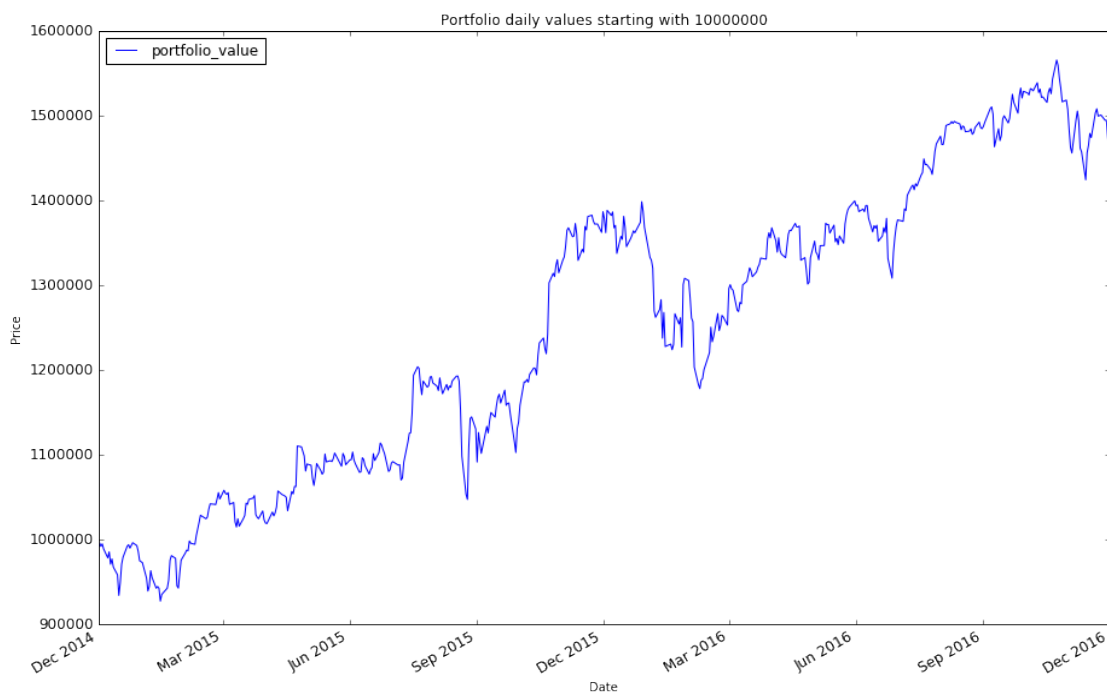
- Sample: Value of each Stock in the portfolio:

	AAPL	GOOG	AMZN	FACEBOOK	GE	MICROSOFT
2014-12-01	166000.000000	166000.000000	166000.000000	166000.000000	166000.000000	166000.000000
2014-12-02	165365.252335	165984.450553	166157.851742	166795.741246	166191.379370	165453.721704
2014-12-03	167240.637232	165228.776100	161162.576687	165513.712823	168296.682013	164156.321562
2014-12-04	166605.889567	167091.544313	161381.530178	166309.454069	166446.574596	166751.132657
2014-12-05	165899.018405	163344.253836	159191.965736	168785.093257	165936.197784	165317.148527

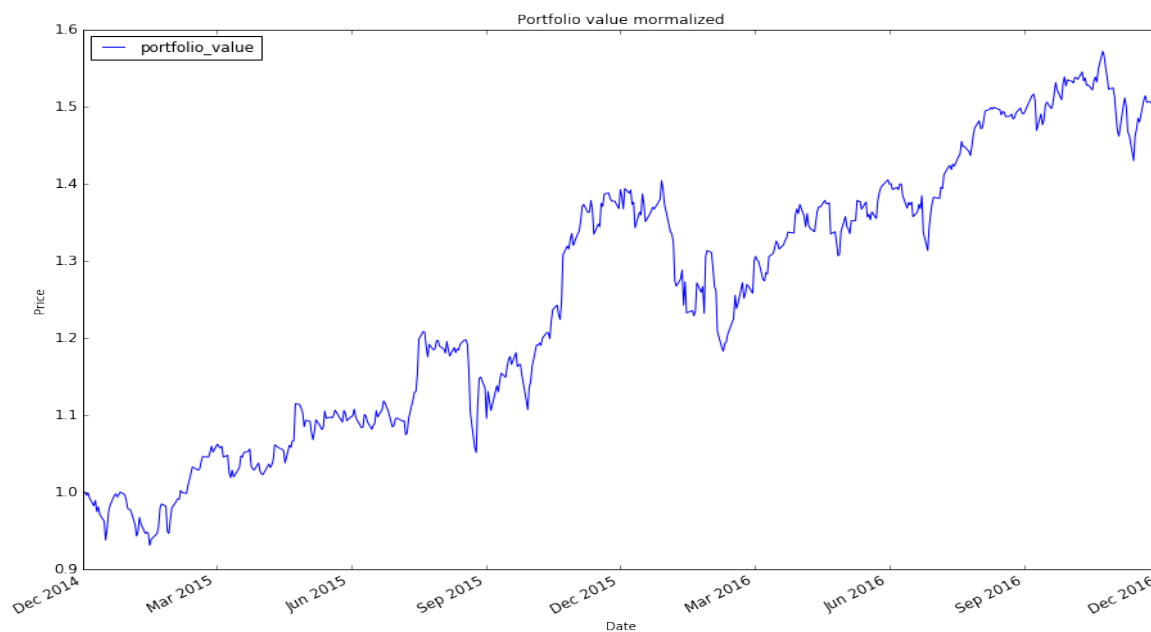
- Portfolio Value – summing up all the stock prices

	portfolio_value
2014-12-01	996000.000000
2014-12-02	995948.396951
2014-12-03	991598.706418
2014-12-04	994586.125379
2014-12-05	988473.677546

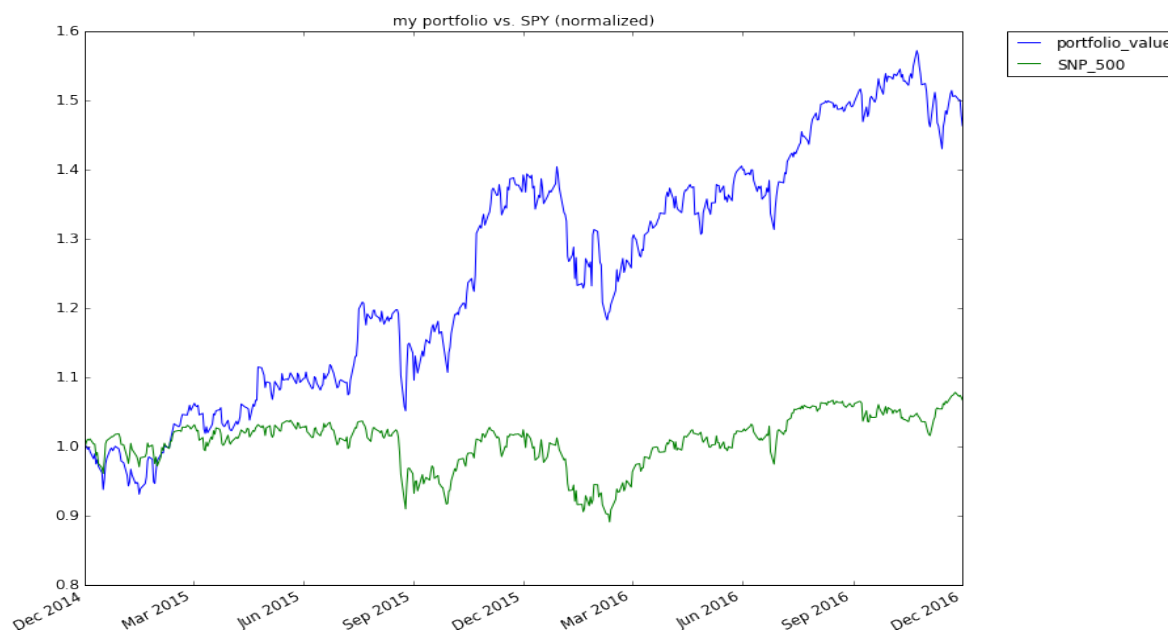
- Following is the plotted portfolio statistics:



- Following is the normalized portfolio statistics:



- Following is the portfolio value normalized plot vs normalized SNP\_500



- Sample Data of Daily returns on the Portfolio

Portfolio Daily returns:

	portfolio_value
2014-12-02	-0.000052
2014-12-03	-0.004367
2014-12-04	0.003013
2014-12-05	-0.006146
2014-12-08	-0.010533

- Portfolio Cumulative Return, Average Daily Return, Standard Deviation, Sharpe Ratio

```
Cumulative return: portfolio_value    0.463121
dtype: float64
Average daily return: portfolio_value    0.000832
dtype: float64
Daily standard deviation: portfolio_value    0.012561
dtype: float64
Sharpe ratio: portfolio_value    1.05206
dtype: float64
```

- Code Snippet to calculate Cumulative Return, Average Daily Return, Standard Deviation, Sharpe Ratio for the Portfolio

```
# Cumulative return
cr = cummulative_return(portfolio_values)
print "Cumulative return: ", cr

# Average daily return
print "Average daily return: ", daily_returns_portfolio.mean()

# standard deviation
print "Daily standard deviation: ", daily_returns_portfolio.std()

# Sharpe ratio
Sharpe_ratio = np.sqrt(252) * (daily_returns_portfolio.mean())/daily_returns_portfolio.std()
print "Sharpe ratio: ", Sharpe_ratio
print "\n"
```

### HEADER 3:

- Optimizing the Portfolio by optimizing or minimizing the negative sharpe ratio
- plotting optimal portfolio statistics
- Scatter plot, co-relation between my portfolio daily return and SNP\_500 daily return, Beta value to see how reactive is my optimized portfolio compared to the market, alpha value to see how well it performs compared to the market
- histogram

Sharpe Ratio is a measure for calculating risk-adjusted return, and this ratio has become the industry standard for such calculations

Source: <http://www.investopedia.com/terms/s/sharperatio.asp>

Code Snippet below to optimize sharpe ratio:

# Capstone Project on Investment and Trading: Build a Stock Price Indicator Machine Learning Engineer Nanodegree

Partha Deka

Feb 12<sup>th</sup>, 2017

```
def find_alloc(df, alloc_func, num_symbols):  
    # guess initial allocation  
    allocgss = [0.166, 0.166, 0.166, 0.166, 0.166, 0.166]  
  
    # has to bound the allocations between 0 and 1. setting the min, max pair for each element in x  
    bounds = [(0,1.0) for i in range(num_symbols)]  
  
    # call optimizer to minimize alloc_function  
  
    result = spo.minimize(alloc_func, allocgss, args=(df,), bounds=bounds, method='SLSQP',  
                          options={'disp':True},  
                          constraints=({ 'type': 'eq', 'fun': lambda inputs: 1.0 - np.sum(inputs)})).x  
  
    return result  
  
def alloc_function(alloc, df):  
    trading_period_days = 252  
    starting_value = 1000000  
  
    #Normalize stock and apply allocation  
    df_normalized = normalize_data(df)  
    df_normalized_allocation = df_normalized * alloc  
  
    # Calculate portfolio value by day  
    starting_value = [1000000]  
    # normalize portfolio  
    df_normed_alloc_value = df_normalized_allocation * starting_value  
    portfolio_values = df_normed_alloc_value.sum(axis=1)  
  
    # Portfolio value by day  
    portfolio_val = pd.DataFrame()  
    portfolio_val['portfolio_value'] = portfolio_values  
    |  
    # Compute daily returns on the portfolio  
    daily_returns_portfolio = compute_daily_return(portfolio_val)  
  
    # Higher the sharp ratio the better (Reward / Risk) , so we should minimize the negative sharp ratio  
    sharpe_ratio = -1 * np.sqrt(trading_period_days) * (daily_returns_portfolio.mean()/daily_returns_portfolio.std())  
    return sharpe_ratio
```

Optimal allocation below after Sharpe Ratio optimization:

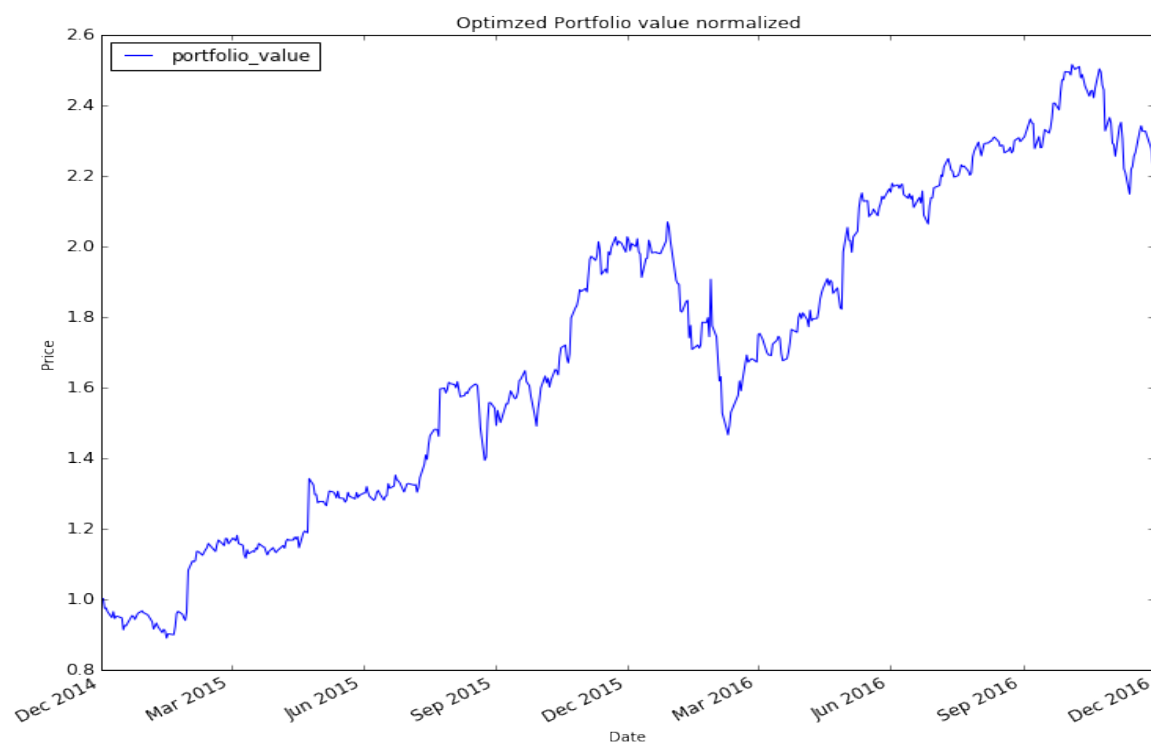
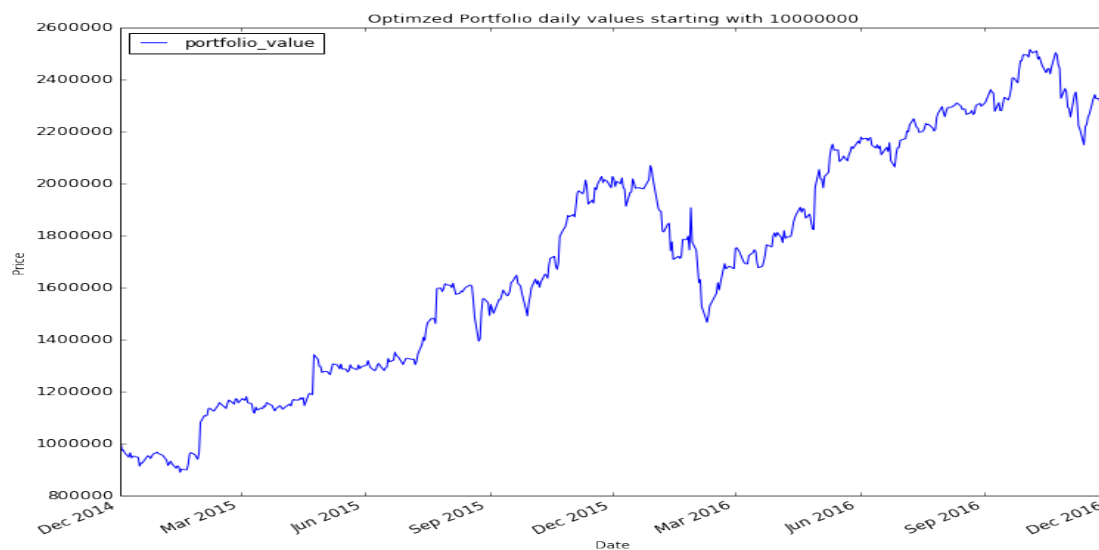
```
optimal allocations:  
AAPL: 0.000  
GOOG: 0.000  
AMZN: 0.914  
FACEBOOK: 0.086  
GE: 0.000  
MICROSOFT: 0.000
```

Plotted Optimal Portfolio Statistics:

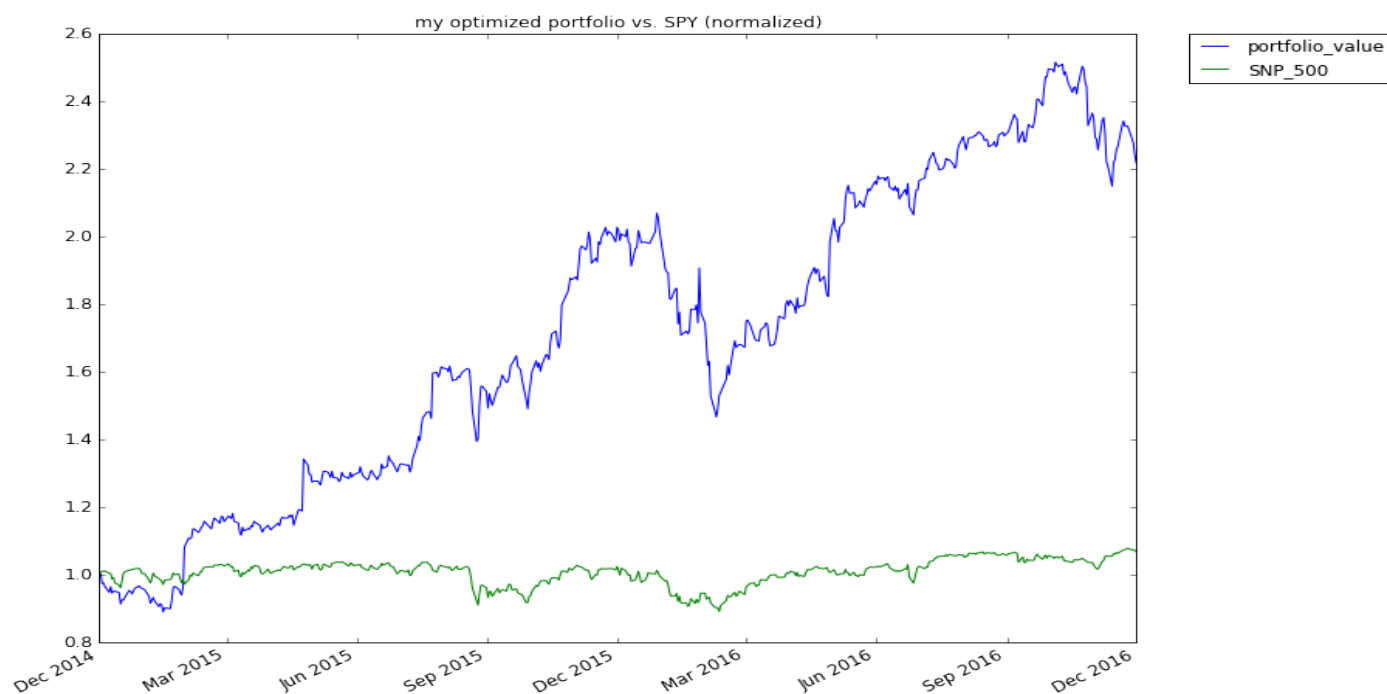
**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**



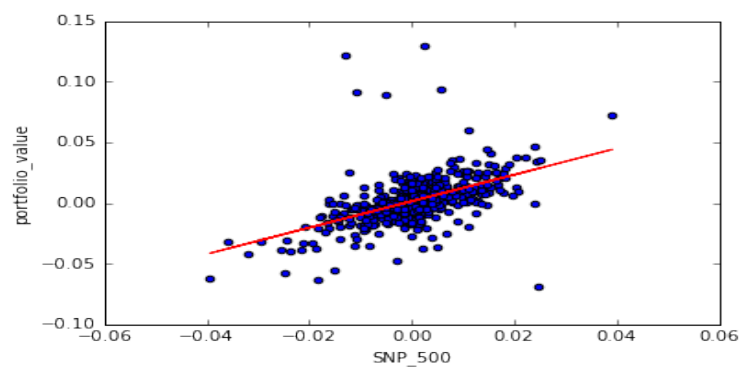




Portfolio Cumulative Return, Average Daily Return, Standard Deviation, Sharpe Ratio:

```
Cumulative return: portfolio_value    1.216763
dtype: float64
Average daily return: portfolio_value    0.00176
dtype: float64
Daily standard deviation: portfolio_value    0.019296
dtype: float64
Sharpe ratio: portfolio_value    1.448061
dtype: float64
```

Scatter plot between my optimized portfolio and SNP\_500:



# Capstone Project on Investment and Trading: Build a Stock Price Indicator Machine Learning Engineer Nanodegree

Partha Deka

Feb 12<sup>th</sup>, 2017

```
beta_AAPL = 1.08940744037  
alpha_AAPL = 0.00157438228967
```

alpha value is positive, so in an average my optimized portfolio is performing little bit better than the SNP\_500

## Algorithms and Techniques:

For this Capstone Project, I am using three supervised based algorithms – Linear Regression, KNN Regression, ARIMA (Auto Regressive integrated moving average for time series data)

### Linear Regression:

In statistics, linear regression is an approach for modeling the relationship between a scalar dependent variable  $y$  and one or more explanatory variables (or independent variables) denoted  $X$ . The case of one explanatory variable is called simple linear regression. For more than one explanatory variable, the process is called multiple linear regression.

Source: [https://en.wikipedia.org/wiki/Linear\\_regression](https://en.wikipedia.org/wiki/Linear_regression)

Multiple Linear Regression represented in the form:  $Y = m_1X_1 + m_2X_2 + m_3X_3 + \dots + \text{Slope}$

In this project for Stock Price Prediction for each stock the  $X$  attributes are: 'Open', 'High', 'Low', 'Close', 'Volume', 'Adj Close'

After finding the Pearson co-relation among these attributes I found out that the attributes 'Open', 'High', 'Low', 'Close' are highly co-related. So, I am keeping only 'Open', 'Volume', 'Adj Close' as the  $X$  attributes

I am adding new column 'Adj\_Close\_req\_Days\_Later' as a target ( $Y$ ) attributed to be predicted by the algorithm after training is done. This attribute is created as follows:

```
d[key]['Adj_Close_req_Days_Later'] = d[key]['Adj Close']  
d[key]['Adj_Close_req_Days_Later'] = d[key]['Adj_Close_req_Days_Later'].shift(-days_shift)
```

Here, I want to build an algorithm that learns from prices of a stock in historical time and predicts the price of the stock for a date in future. In other words, the algorithm should model the relationship between the dependent variables such as 'Open', 'Volume', 'Adj Close' with the Stock price. So, the linear regression function or algorithm should learn from the relationship between the dependent variables and the stock prices producing the coefficients for each of the stock and the slope. For new values of dependent variables, we can then plug them in the linear regression equation and get the price or the predicted stock price. But the caveat in using linear regression for stock price prediction is it does not track the seasonality and the trend in the data well, it basically thinks the data points follow a linear path without paying attention to ups and dips in the stock prices. Another option is to use KNN regression, where it can relate to the behavior of the local data points rather than following a linear path.

### KNN Regression:

The k-Nearest Neighbors algorithm (or k-NN for short) is a non-parametric method used for classification and regression. In both cases, the input consists of the  $k$  closest training examples in the feature space. In k-NN regression, the output is the property value for the object. This value is the average of the values of its  $k$  nearest neighbors. k-NN is a type of instance-based learning, where the function is only approximated locally and all computation is deferred until classification or Regression. The k-NN algorithm is among the simplest of all machine learning algorithms. Both for classification and regression, it can be useful to assign weight to the contributions of the neighbors, so that the nearer neighbors contribute more to the average than the more distant ones.

Source: [https://en.wikipedia.org/wiki/K-nearest\\_neighbors\\_algorithm](https://en.wikipedia.org/wiki/K-nearest_neighbors_algorithm)

In this case, for a given date in future the function would approximate the stock price based on the nearby values in the nearest past.

\*\*Another reason I tried these two models as they were well discussed in the “Machine Learning for Trading” Course

#### **ARIMA model:**

ARIMA stands for Auto Regressive Integrated Moving Average. There are two types ARIMA models – Seasonal and Non-Seasonal. I am using Non-seasonal ARIMA model in this Stock Price Prediction problem. Non-seasonal ARIMA models forecast future points based on the construction of three components, an autoregressive, AR, a differencing, I, and a moving average, MA, component. When we put it all together non-seasonal ARIMA models are displayed as ARIMA  $p, d, q$ . The  $p$ ,  $d$  and  $q$  values represent the number of periods to lag for in our ARIMA calculation. For e.g. if we say  $p=2$ , we will be using the two previous periods of the time series in the autoregressive portion of the calculation. This helps adjust the line fitted to forecast the series. Purely, AR models would look like a linear regression where the predictive variables are a certain number of previous periods which is represented by  $p$ . The differencing term refers to the process we use to transform a time series into a stationary one, which is a series without trend or seasonality. This process is called differencing and the ‘ $d$ ’ refers to the number of transformation used in the process. The moving average term (MA) refers to the lag of the error component. The error component refers to the part of the time series not explained by trend or seasonality. MA models look like linear regression models where the predictive variables are the previous  $q$  periods of errors.

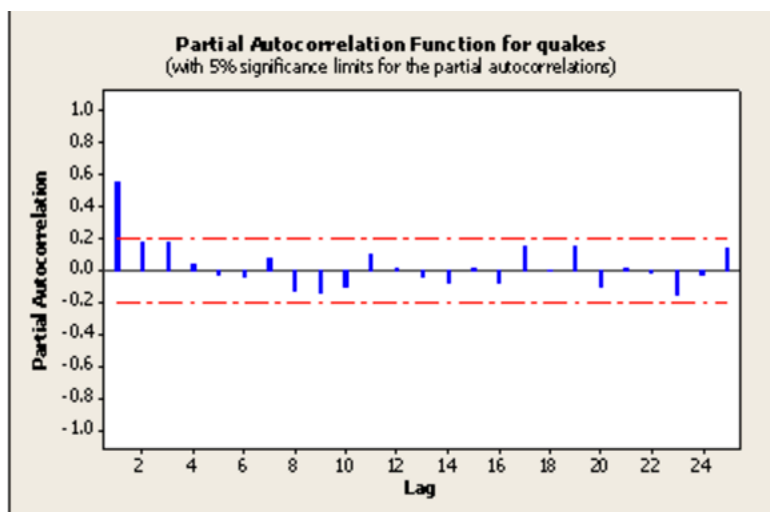
Source: Udacity “**Time Series Forecasting**” Course

Three items should be considered to determine a first guess at an ARIMA model: a time series plot of the data, the ACF, and the PACF.

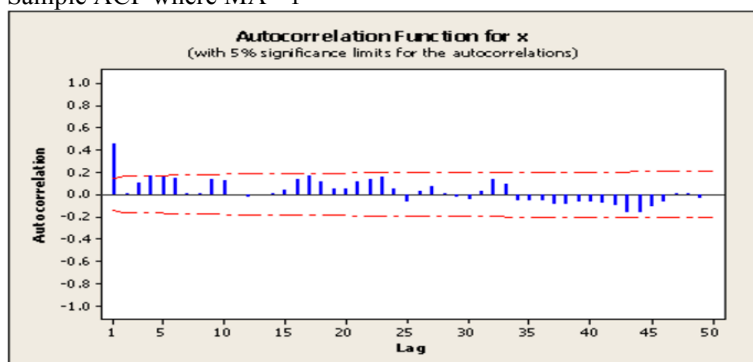
AR model: **Identification of an AR model is often best done with the PACF.** the theoretical PACF “shuts off” past the order of the model. The phrase “shuts off” means that in theory the partial autocorrelations are equal to 0 beyond that point

MA model: For an MA model, the theoretical PACF does not shut off, but instead tapers toward 0 in some manner. A clearer pattern for an MA model is in the ACF. The ACF will have non-zero autocorrelations only at lags involved in the model.

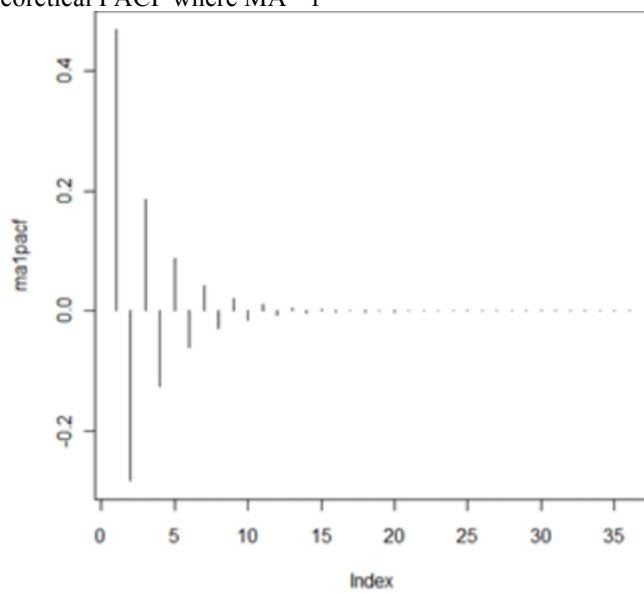
Sample PACF for an AR model where the PACF shuts off past the 1<sup>st</sup> order so,  $AR=1$  in this case



Sample ACF where MA =1



Theoretical PACF where MA =1



Sources:

**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**

Udacity Free Course: Time Series Forecasting

<https://onlinecourses.science.psu.edu/stat510/node/62>

<https://onlinecourses.science.psu.edu/stat510/node/49>

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

<http://www.inertia7.com/projects/time-series-stock-market-python>

### **Benchmark**

For all the three models my benchmark is a R2 Score of 0.6 on test data(unseen data) and a R2 Score of 0.85 on train data. Please note if we pass a list of stocks in the “query\_regressor” function for both Linear and KNN regression, my code builds individual model for each stock. Each individual model for each must meet or exceed the benchmark.

Additionally, on an average the predicted prices from all the models **must be more or less within +/- 5%** of the actual stock prices.

## **3.Methodology:**

### **Data Preprocessing: (Please refer the All\_Analysis.ipyn ipython notebook “HEADER 2”)**

I have chosen these stocks for training and prediction: 'AAPL','GOOG', 'AMZN', 'FACEBOOK', 'GE', 'MICROSOFT' as well as 'SNP\_500' index and 'GLD' for comparative analysis. I have manually downloaded the data from Yahoo finance, I went back as far as December 2011 until 1<sup>st</sup> week of December 2016. Please note that some stocks like 'Facebook' did not even exist back in 2011.

I have built the following data preprocessing utility functions:

\*\*\* I have taken extensive help from the udacity free course “Machine Learning for Trading”

Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree

Partha Deka

Feb 12<sup>th</sup>, 2017

```
def get_data(symbols, dates):
    """Read stock data (adjusted close) for given symbols from CSV files."""
    df = pd.DataFrame(index=dates)
    if 'SNP_500' not in symbols: # add SPY for reference, if absent
        symbols.insert(0, 'SNP_500')

    for symbol in symbols:
        # TODO: Read and join data for each symbol
        df_temp=pd.read_csv("Data/{}.csv".format(symbol), index_col='Date', parse_dates=True, usecols=['Date',
                                                                                                     'Adj Close'],
                             na_values=['nan'])
        df_temp = df_temp.rename(columns={'Adj Close': symbol})
        df = df.join(df_temp) #use default how='left'
        if symbol == 'SNP_500':
            df = df.dropna(subset=["SNP_500"])
    return df

def normalize_data(df):
    """Normalize data by using the first row of the dataframe"""
    return df/df.ix[0,:]

def compute_daily_return(df):
    daily_returns = (df / df.shift(1)) - 1
    daily_returns.ix[0, :] = 0
    return daily_returns[1:]

def cumulative_return(df):
    cumulative_return = (df.ix[-1] / df.ix[0]) - 1
    return cumulative_return

def plot_data(df, title="Stock prices", xlabel="Date", ylabel="Price"):
    ax = df.plot(title=title, figsize=(12,10), fontsize=12)
    ax.set_xlabel(xlabel)
    ax.set_ylabel(ylabel)
    plt.show()
```

```
def plot_selected(df, columns, start_index, end_index, title="Stock prices"):
    """Plot the desired columns over index values in the given range."""
    df_new = df.ix[start_index : end_index, columns]
    ax = df_new.plot(title=title, fontsize=10, figsize = (15,10))
    ax.set_xlabel("Date")
    ax.set_ylabel("Price")
    plt.show()

def normalize_data(df):
    """Normalize data by using the first row of the dataframe"""
    return df/df.ix[0,:]

def get_rolling_mean(values, window):
    """Return rolling mean given values, using specified window size"""
    return pd.rolling_mean(values, window=window)

def get_rolling_std(values, window):
    """Return rolling standard deviation of given values, using specified window size"""
    return pd.rolling_std(values, window=window)

def get_bollinger_bands(rm, rstd):
    """Return upper lower bollinger bands."""
    upper_band = rm + rstd * 2
    lower_band = rm - rstd * 2
    return upper_band, lower_band

def compute_daily_return(df):
    daily_returns = (df / df.shift(1)) - 1
    daily_returns.ix[0, :] = 0
    return daily_returns[1:]

def cumulative_return(df):
    cumulative_return = (df.ix[-1] / df.ix[0]) - 1
    return cumulative_return
```

**Explanation of some of the utility functions:**

**get\_data:** This function reads the csv files from the 'Data' directory and creates a DataFrame with 'Adj Close' price for each stock with Date as the index. It also replaces the columns with the Stock name with the Adj Close as the values for the respective stocks. Additionally, I am removing all the records from the DataFrame for which value for SNP\_500 is 'nan' (we are already stock values when the market was open). I am also sorting the stocks with ascending order of dates. Within the main function I am so doing forward fill first and then backward for stocks which have 'nan' values (We cannot keep 'nan' values e.g. non-trading days like holidays or weekends):

```
# Fill empty trade dates with forward filling dates first and then backward filling  
df.fillna(method="ffill", inplace="True")  
df.fillna(method="bfill", inplace="True")
```

**normalize\_data:** This function returns the normalized DataFrame by normalizes the data by dividing each value of each stock by the starting value (starting date of the stock) or in a word diving each row by the first row for the corresponding columns.

**get\_rolling\_mean:** This function returns rolling mean given values, using specified window size.

**get\_rolling\_std:** This function returns rolling standard deviation of given values, using specified window size.

**get\_bollinger\_bands:** This function returns the Upper and lower Bollinger bands using the rolling mean and rolling standard deviation.

**compute\_daily\_return:** This functions returns the daily return of the stock compared to the previous date's price. Please note this will return 'nan' for the starting date as we don't 've the value for that stock in the previous date, so we need to set the first row of the Daily return DataFrame to '0'.

Please refer the ipython notebook book: **All\_Analysis.ipynb** – **HEADER 2** and **HEADER 3**

**Implementation:**

Please refer the following headers in the ipython notebook book "All\_Analysis.ipynb" :

**# HEADER 4**

**- STOCK PRICE PREDICTION INTERFACE --- LINEAR REGRESSION**

**# HEADER 5**

**- STOCK PRICE PREDICTION INTERFACE --- KNN REGRESSION**

**# HEADER 6**

**- STOCK PRICE PREDICTION INTERFACE: ARIMA**

**LINEAR REGRESSION:**

I have used SKLEARN's Linear Regression to implement this algorithm.

```
class sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False, copy_X=True, n_jobs=1)
```

I have used the get\_data function to preprocess the data fames for each stock. For training the algorithms, I have implemented the 'get\_data' function in a little different way than the one I implemented for the initial data analysis as discussed earlier in this document:

# Capstone Project on Investment and Trading: Build a Stock Price Indicator Machine Learning Engineer Nanodegree

Partha Deka

Feb 12<sup>th</sup>, 2017

```
def get_data(symbol, dates):
    df = pd.DataFrame(index=dates)
    df_temp = pd.read_csv("Data/{}.csv".format(symbol), index_col='Date', parse_dates=True, usecols=['Date', 'Open',
                                                                                                   'High', 'Low',
                                                                                                   'Close', 'Volume',
                                                                                                   'Adj Close'],
                        na_values=['nan'])

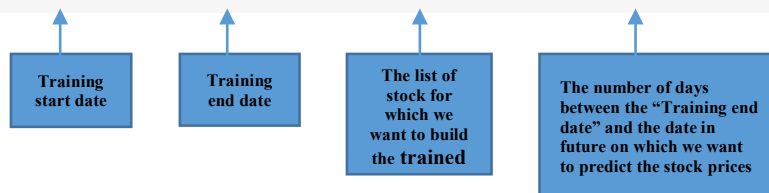
    df = df.join(df_temp)
    df.fillna(method="ffill", inplace=True) # Forward fill empty trade dates
    df.fillna(method="bfill", inplace=True) # backfill empty trade dates
    return df
```

Please note that I am passing a stock and not a list of stocks to this function, pulling all attributes for that stock in the DataFrame. I want to see how the attributes are co-related to one another and use the attributes which are least co-related to build the model. My goal is to build a model for each stock for the list of stocks passed in the “query\_regressor” function. I’ve also implemented few extra functions: “train\_regressor”, “query\_regressor” & “train\_classifier\_gridsearch”

Explanations of the new functions implemented:

“train\_regressor”:

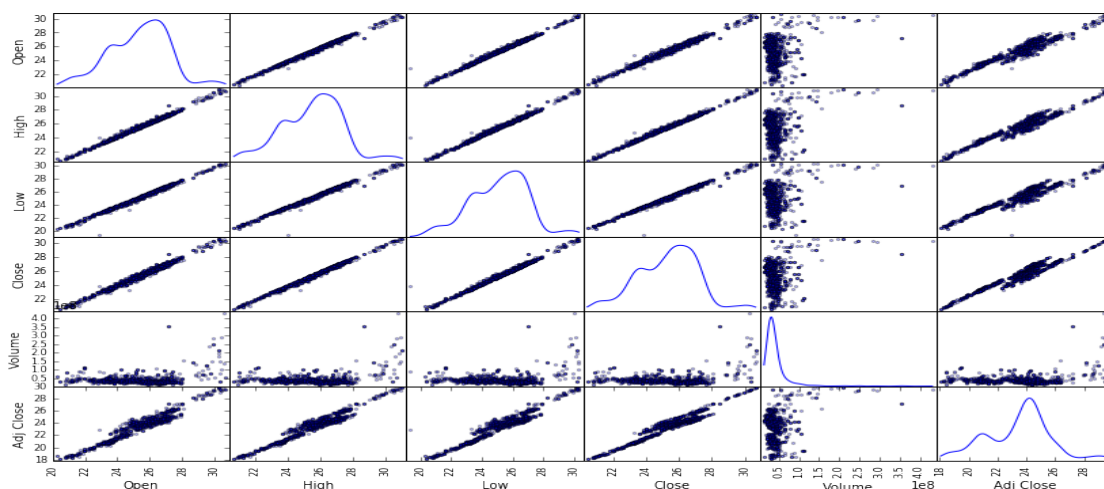
```
def train_regressor(start_date, end_date, stock_list, days_shift):
```



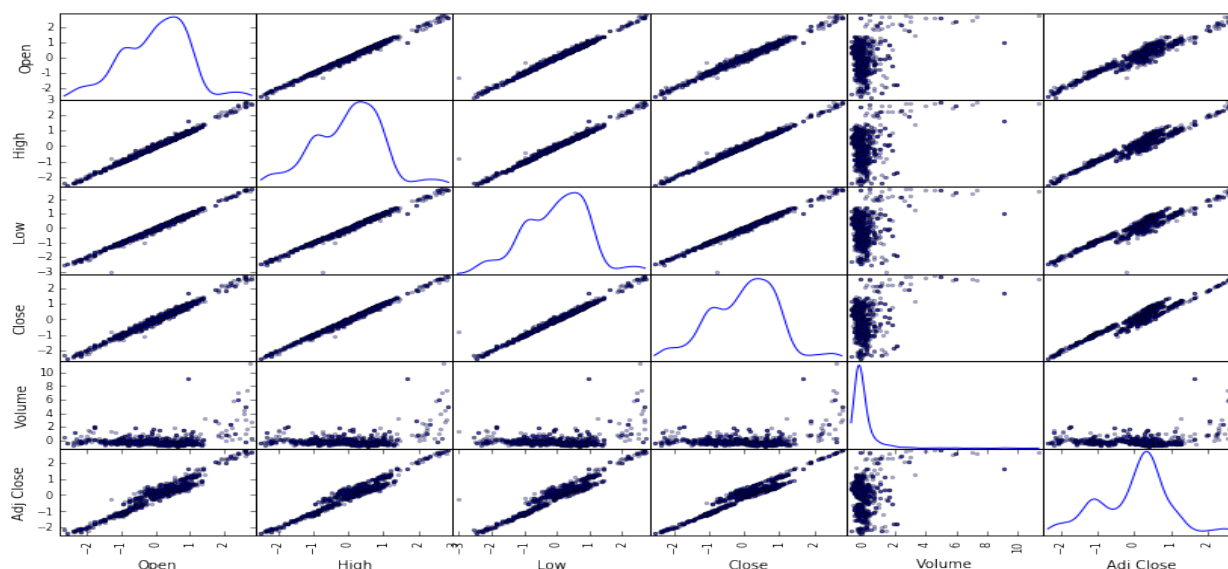
This function primarily returns a **dictionary of trained models**(with key as the stock name and value as the **trained regressor**) for all stocks in the list and also performs analysis on the attributes such as how are they co-related to one another and as well as their distributions. I could find out that the attributes ['Open', 'High', 'Low', 'Close'] are highly correlated so I decided to keep ‘Open’, ‘Volume’, ‘Adj Close’ for each stock. The distributions for most of the attributes for most of the stocks are non-Gaussian, although few of them are close to normal. Nevertheless, I used the SKLEARN’s “StandardScaler” for all of them to convert them to zero mean normal distributions.

Below is the distribution of GE’s attributes before performing standardization:





Below is the distribution of GE's attributes after performing Standardization:



I am doing a manual test train on the data as it is advised in the course – “Machine Learning for trading” that we need to split in the sequential in the stock price prediction world. The Sklearn’s test train split module performs randomized test train split, we must not get into a situation where we are testing on datasets which is prior to trained datasets. Please refer the below the code snippet to perform the manual test train split:

```
X_train = d[key].iloc[0:train_num,:-1]
scaler.fit(d[key].iloc[0,:-1])
X_train = scaler.transform(X_train)
y_train = d[key].iloc[0:train_num, -1]
X_test = d[key].iloc[train_num:num_days,:-1]
X_test = scaler.transform(X_test)
y_test = d[key].iloc[train_num:num_days, -1]
```

Just re-iterating, please also note that I am engineering the feature “Adj\_Close\_Req\_Days\_Later”. Here “days\_shift” is the number of days between the last training date and the date on which I want to predict the stock price in future :

```
d[key]['Adj_Close_req_Days_Later'] = d[key]['Adj_Close']
d[key]['Adj_Close_req_Days_Later'] = d[key]['Adj_Close_req_Days_Later'].shift(-days_shift)
```

I also tried to perform hyperparameter tuning for few parameters such as ‘fit\_intercept’, ‘normalize’, ‘copy\_X’ on the Sklearn’s Linear Regression module but didn’t see any difference, please refer below the parameters:

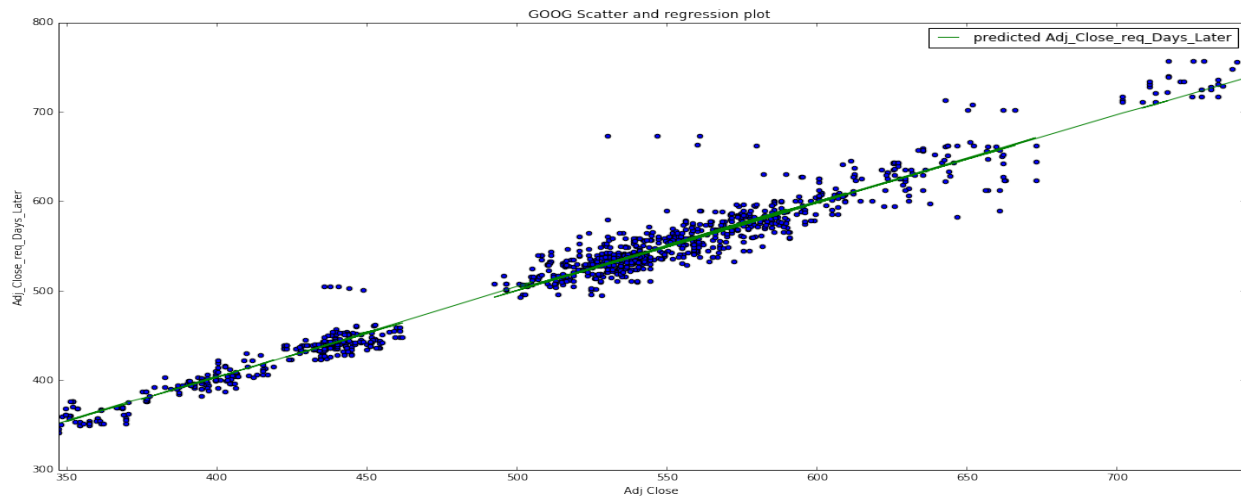
```
params = {'fit_intercept':[True, False],
          'normalize' :[True, False],
          'copy_X': [True, False],
          'n_jobs':[-1]}
```

I printing out the training and test scores on the stocks, and the % +/- variation of the predicted stock value compared to the actual values. Additionally, I am also plotting the scatter plot between ‘Adj Close’ and ‘Adj\_Close\_req\_Days\_laters’ and the regression fitted line on the predicted ‘Adj\_Close\_req\_Days\_laters’. For e.g:

```
Score on training data for GOOG :0.97
Score on test data for GOOG : 0.89
```

```
The prediction on the test dataset 5 days after the training date for GOOG is +/- 2.499% compared to the actual values
```

Scatter plot between ‘Adj Close’ and ‘Adj\_Close\_req\_Days\_laters’ and the regression fitted line on the predicted ‘Adj\_Close\_req\_Days\_laters’ for Google



## KNN REGRESSION:

The steps for KNN Regression is ditto to Linear Regression except that I am using SKLEARN's KNN Regression with default parameters:

```
class sklearn.neighbors. KNeighborsRegressor (n_neighbors=5, weights='uniform', algorithm='auto',
leaf_size=30, p=2, metric='minkowski', metric_params=None, n_jobs=1, **kwargs)
```

Additionally, since the R2 scores on test data for most of the stocks using the default parameters are terrible, I thought of using gridsearch to find the best tune parameters for KNN Regressor, please refer below the code snippet for hyper parameter tuning:

```
params = {'n_neighbors': range(1, 50, 5),
          'leaf_size': range(1, 50, 5),
          'algorithm' : ['auto'] ,
          'weights': ['uniform', 'distance']}

scorer = make_scorer(r2_score)

regr = train_classifier_gridsearch(regr, params, scorer, X_train, y_train)
```

R2 Score on Google to predict the date which is one day ahead in future compared to the last training date without gridsearch:

```
Score on training data for GOOG :0.76
Score on test data for GOOG : -21.84
```

```
The prediction on the test dataset 1 days after the training date for GOOG is +/- 10.138% compared to the actual v
alues
```

After using grid search as described above, astonishingly my R2score on test data is even poor but the score on trained data has improved a lot, please refer below. This typical sign of overfitting. So, overall I would think KNN is a bad algorithm for Stock price prediction:

# Capstone Project on Investment and Trading: Build a Stock Price Indicator Machine Learning Engineer Nanodegree

Partha Deka

Feb 12<sup>th</sup>, 2017

Score on training data for GOOG :0.98  
Score on test data for GOOG : -21.83  
The prediction on the test dataset 1 days after the training date for GOOG is +/- 10.111% compared to the actual values

**ARIMA:** Please refer the “Algorithms and techniques” section of this document for an explanation of ARIMA:

I have tried ARIMA with APPLE stock to predict the stock price one day ahead in future, please refer the function “`exec_arima`” in the header below in the “All\_Analysis.ipynb” notebook:

## # HEADER 6

### - STOCK PRICE PREDICTION INTERFACE: ARIMA”

Please note the libraries imported:

```
import pandas as pd
import numpy as np
from statsmodels.tsa.arima_model import ARIMA, ARIMAResults
#from statsmodels.tsa.stattools import acf, pacf
from statsmodels.graphics.tsaplots import plot_acf, plot_pacf
import matplotlib.pyplot as plt
#import matplotlib.pyplot as plt
import matplotlib.dates as dates
from sklearn.metrics import r2_score
```

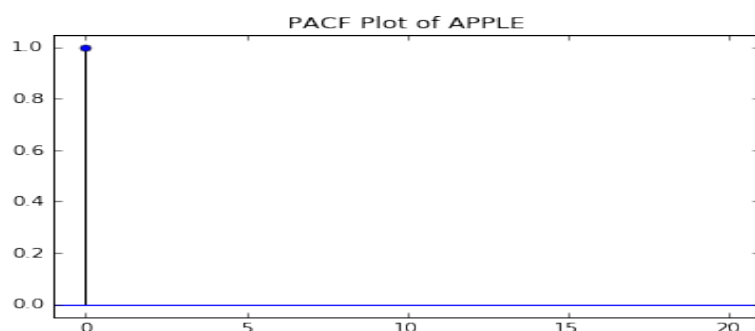
Please note for ARIMA I am just using ‘ADJ CLOSE’ as the attribute and since I am trying to predict the stock price one ahead I am adjusting the ‘ADJ CLOSE’ as below:

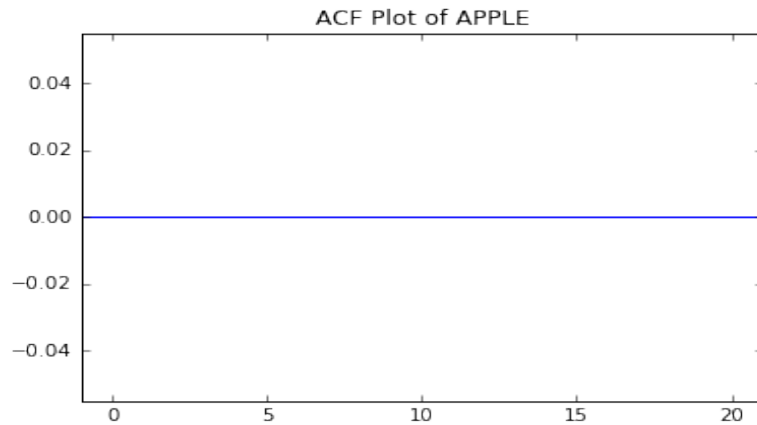
```
df['Adj Close'] = df['Adj Close'].shift(-1)
```

I started with diagnosing the AR, I and MA terms for the ARIMA model, please refer below the code snippet as well the ACF and PACF plots:

```
##DIAGNOSING THE ACF AND PACF PLOTS
acf = plot_acf(df, lags = 20)
plt.title("ACF Plot of APPLE")
acf.show()

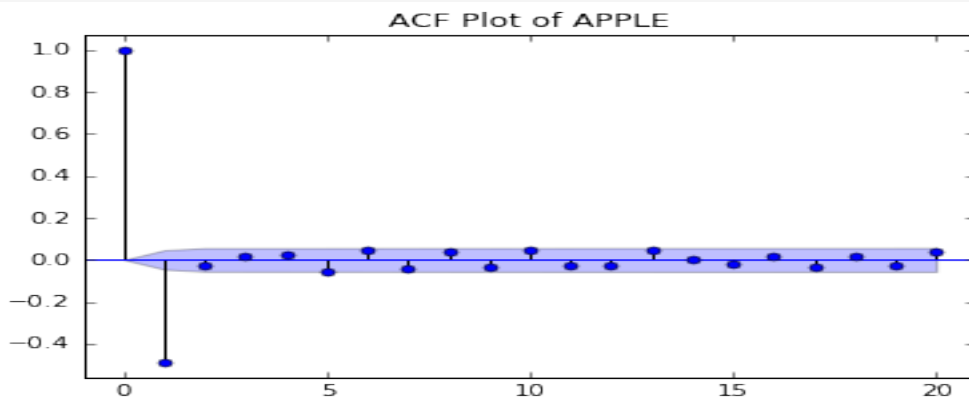
pacf = plot_pacf(df, lags = 20)
plt.title("PACF Plot of APPLE")
pacf.show()
```

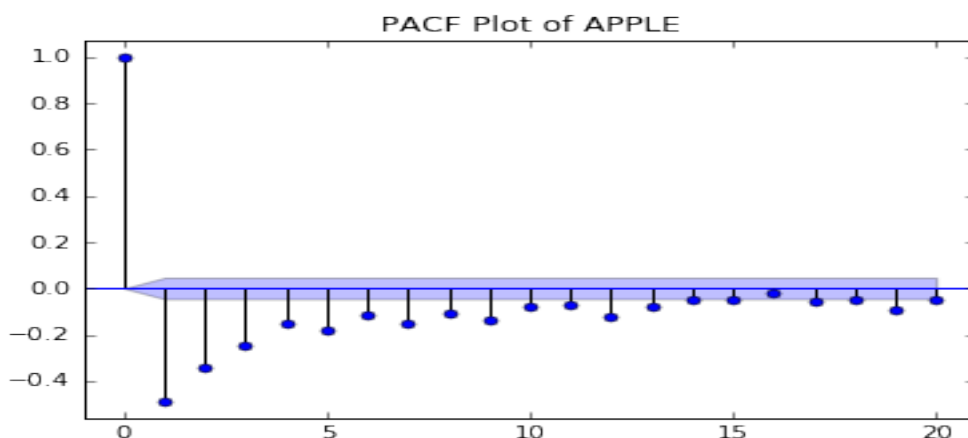




The ACF plot above suggests there is no auto-correlation, so I decided to do 1 lag differencing which is also, not helping much and went ahead with 2 lag differencing. Please refer the code snippet as well the ACF and PACF plots below:

```
#The ACF graph suggest there is constant auto co-relation. Lets do 1 lag differencing and dig further to confirm th  
#1 lag differencing is not helping much and so decided to do 2 lag differencing  
  
df_diff = df - df.shift()  
df_diff = df_diff - df_diff.shift()  
df_diff.fillna(method="ffill", inplace="True") # Forward fill empty trade dates  
df_diff.fillna(method="bfill", inplace="True") # backfill empty trade dates  
  
##DIAGNOSING THE ACF AND PACF PLOTS for DIFF  
acf = plot_acf(df_diff, lags = 20)  
plt.title("ACF Plot of APPLE")  
acf.show()  
  
pacf = plot_pacf(df_diff, lags = 20)  
plt.title("PACF Plot of APPLE")  
pacf.show()
```





As we can see, the PACF plot is slowly tapering towards '0'. So it's basically a MA model based on the ACF plot with  $q=2$  as the non-zero auto correlation disappears after lag 2. So, this ARIMA has  $p=0$ ,  $d=2$ ,  $q=1$ .

Please refer the code snippet below to observe how I kept a slice of the dataset for training, predicted on a slice of the dataset part of which is also include in the training set. This would help me observe how the ARIMA model performs in both train and test dataset:

```
#Building the ARIMA model

df_train = df['2016-09-01':'2016-11-21']

mod = ARIMA(df_train, order = (0, 2, 1))

results = mod.fit()

predVals = results.predict('2016-10-01', '2016-12-01', typ='levels')

google_pred = pd.concat([df['2014-05-01':'2016-12-01'], predVals], axis = 1)

google_pred = google_pred.rename(columns={'Adj Close':'Original',0:'Predicted'})

google_pred.fillna(0, inplace=True)

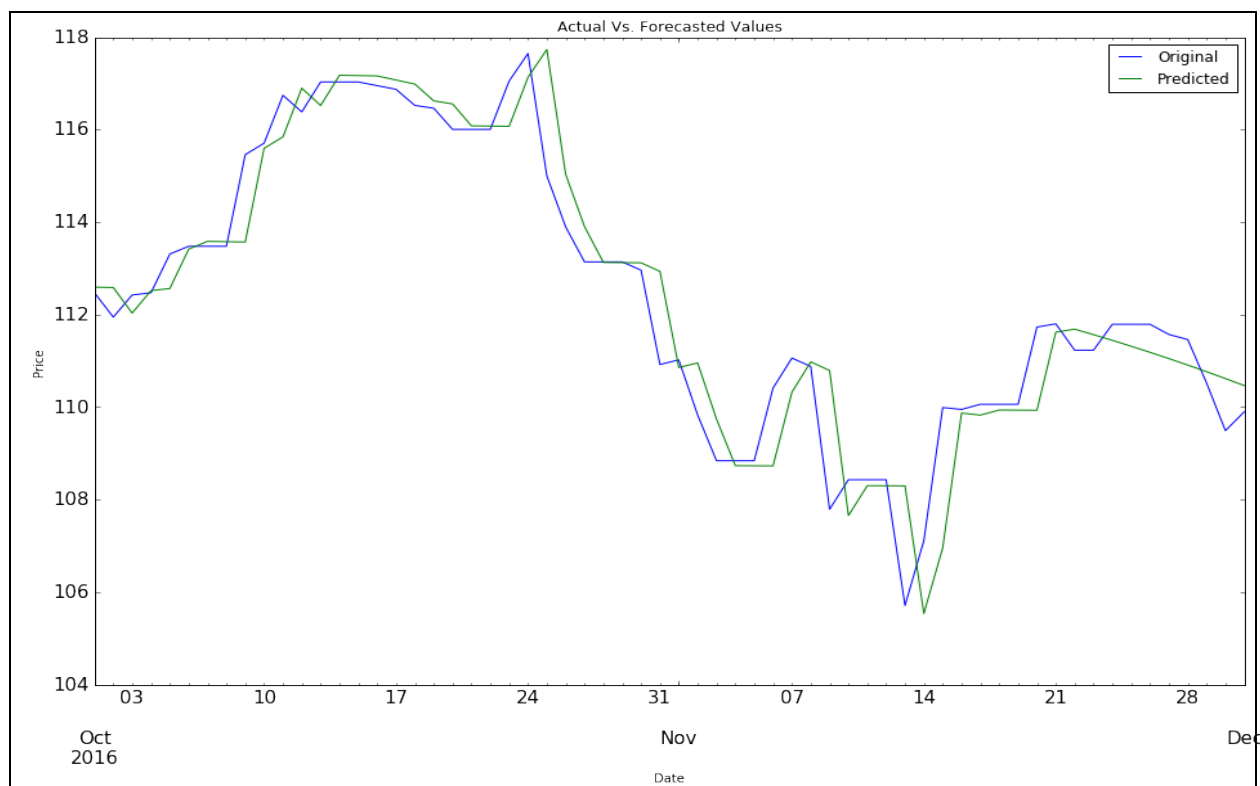
##plot actual vs predicted

plot_selected(google_pred,['Original','Predicted'], '2016-05-01', '2016-12-01',title="Actual Vs. Forecasted Values")

#Calculate R squared on train and test datasets

r2_score_train = r2_score(google_pred['Original']['2016-10-01':'2016-11-21'], google_pred['Predicted']['2016-10-01':
print "Train score ", r2_score_train

r2_score_test = r2_score(google_pred['Original']['2016-11-22':'2016-12-01'], google_pred['Predicted']['2016-11-22':
print "Test score ", r2_score_test
```



Please observe the plot above - the 'green' plot almost perfectly follows the original 'blue' plot for the trained dataset. Towards the end of the plot, please observe how the green plot follows the 'blue' plot for the test dataset.

The train score: 0.887920343472

The test score: 0.472620109556

Sources:

Udacity Free Course: Time Series Forecasting

<https://onlinecourses.science.psu.edu/stat510/node/62>

<https://onlinecourses.science.psu.edu/stat510/node/49>

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

<http://www.inertia7.com/projects/time-series-stock-market-python>

### Refinement:

The process of improving upon the algorithms and techniques used is clearly documented. Both the initial and final solutions are reported, along with intermediate solutions, if necessary.

Among the three algorithms, the "Linear Regression" model has the best score. I tried to improve the R2 score by performing hyper parameter tuning on the parameters, but I am getting same scores.

Scores without grid search:

Score on training data for GOOG :0.97

Score on test data for GOOG: 0.89

Score with grid search:

Score on training data for GOOG :0.97

Score on test data for GOOG: 0.89

This observation is expected for Linear Regression as there not parameters to tune. Please refer the code snippet below:

```
params = {'fit_intercept':[True, False],
          'normalize' :[True, False],
          'copy_X': [True, False],
          'n_jobs':[-1]}

scorer = make_scorer(r2_score)

regr = gridsearch(regr, params, scorer, X_train, y_train)
```

As discussed above in the implementation section “KNN regression” model is terrible (score on test data is negative), gridsearch for hyper parameter tuning didn’t help at all.

For ARIMA (please refer the ARIMA code), I got a score of 0.47 for ‘APPL’ to predict the stock price one day ahead in future. This score is below the benchmark score of 0.6

## **4.Results: Model Evaluation and Justification:**

By far, based on R2 score and +/- % average prediction accuracy compared to actual values, I found the **Linear Regression** model is much better off compared to KNN Regression for all the stocks I am chosen for this project. Although the Linear Regression model’s accuracy decreases when the number of days from the last training date to the date to be predicted increases, up to **10 days in future**, the R2 score on trained datasets for all the stocks taken into consideration is more or less 0.9 and R2 score on the test datasets is more or less 0.6. On an average the predicted prices are also **more or less within +/- 5%** of the actual stock prices up to **10 days** stock price prediction in future.

Non-Seasonal ARIMA showed some promising results, for Apple Stock its R2 score on the train dataset **1 day** out in future is **0.88** but the score on the unseen test dataset is **0.47** and missed the benchmark of 0.6 on test dataset. I could have tried out the Season ARIMA to bring in seasonality into consideration but I would keep it for a later time.

## **5. Conclusion: Free form visualization, Reflection and Improvement:**

As discussed in the Data Exploratory and Exploratory visualization section, the following two visualizations on: normalized historical Portfolio Data vs normalized SNP\_500 index seems to be very interesting for me. The first one is without portfolio optimization but equal allocation to all stocks- Fund allocations by symbol: symbols = ['AAPL', 'GOOG', 'AMZN', 'FACEBOOK', 'GE', 'MICROSOFT'], allocation = [0.166, 0.166, 0.166, 0.166, 0.166, 0.166]

The second is based on allocation using Sharpe ratio optimization: Fund allocations by symbol: symbols = ['AAPL', 'GOOG', 'AMZN', 'FACEBOOK', 'GE', 'MICROSOFT'], allocation = [0.0, 0.0, 0.914, 0.086, 0.0, 0.0]

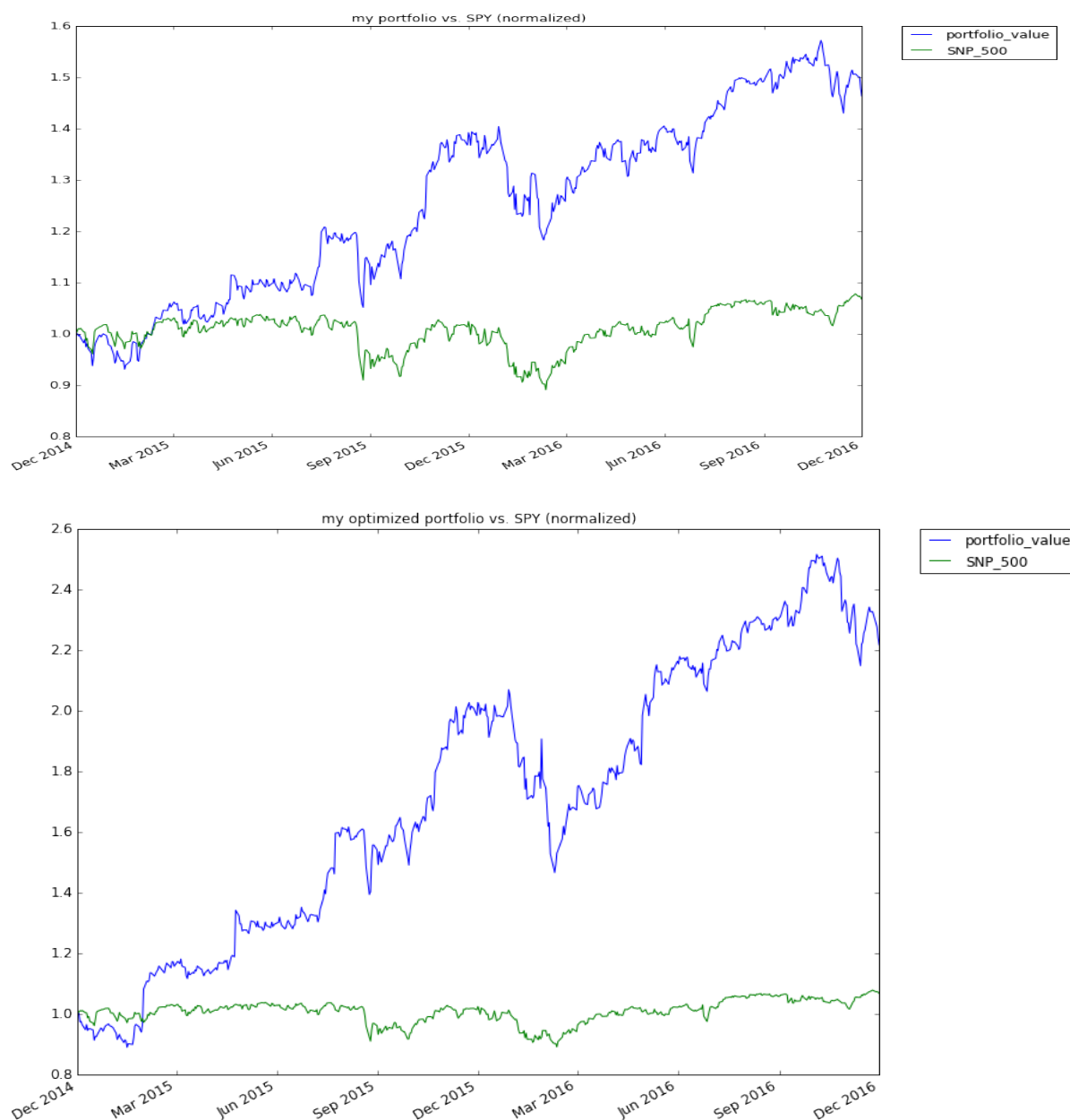
It’s evident that performance of the optimized portfolio is better off compared to the equally allocated portfolio:



**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**



Overall, I enjoyed throughout this project, there are two interesting aspects I enjoyed the most – i) the financial data analysis and its implementation in Python which I learnt throughout the course “Machine Learning for Trading” ii) I got the chance to get my hands dirty on Time Series Forecasting. I learnt techniques like ARIMA and exponential smoothing. Courtesy: Udacity course: “Time Series Forecasting”

I believe stock price prediction is a mature and old problem, I am certain many companies have built mature machine learning models and operationalized them on real-time software engineering infrastructure.

**Capstone Project on Investment and Trading:  
Build a Stock Price Indicator  
Machine Learning Engineer Nanodegree**

**Partha Deka**

**Feb 12<sup>th</sup>, 2017**

For me, to build a better model in future so that I can predict stock prices accurately for months ahead and not just days, I am thinking of using **Exponential Smoothing** as discussed in the “**Time Series Forecasting**” by apply **error, trend** and **seasonality** in the smoothing method calculation. I feel stock prediction is more of a time series prediction problem than just a regression problem, each stock would have its own trend & seasonality. I also have plans to implement an **ensembled** learning technique which mostly gives better results than a single learner because it cancels the biases of the individual learners. I also have great enthusiasm to use the techniques of **Reinforcement** learning, as there many other external parameters which a smart self-learned model would want to learn from – such as current news on economic, political, geographical, climatic situations around the country, world etc. which affects the stock market. I believe a real time self-learning robotic system which take inputs from its environment is more relevant in the world of stock price prediction than an algorithm which only learns from history and predicts the future. I am starting with the “**Self- Driving Car**” Nano degree beginning 16<sup>th</sup> of Feb, I am confident that I would gain good expertise in **Deep learning** and **Robotics** after completing the Nanodegree, which I can definitely apply in the stock price prediction problem as well.