



**STEVENS**  
INSTITUTE *of* TECHNOLOGY  
THE INNOVATION UNIVERSITY®

# **FE-582 Foundations of Financial Data Science**

**Fall 2019  
Final Project Report**

**Stock analysis for 30 Dow Jones Stocks**

**1870**  
**Junzhe Wang  
Ashish Negi  
Rohnit Shetty**



## Introduction

The enormous availability of data in today's time enables investors, at any scale to make better investment decisions. Keeping such a theory in mind, our goal is to use historical price of 30 Dow Jones stocks to analyze the stock market and predict the performance of stocks.

## Research Question

The main question we wish to find an answer is, if we can use the content of news analytics to predict stock price performance?

## Dataset Description

We chose the dataset which provides free end of day data for all stocks currently in the Dow Jones Industrial Average, from 2015-01-02 to 2019-06-11. There are 30 csv files in the current version of the dataset. For each of the 30 components of the index, there is one CSV file named by the stock's symbol (e.g. AAPL for Apple). Each file provides historically adjusted market-wide data (daily, max. 5 years back).

List of stocks and symbols as per:

[https://en.wikipedia.org/wiki/Dow\\_Jones\\_Industrial\\_Average](https://en.wikipedia.org/wiki/Dow_Jones_Industrial_Average)

---



5 rows  $\times$  30 columns

5 rows x 30 columns

Raw data sample - Tail



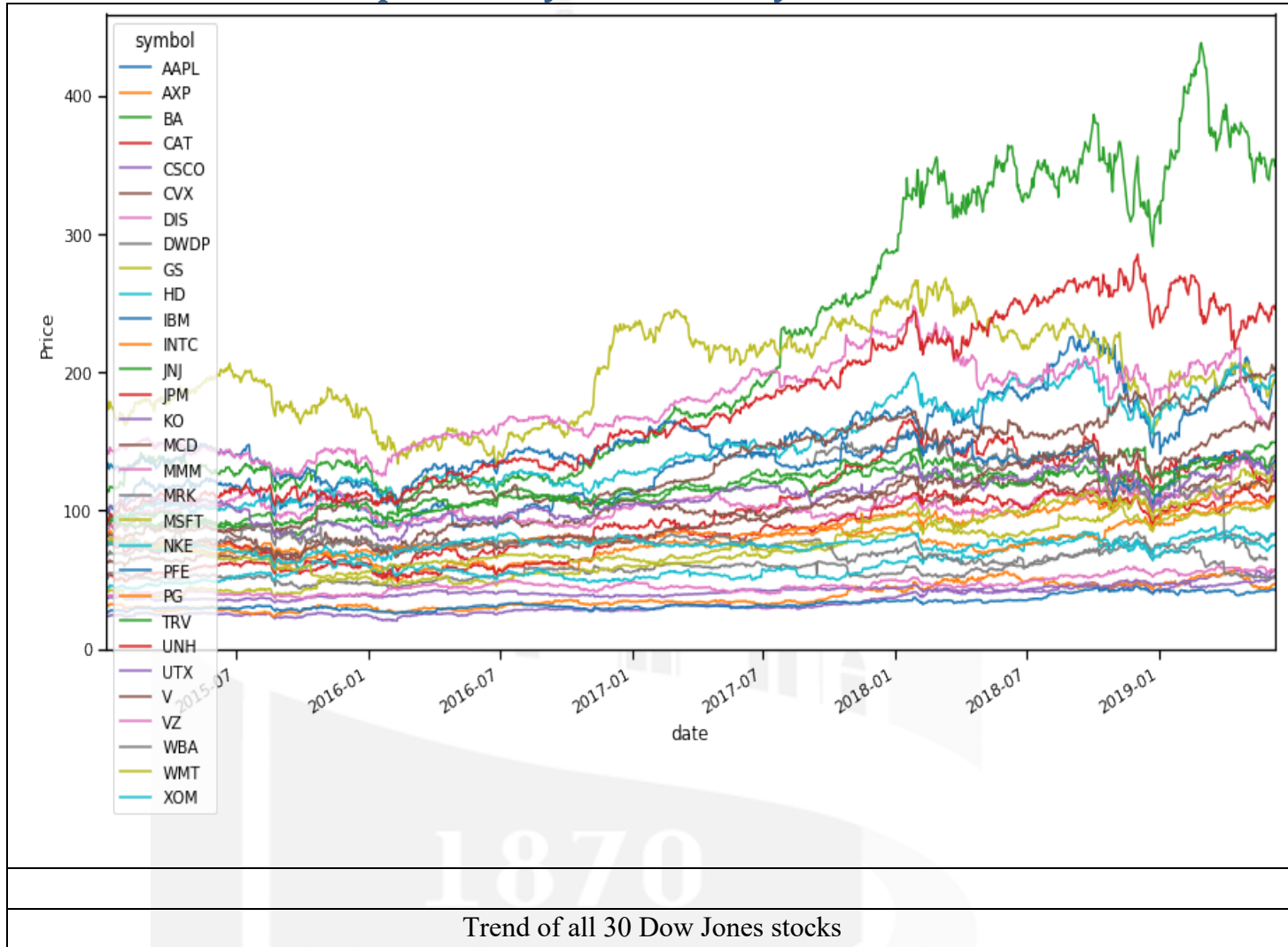
## Objectives and methods of Exploratory Data Analysis

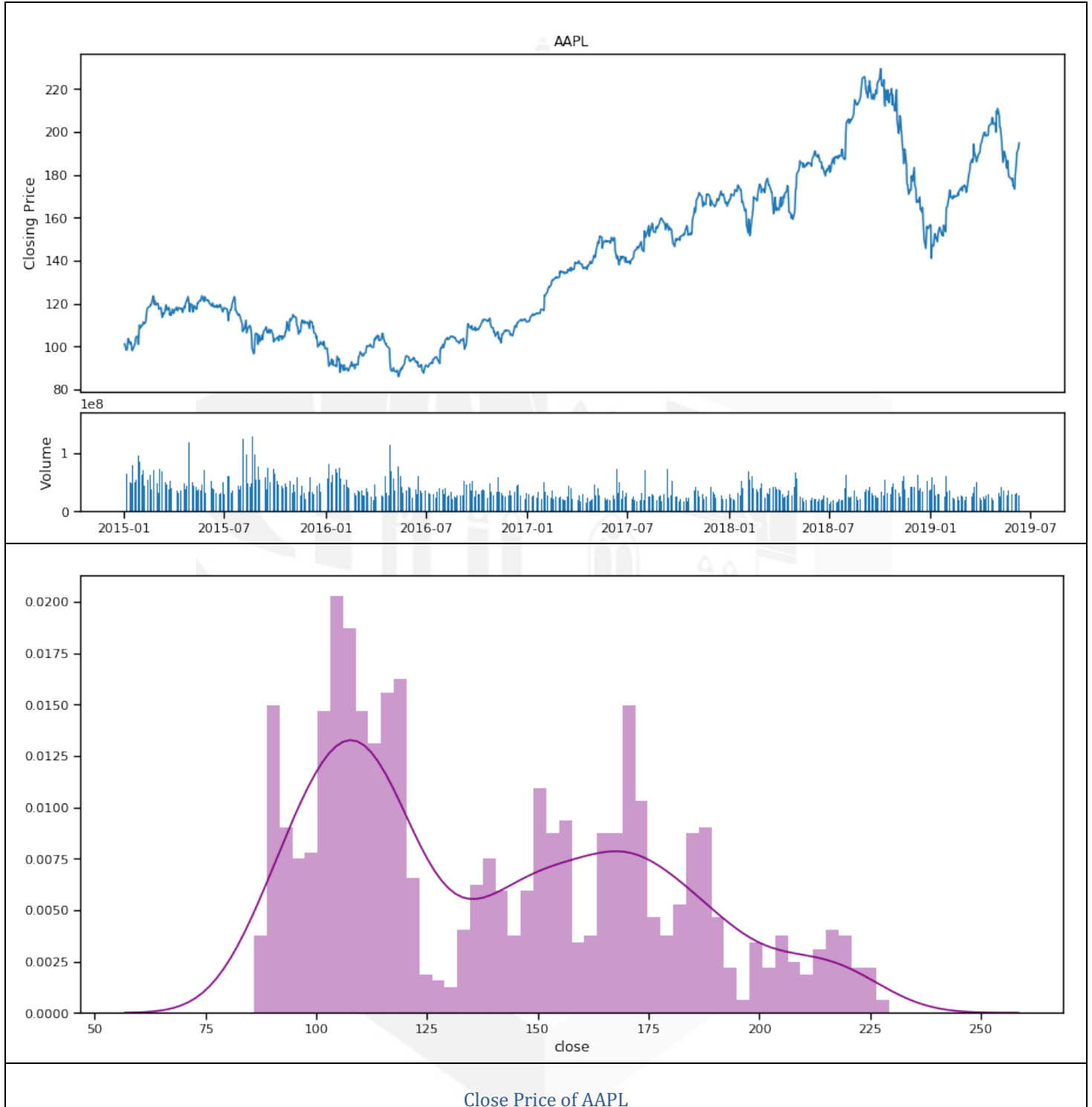
The objectives of the EDA are as follows:

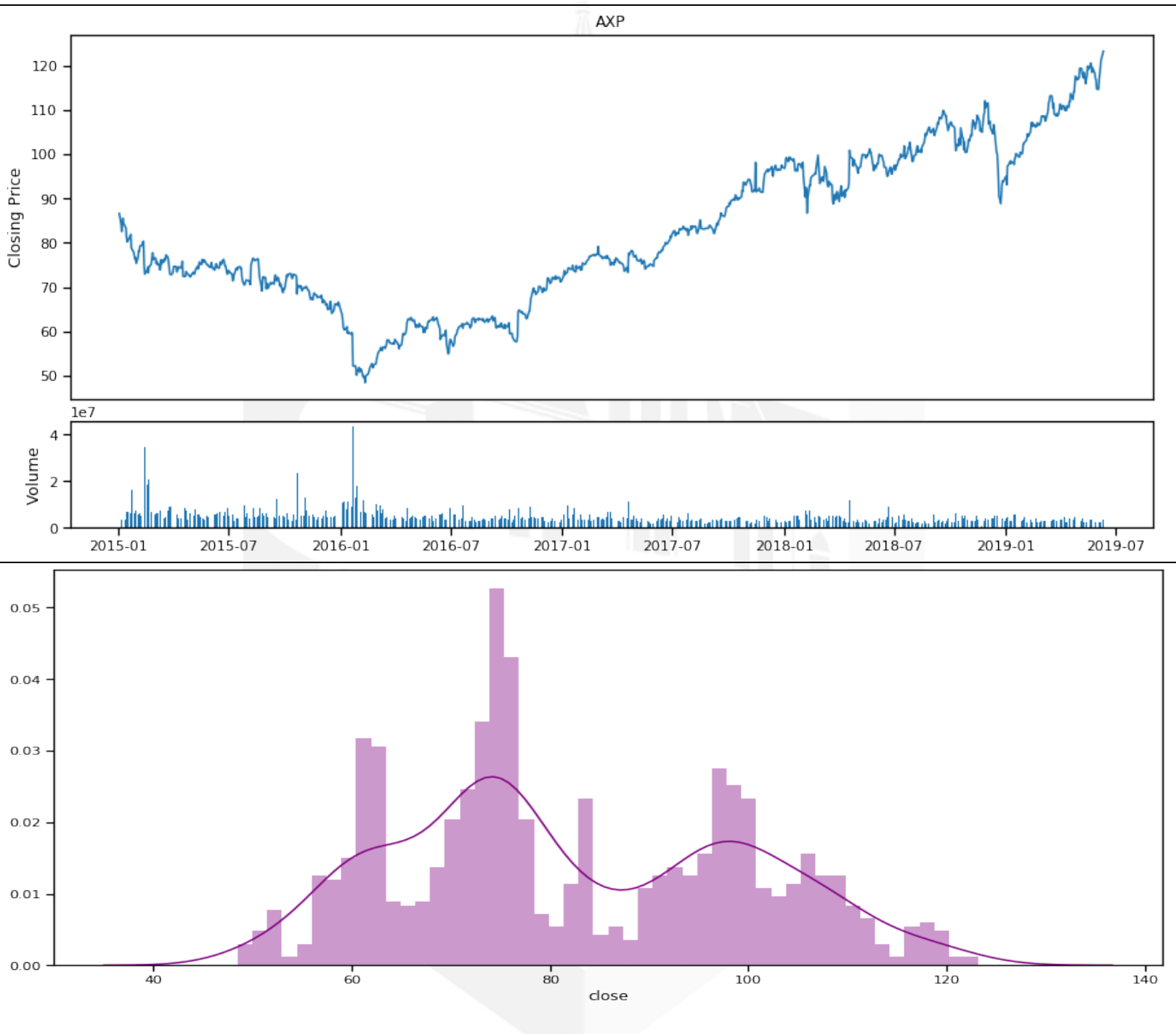
1. To get an overview of the data and distribution of the dataset, including data type, the size(shape) of the whole data frame, the highest price, the lowest price, the volume and the mean price by different window size in two months or 1 year(41 or 252 trading day).
2. Check for missing numerical values, outliers or other anomalies in the dataset. If there is any missing data, we can fill forward or remove the invalid data.
3. Discover patterns and relationships between variables in the dataset. We can solve this question by ranking the correlation, or visualizing the correlation by heap map.
4. Check the underlying assumptions in the dataset. In this case, we hope to come up with a trading strategy and back test with the dataset.

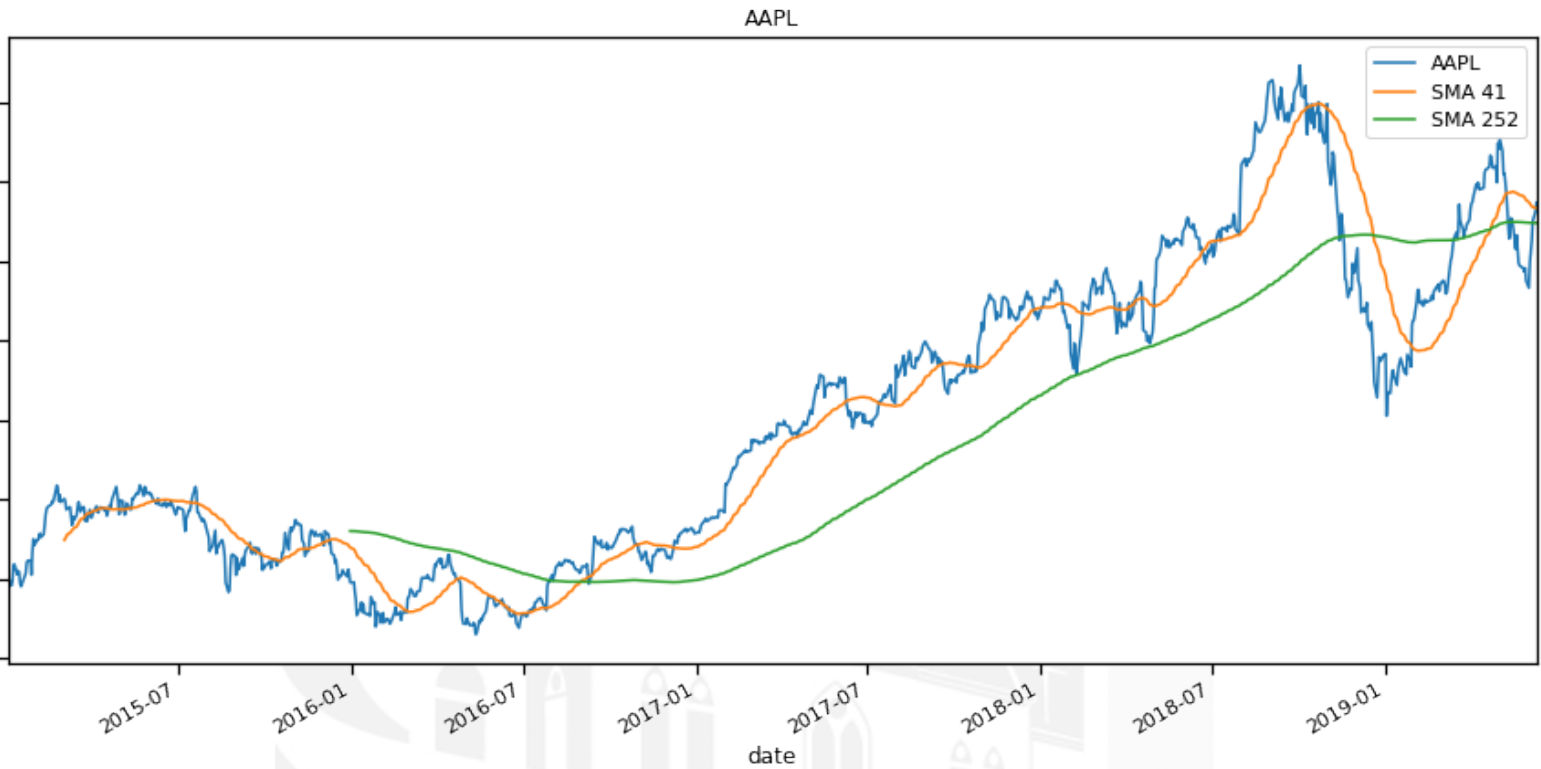


## Some results of Exploratory Data Analysis





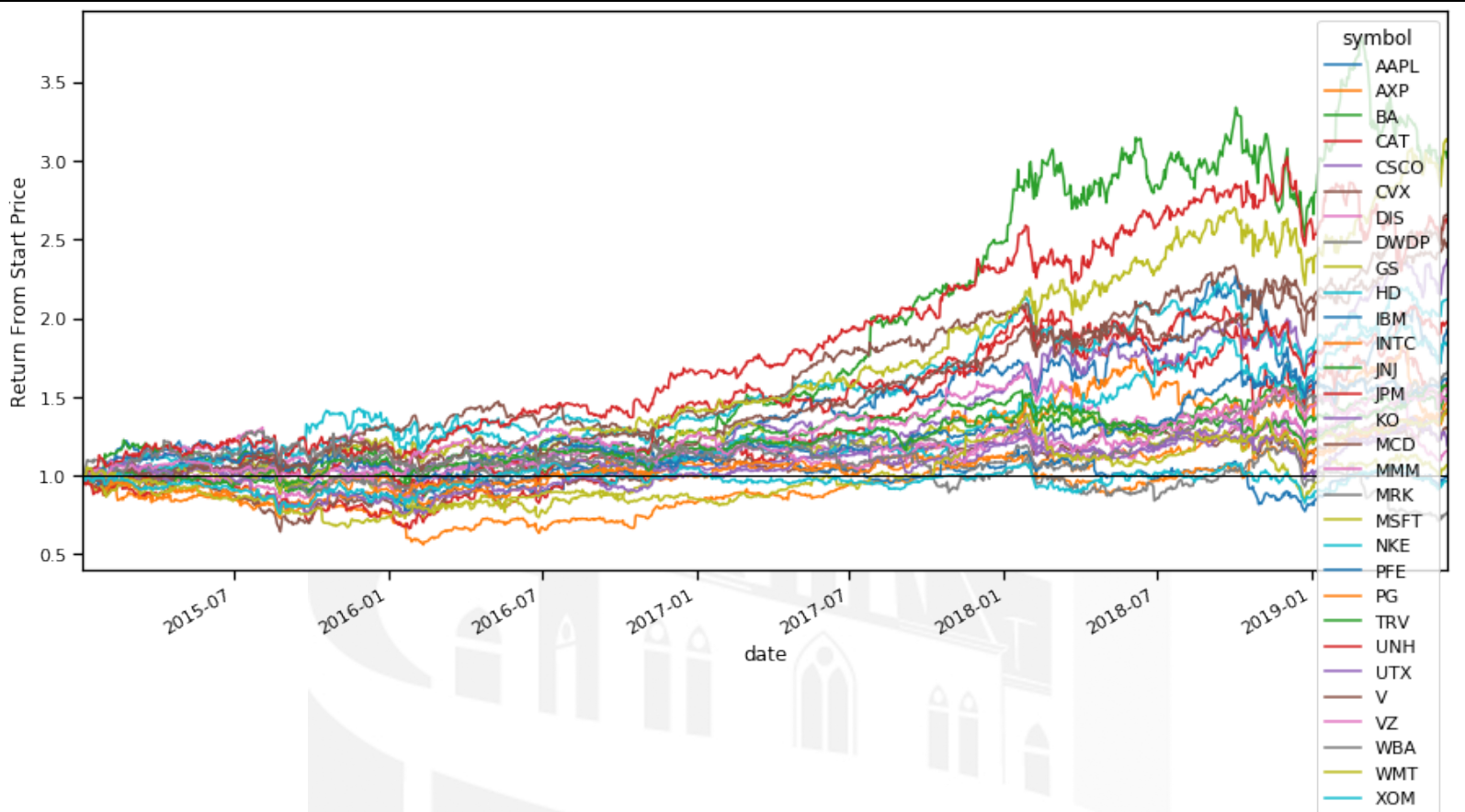




42 and 252 days rolling mean of AAPL

1870





Daily Returs of all 30 Dow Jones Stocks



## Algorithms and methods

### Linear Regression

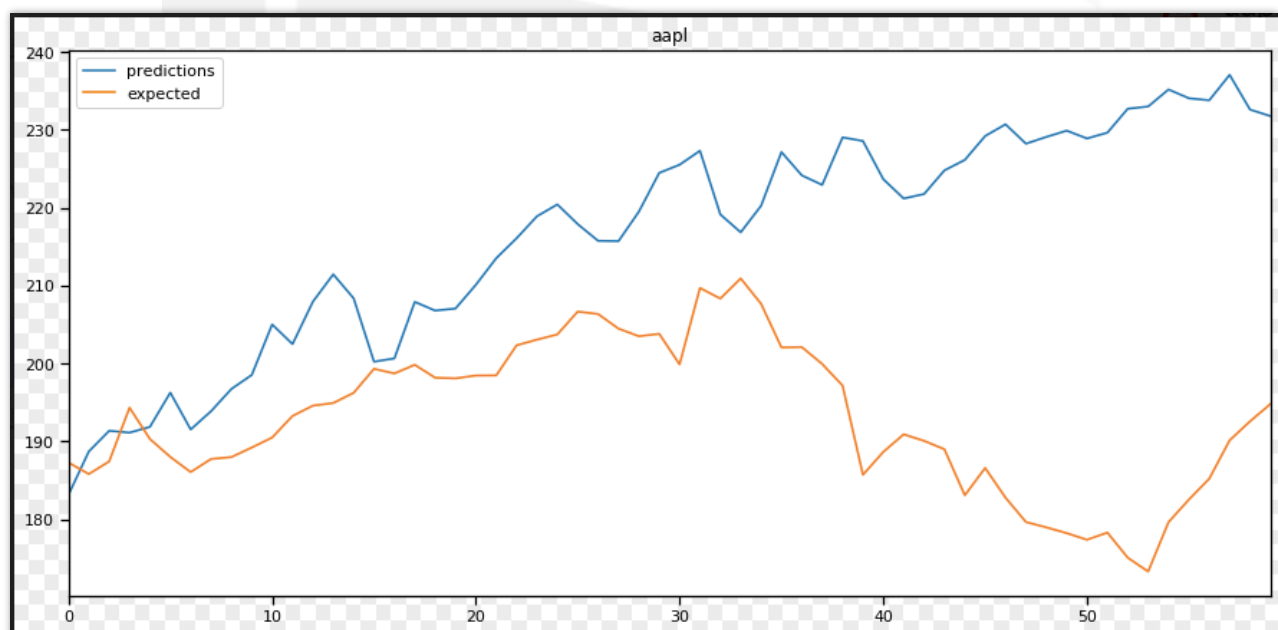
Linear regression is a linear approach to modeling the relationship between a scalar response and one or more explanatory variables.

Considering the data, we have and the prediction we wish to do, starting the process with linear regression seems to be the best bet.

We know that, linear regression would not perform well, keeping up with the trends of the stock market, but constructing a good model will help us establish a well-formed foundation to carry out further predictions and achieve perfect models.

We split the data into a 66% training and 34% test set, before continuing with the algorithm.

Here's how linear regression performed for Apple.





## Arima

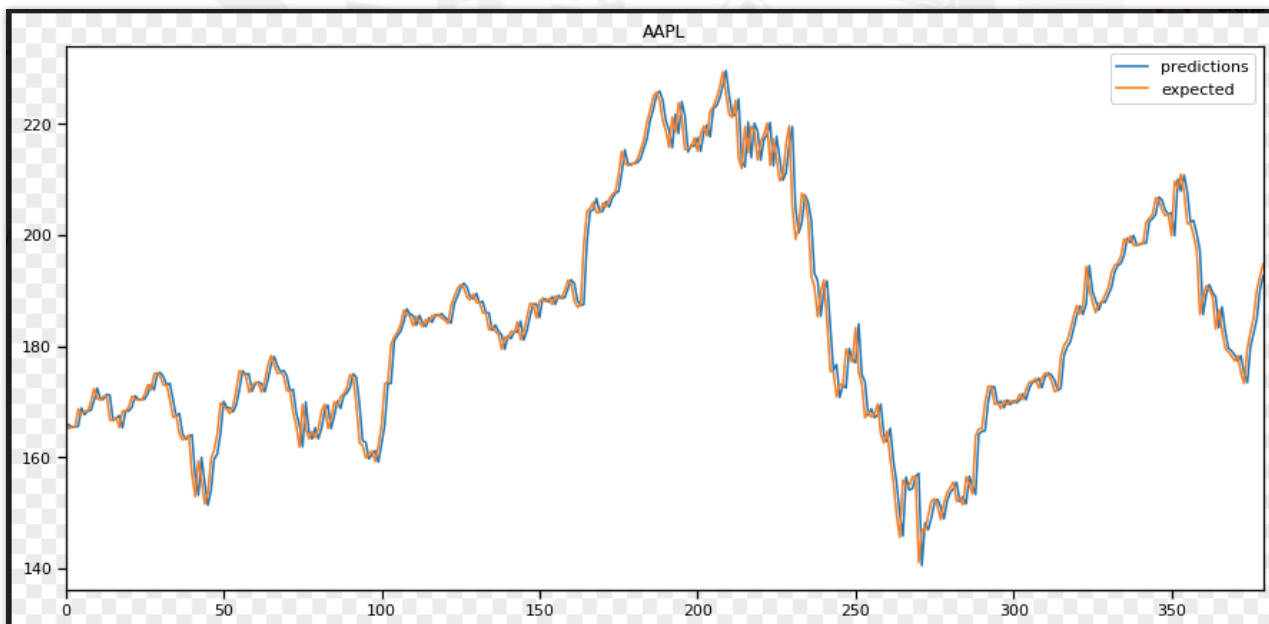
The obvious next bet would be to try Arima prediction. Using ARIMA model, you can forecast a time series using the series past values.

Since the Arima function doesn't demand a separate training and testing set. We gave it an initial training window of 66%

This means that the algorithm started working with a window of 66% of the data (around 737 values) from one column and predicted the 738<sup>th</sup> value. Next it considered the 738 values from the raw data set and predicted the next value. This continued till the dataset was eventually exhausted with all the values being tested for.

Arima turned out to be a very favorable model for our project.

Here's how Arima turned out for Apple,



```
In [94]: df_aapl = arima_prediction(ts_aapl)
```

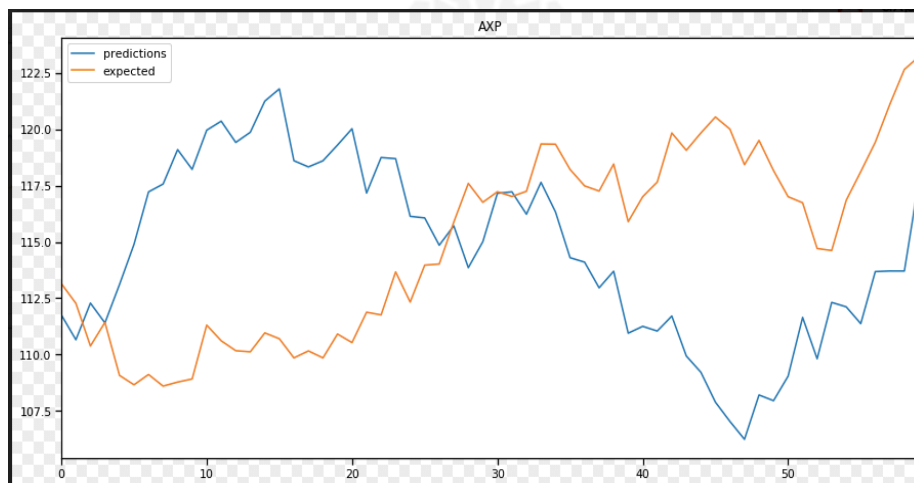
Test MSE: 10.900



## Comparison between Linear Regression and ARIMA

# American Express

### Linear Regression



### Arima



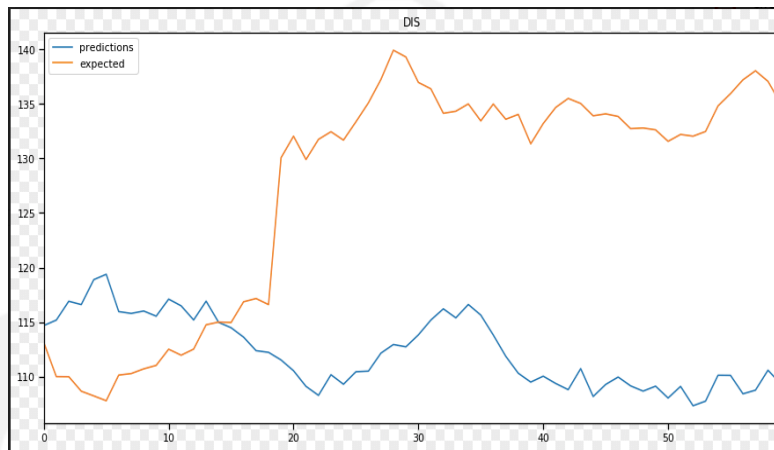
```
In [87]: df_axp = arima_prediction(ts_axp)
```

Test MSE: 1.866

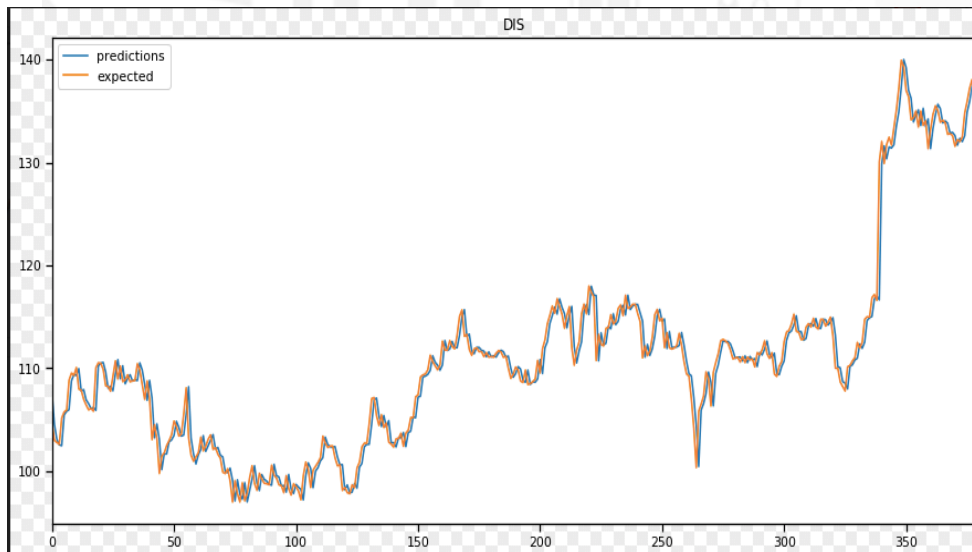


# Walt Disney Company

## Linear Regression



## Arima



```
In [87]: df_axp = arima_prediction(ts_axp)
```

Test MSE: 1.866



## Conclusion

- Considering the fact that we have to deal with time-series data here and the trend can go up or down with time, choosing Arima is a better option for our dataset instead of models like linear regression
- We also have to understand here that the accuracy we achieved is so high because the time frame is only worth a variably short amount of time
- By analyzing decades worth of data there might be introduced certain volatility, change in variance
- For such situations, Arima falls short. To extend our accuracy we need to research working of more models

## Future Plans

We researched a few models which can help us tackle issues related to volatility

- ARCH or Autoregressive Conditional Heteroskedasticity method provides a way to model a change in variance in a time series that is time dependent, such as increasing or decreasing volatility.
  - GARCH or Generalized Autoregressive Conditional Heteroskedasticity allows the method to support changes in the time dependent volatility, such as increasing and decreasing volatility in the same series
-



**STEVENS**  
INSTITUTE *of* TECHNOLOGY  
THE INNOVATION UNIVERSITY®

## References

- Dow Jones Industrial Average  
([https://en.wikipedia.org/wiki/Dow\\_Jones\\_Industrial\\_Average](https://en.wikipedia.org/wiki/Dow_Jones_Industrial_Average))
- Dow Jones Stock dataset (<https://www.kaggle.com/timoboz/stock-data-dow-jones>)
- Seasonal Adjustment  
([https://en.wikipedia.org/wiki/Seasonal\\_adjustment](https://en.wikipedia.org/wiki/Seasonal_adjustment))
- Linear Regression Example ([https://scikit-learn.org/stable/auto\\_examples/linear\\_model/plot\\_ols.html](https://scikit-learn.org/stable/auto_examples/linear_model/plot_ols.html))



1870