

Data 698 - Final presentation

Tze Fung Lung, Jim

# Topic: **Portfolio optimization and Machine learning with visualization analysis for S&P 500**



# Introduction

- We are looking at the S&P 500, an index of the largest US companies. The S&P 500 is an American stock market index based on the market capitalization of 500 large companies having common stock listed on the NYSE, NASDAQ Exchange.
- I loaded all 500 dataset in S&P 500 for analysis by using portfolio optimization to get the possible several stocks with higher return and lower risk. And using the machine learning predict the investment trend for S&P 500 index.

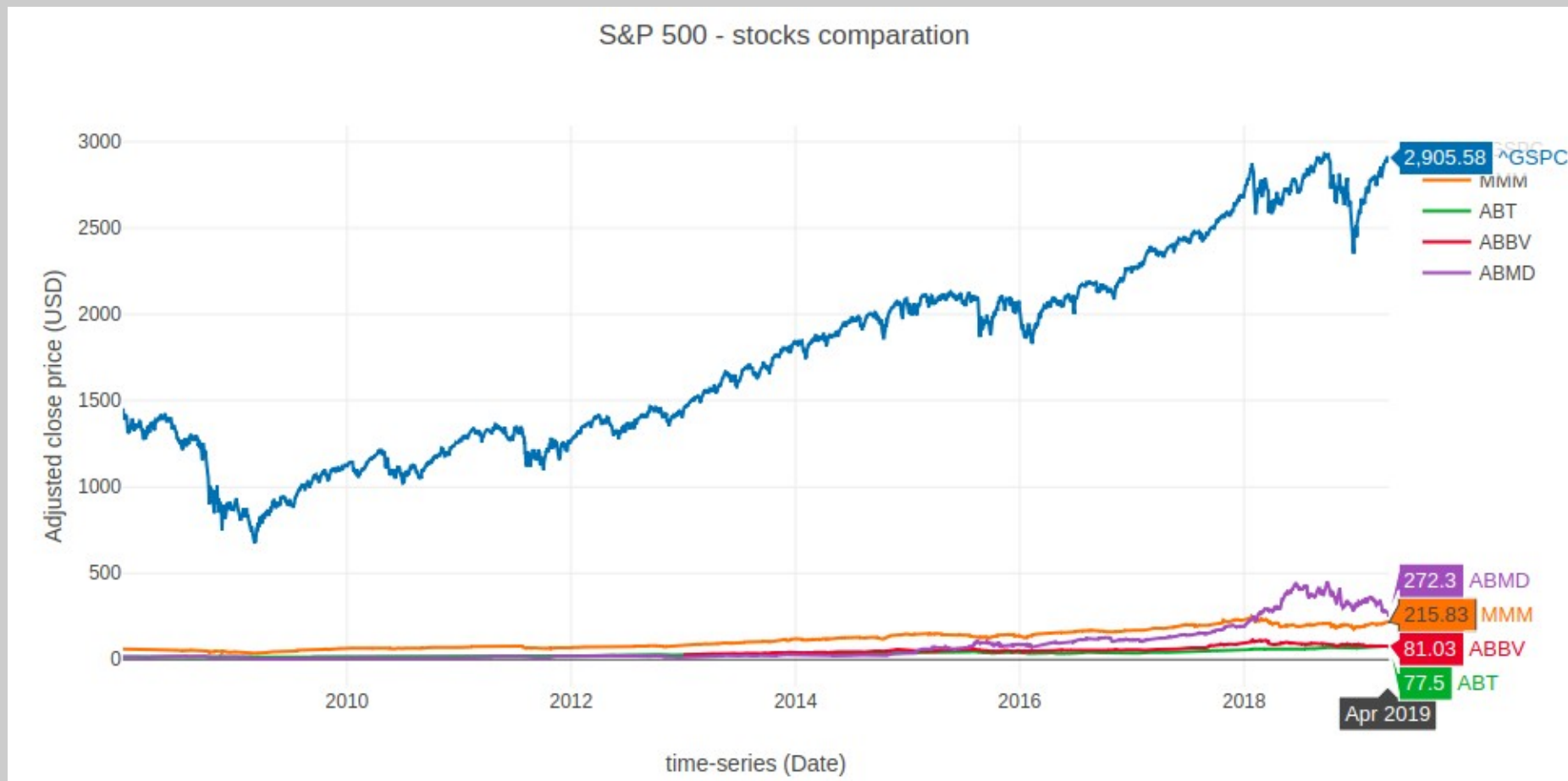
# Methodology

- This project are separated parts of analysis from data exploration, visualization, correlation and monthly return for data extraction by mathematical programming, portfolio optimization and machine learning .
- 1. Data Exploration
- 2. Higher monthly return
- 3. Portfolio Optimization
- 4. Machine Learning

# Data Exploration

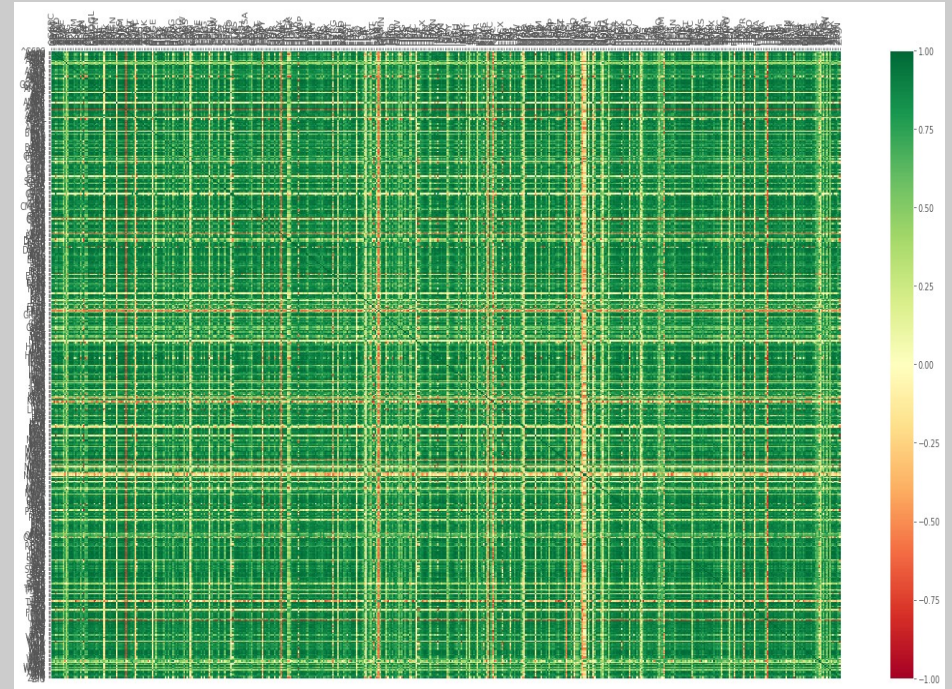
- We considered if it plot all 500 stock shares, it can't appear a meaningful comparison graph. Initially, we view the first 5 column including S&P 500 index in the time series from 2008 to 2019 (2019-04-18).
- We will further see the top 10 monthly return, top 10 for higher correlation with S&P 500 index, and top 10 investing stocks strategic from the portfolio optimization in our data time series.

# Data Exploration



# Correlation for all 500 stocks

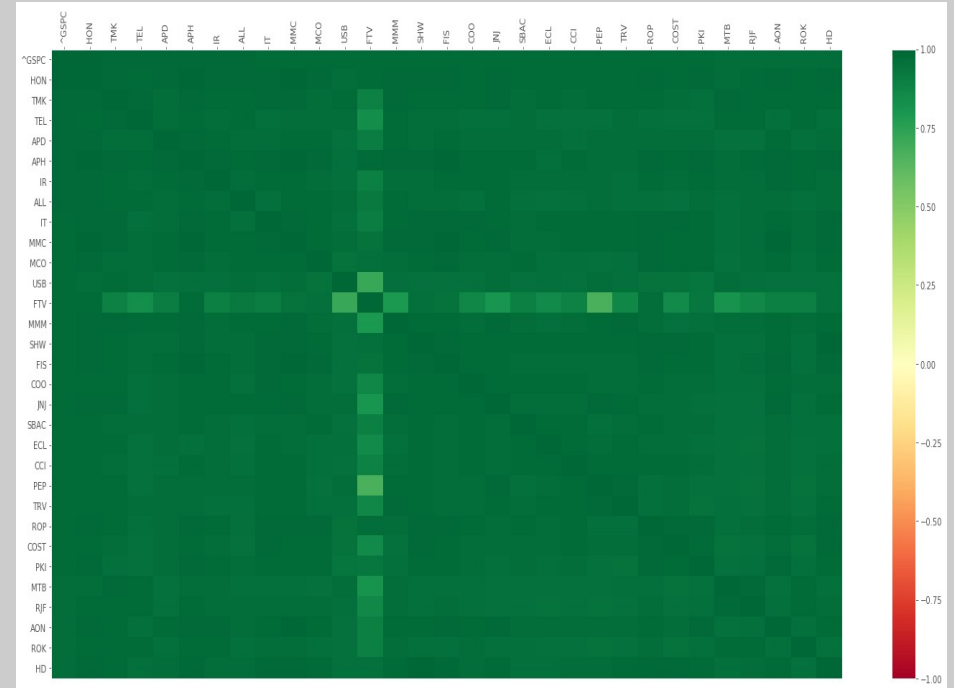
Most people think that they will get higher return when the main index going up, so we can calculate the correlation to check which stocks share are strong relationship and low relationship.





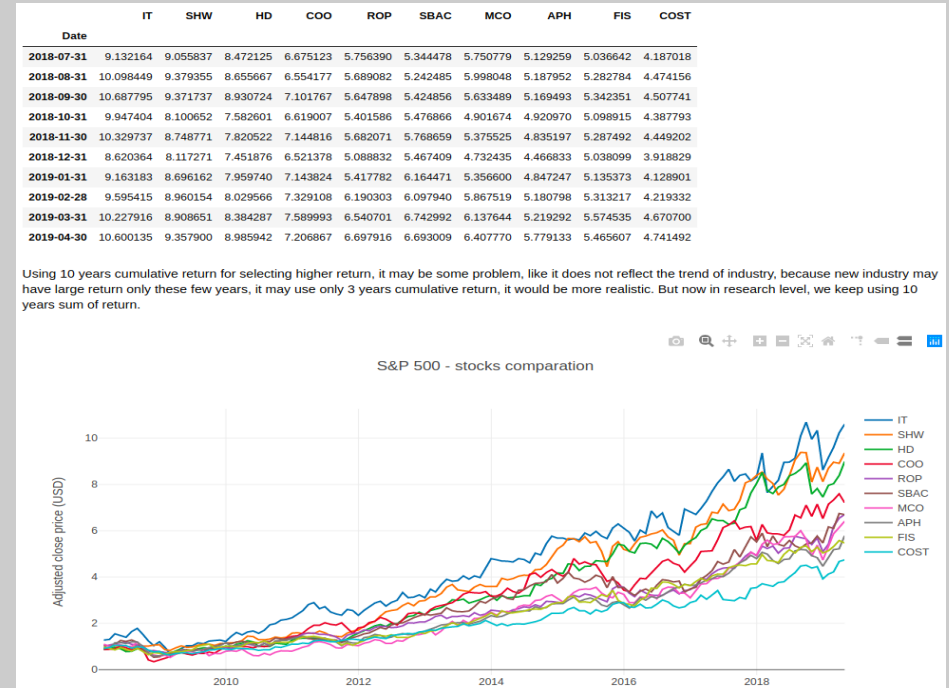
# Highest Correlation (Sort out top 30 from all 500 stocks share)

- to find out the top 30 stocks which are higher correlation with S&P 500 index from 500 stocks share.
- Last part of this project will conduct the prediction of machine learning for S&P 500 index



# Highest monthly return (Sort out top 10 from 30 stocks share)

- Apart from the risk consideration, the percentage of increasing return can be 9.8 times over 10 year if just only invest "IT" tickers of stock.



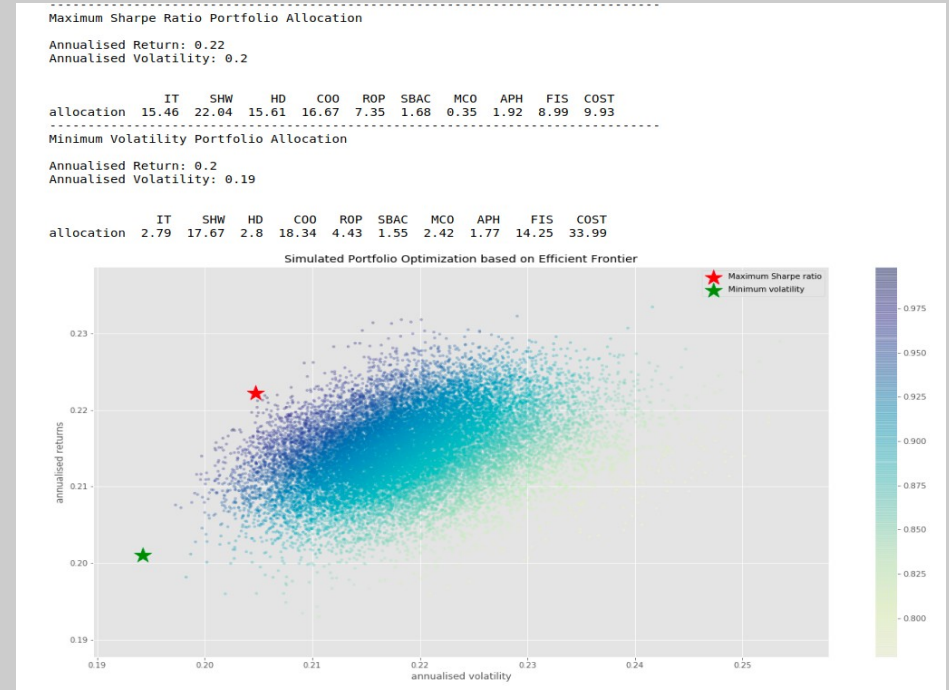


# Highest monthly return (Sort out top 10 from 30 stocks share)

- Gartner Inc
- Sherwin-Williams
- Home Depot
- The Cooper Companies
- Roper Technologies
- SBA Communications
- Moody's Corp
- Amphenol Corp
- Fidelity National Information Services
- Costco Wholesale Corp.

# Portfolio Optimization - Random Portfolios

- In order to choose from 10 number of stocks share from the result of last part which are higher monthly return and lower risk as our final medium or long term investment strategic.

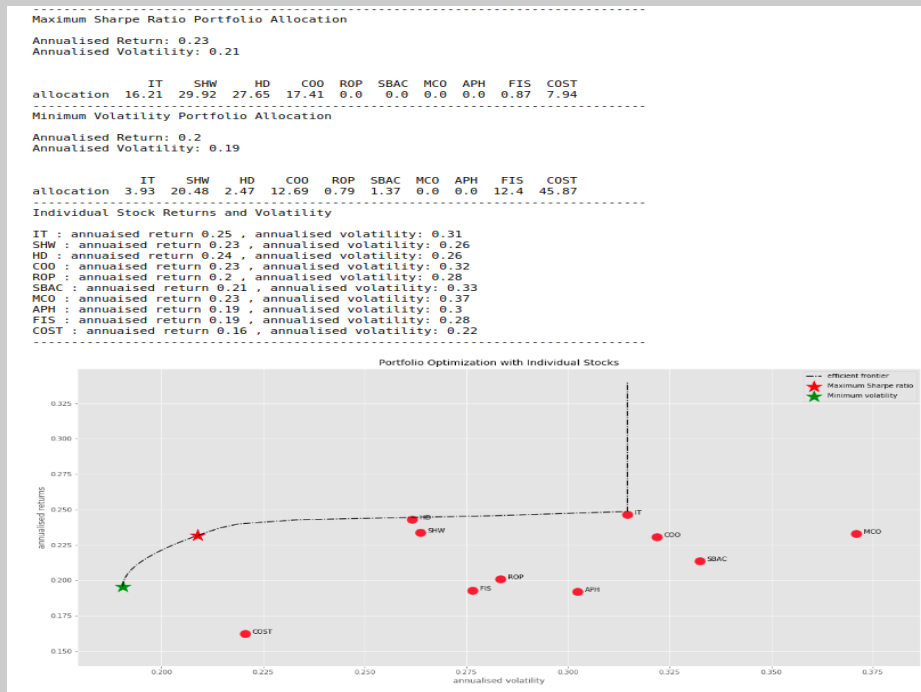


# Portfolio Optimization - Random Portfolios

- For minimum risk portfolio, we can see around 30% of our budget is allocated to "Cost" - Costco. If you take another look at the daily return plot from earlier, the Costco is the least volatile among these stocks, so allocating a large percentage to Costco for minimum risk portfolio makes sense.
- In this scenario, we are allocating a significant portion to "SHW" - Sherwin-Williams and "IT" - Gartner Inc, which are quite volatile stocks from the previous plot of daily returns. And "Cost" - Costco which had around 30% in the case of minimum risk portfolio, has only 10% budget allocated to it.

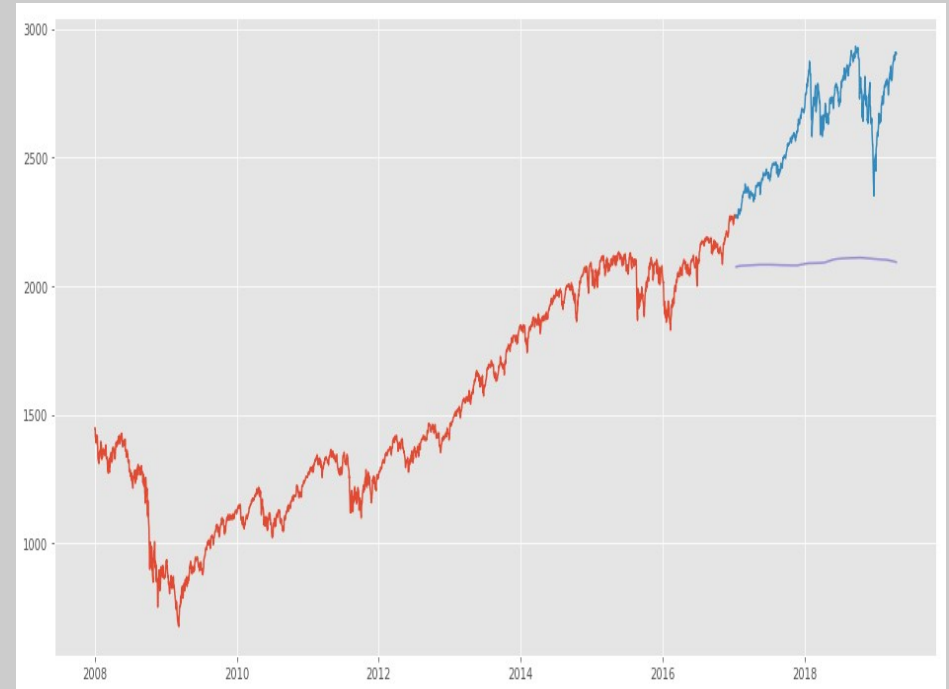
# Portfolio Optimization - Efficient Frontier

- Stocks with the least risk is COST at around 0.22, but the return is only around 0.16. If we will to take slightly more risk around 0.225 return, we should consider to choose HD and SHW rather than IT with higher risk with portfolio optimisation.



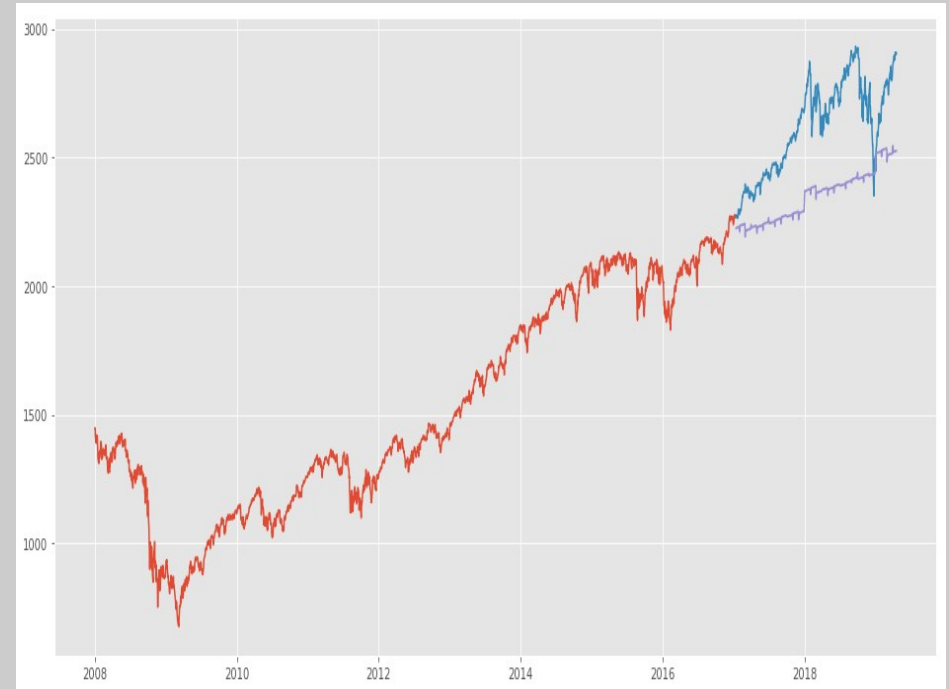
# Machine Learning - Moving Average

- The root mean square error is 556.7202, prediction is 2092.7282
- the results are not very promising, the predicted values are not the same range as the observed values in the train set
- (there is not obvious an increasing trend and then a slow decrease)



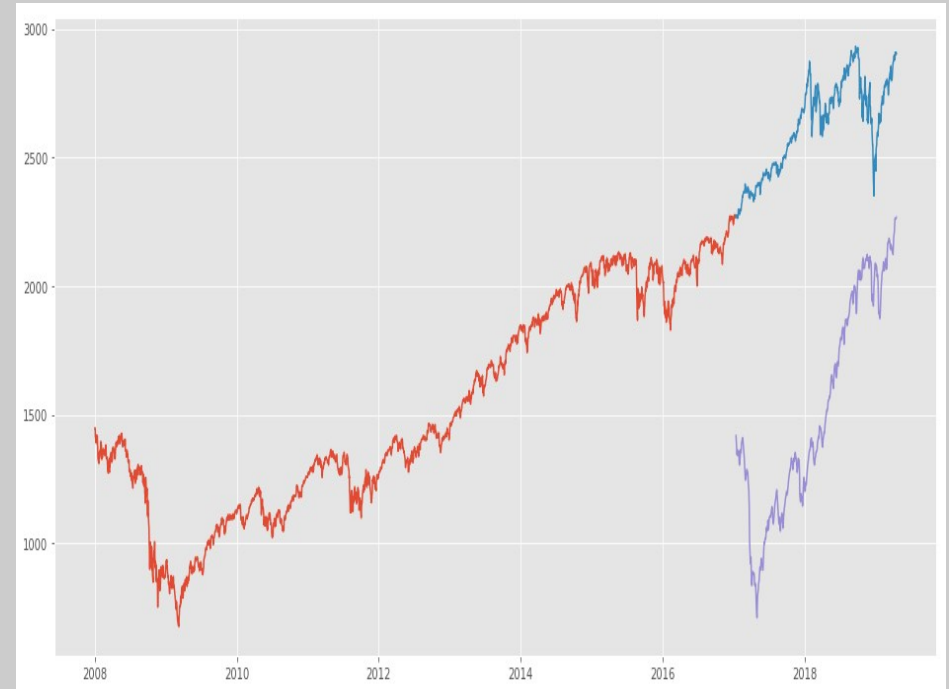
# Machine Learning - Linear Regression

- The root mean square error is 296.1361, prediction is 2525.14814
- The RMSE value is lower than the previous technique, which clearly shows that linear regression has performed better.



# Machine Learning - k-Nearest Neighbours

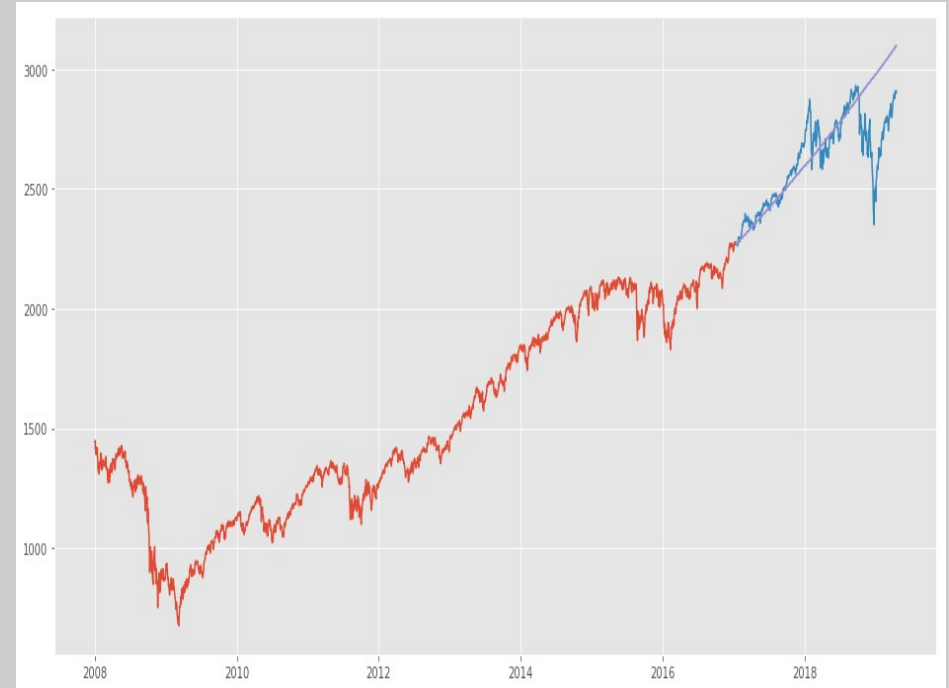
- The root mean square error is 1142.9905, prediction is 2266.327
- The RMSE value of KNN is higher than the linear regression model and the plot shows the same pattern
- kNN also identified a raising trend from January 2019 since that has been the pattern for the past years





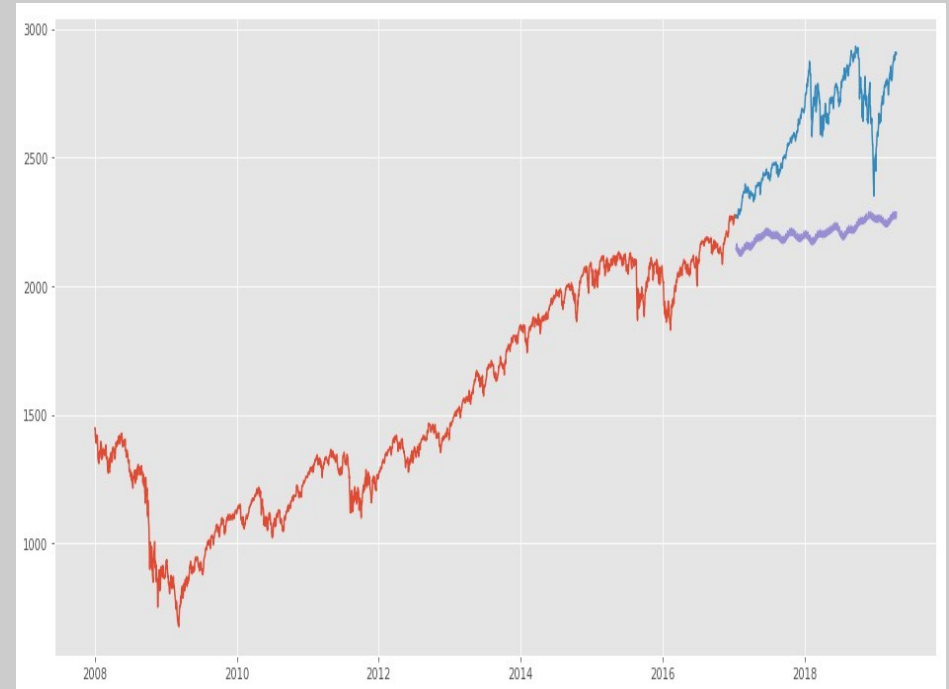
# Machine Learning - Auto ARIMA

- The root mean square error is 149.4734, prediction is 3098.4159
- Although the predictions using this technique are far better than that of the previously implemented machine learning models, these predictions are still not close to the real values.



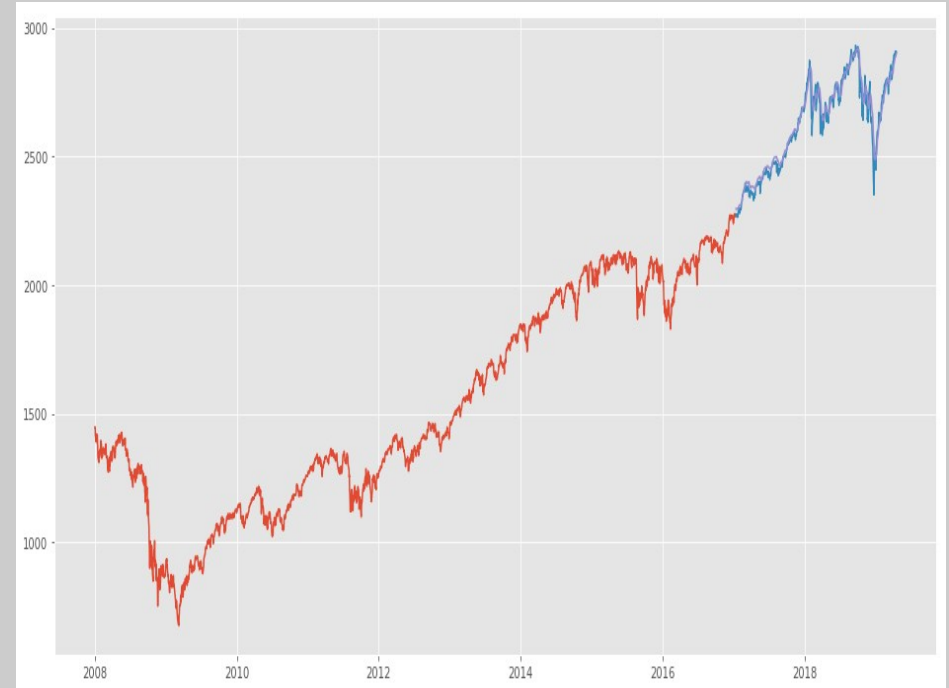
# Machine Learning - Prophet

- The root mean square error is 446.4384233932304, prediction is 2285.3302
- Forecasting techniques like ARIMA and Prophet would not show good results for this particular problem.



# Machine Learning - Long Short Term Memory (LSTM)

- The root mean square error is 36.8895, prediction is 2899.4375
- LSTMs are widely used for sequence prediction problems and have proven to be extremely effective.
- the model seem to predict well with real data



# Model Comparsion

- The S&P 500 index is 2905.030029 on 2019-04-18
- the most lowest rmse of model is by using of Long Short Term Memory method, prediction value is 2899.4375 As we mentioned as last part, we should also consider the news about the company and other factors like demonetization or merger/demerger of the companies.

	RMSE	Prediction
ma	556.720265	2092.728205
linear	296.136189	2525.148140
knn	1142.990569	2266.327772
arima	149.473499	3098.415958
prophet	446.438423	2285.330298
lstm	36.889548	2899.437500

The last S&P 500 index is 2905.030029