

# Machine Learning Engineer Nanodegree

## Capstone Proposal

Mincent Lee January 31st, 2018

## Proposal

### Domain Background

This project will develop a stock price predictor by machine learning. The proposal is historically simplified from the [Project Description ~ Investment and Trading](#)<sup>[1]</sup> and based on the [Course ~ Machine Learning for Trading](#)<sup>[2]</sup> for my first solid step to study machine learning for trading. Because the risk free rate of return from a bank account or a very short-term treasury bond is about 0 lately, folks have put so much money into the stock market<sup>[2:1]</sup>. The stock prediction can help us to understand market behaviour and trade profitable investments according to the wealthy information in the stock history and company data which is suitable for machine learning process<sup>[1:1]</sup>. There are lot related academic research support the stock prediction<sup>[2:2][3]</sup> while there are also opponent Efficient-Market Hypothesis<sup>[4]</sup>.

### Problem Statement

For reality and accuracy<sup>[4:1]</sup> concerns, the target problem of my first stock study is simplified to predict whether the *Adjusted* (for stock splits and dividends) *Closing price* rises or falls. The stock price predictor is inputted a certain range of daily trading data and outputs whether the *Adjusted Closing price* rises or falls (might ignore the rare flat cases at the first step) **next to the certain range, i.e., the predicted day is the *next day*, e.g., predicting the last Thursday according to the data of the last Monday to Wednesday. The next day is supposed to have the highest correlation and predictability according to the input features**, and suitable to be the basic first step. This is quantifiable, measurable, and replicable. The relevant potential solution are the Classifiers of the [scikit-learn](#)<sup>[5]</sup>, e.g., the [ensemble Gradient Boosting Classifier](#)<sup>[6]</sup>.

### Datasets and Inputs

The datasets used in this project is obtained from the [Yahoo! Finance](#)<sup>[7]</sup> by the python module [yahoo-finance](#)<sup>[8][9]</sup>. The target stock might be the [S&P 500 Index](#)<sup>[10]</sup> that might be the best representation of the U.S. stock market<sup>[10:1]</sup>. The inputs include *Opening price*, *Highest price*, *traded Volume*, *Adjusted Closing price*, and so on. Each price prediction is according to the trading data of a consistent **day range**, e.g., considering 2+1-day range in a trading week, the input ( $X_1, X_2, X_3 \dots$ ) and predicted ( $y_1, y_2, y_3 \dots$ ) days are:

	X (2-day range)	y (the next day of the range)
1	Mon. Tue.	Wed.
2	Tue. Wed.	Thu.
3	Wed. Thu.	Fri.

The sampled days for this project should include the current day for practicality and then trace back to find a balanced day range in which the distribution of the target classes (price *Rise/Fall*) is balanced for balanced evaluation metrics. The balanced day range could be searched from the same-price ranges in which the prices of the first and last day are the same to have balanced probability of *Rises* and *Falls*. The sampled day range might also larger than one year to cover annual and monthly characteristics. Some data also might be dropped out to force balanced if the balanced day range cannot be found.

## Solution Statement

The potential solution is training a Classifier by daily trading data within specific ranges of days to predict *Rising* or *Falling* of the *Adjusted Closing price* following the range. The daily trading data are obtained from the [Yahoo! Finance](#)<sup>[7:1]</sup> by the python module [yahoo-finance](#)<sup>[8:1]</sup>. The machine learning libraries and Classifiers might come from [scikit-learn](#)<sup>[5:1]</sup>, e.g., the [ensemble Gradient Boosting Classifier](#)<sup>[6:1]</sup>, and the parameter [random\\_state](#)<sup>[11]</sup> will be recorded. Therefore, the solution is quantifiable, measurable, and replicable.

## Benchmark Model

The predicted *Rising* and *Falling* results will be evaluated in fact with the exact benchmark of specific daily prices from the [Yahoo! Finance](#)<sup>[7:2]</sup> by the python module [yahoo-finance](#)<sup>[8:2]</sup>. Might also build a [Naïve Predictor](#)<sup>[12]</sup> which always predict *True/False* (*Rising/Falling*), if necessary.

## Evaluation Metrics

The solution model will be evaluated with the exact benchmark of specific daily prices from the [Yahoo! Finance](#)<sup>[7:3]</sup> by the  $F_\beta - score$ <sup>[13]</sup> with the `fbeta_score` function of `scikit-learn`<sup>[14]</sup>. The mathematical representations is:

$$F_\beta = (1 + \beta^2) \frac{precision \cdot recall}{\beta^2 precision + recall}$$

The  $\beta$  might be 1 for balanced precision and recall<sup>[13:1]</sup>.

## Project Design

- Analysis
  - Data Collection
    - Will try [S&P 500 Index](#)<sup>[10:2]</sup> first
    - Data include *Opening price, Highest price, traded Volume, Adjusted Closing price*, and so on
    - Will try the python modules
      - ["yahoo-finance 1.4.0," PyPI - the Python Package Index](#)<sup>[8:3]</sup>
      - ["fix-yahoo-finance 0.0.21," PyPI - the Python Package Index](#)<sup>[9:1]</sup>
  - Data Structure
    - Format data to the [DataFrame of pandas](#)<sup>[15]</sup>
  - Data Cleaning
    - Basic abnormal trading data handling<sup>[2:3]</sup>
    - [Imputation of missing values](#)<sup>[16]</sup>
  - Feature-set Exploration
    - Feature-set include *Opening price, Highest price, traded Volume, Adjusted Closing price*, and so on
    - Data Mining
      - [pandas.DataFrame.describe](#)<sup>[17]</sup>
      - [NumPy](#)<sup>[18]</sup>
    - Exploratory Visualization
      - [matplotlib](#)<sup>[19]</sup>
      - [seaborn ~ statistical data visualization](#)<sup>[20]</sup>
  - Benchmark
    - Models
      - Exact benchmark of specific daily prices from the [Yahoo! Finance](#)<sup>[7:4]</sup> by the python module [yahoo-finance](#)<sup>[8:4]</sup>

- Might build a [Naïve Predictor<sup>\[12:1\]</sup>](#) which always predict *True/False* (*Rising/Falling*), if necessary
- Evaluation Metrics
  - $F_\beta - score^{[13:2]}$ 
    - `fbeta_score` function of `scikit-learn`<sup>[14:1]</sup>
    - The  $\beta$  might be 1 for balanced precision and recall<sup>[13:3]</sup>
- Methodology
  - Data Pre-processing
    - Outlier detection
      - [Methods/algorithms of scikit-learn<sup>\[21\]\[22\]</sup>](#)
    - Normalizing Numerical Features
      - [Scaling/Standardization<sup>\[23\]</sup> scalars of scikit-learn<sup>\[24\]</sup>](#)
  - Encode Stock Price Changings for Classification
    - Encode the predicted *Adjusted Closing price* to *True* if it is *Rising* than the previous trading day (skip the non-trading days) and vice versa (*False* if it is *Falling*)
  - ~~Shuffle and Split Data~~
    - ~~Apply `sklearn.model_selection.train_test_split`<sup>[25]</sup> with recorded `random_state`<sup>[11:1]</sup> for replicability~~
    - The *Shuffle* might be avoided to avoid the [look ahead bias in time series<sup>\[26\]</sup>](#) as the reviewed comments.
    - The training date range might be the last year and the testing day range might be this year.
    - However, there is a *Shuffle* experiment that the data are constructed into many isolated day range packages in which the days inside are kept continuous for each prediction, e.g., considering 2+1-day range packages in a trading week, the input ( $X_1, X_2, X_3 \dots$ ) and predicted ( $y_1, y_2, y_3 \dots$ ) days are as the table below. Therefore, the prediction packages can keep continuous day range inside, and the outside index 1~3 can be shuffled.

	X (2-day range)	y (the next day of the range)
1	Mon. Tue.	Wed.
2	Tue. Wed.	Thu.
3	Wed. Thu.	Fri.

- Supervised Learning Models
  - [Classifiers in scikit-learn<sup>\[5:2\]</sup>](#), e.g., the [ensemble Gradient Boosting Classifier<sup>\[6:2\]</sup>](#)

- Training and Predicting Pipeline
  - Build normalizing, training, predicting, scoring functions
  - Make the normalizing [Scaler](#)<sup>[24:1]</sup> and [Classifier](#)<sup>[5:3]</sup> in a [Pipeline](#)<sup>[27]</sup>
- Initial Model Evaluation
  - Apply default and coarse-grained parameters
  - Record all the available parameter `random_state`<sup>[11:2]</sup> for replicability
  - The specified day range of each daily trading data might be a week initially
- Refinement
  - Fine-tune the hyper-parameters<sup>[28]</sup> (Exhaustive Grid Search<sup>[28:1]</sup> by Cross-validation<sup>[29]</sup> with [TimeSeriesSplit](#)<sup>[30]</sup>)
  - Tune the day range of each input might by [Feature importances](#)<sup>[31]</sup> or [Feature selection](#)<sup>[32]</sup>
- Results
  - Model Evaluation and Validation
    - Evaluate with the  $F_1$  — *score*<sup>[13:4]</sup>
    - Validate by the [Cross-validation](#)<sup>[29:1]</sup> with the [TimeSeriesSplit](#)<sup>[30:1]</sup> and long period of days
  - Justification
    - Compare with the exact benchmark of daily prices<sup>[8:5]</sup> in fact
    - Might compare with the [Naïve Predictor](#)<sup>[12:2]</sup>, if necessary

---

1. "MLND Capstone Project Description - Investment and Trading," [Udacity](#) ↔ ↔

2. Tucker Balch, "Machine Learning for Trading," [Georgia Tech](#) and [Udacity](#) ↔ ↔ ↔ ↔

3. "Stock market prediction," [Wikipedia](#) ↔

4. "Efficient-market hypothesis," [Wikipedia](#) ↔ ↔

5. "Choosing the right estimator," [scikit-learn.org](#) ↔ ↔ ↔ ↔

6. "sklearn.ensemble.GradientBoostingClassifier," [scikit-learn.org](#) ↔ ↔ ↔

7. Yahoo! Finance ↔ ↔ ↔ ↔ ↔

8. "yahoo-finance 1.4.0," [PyPI - the Python Package Index](#) ↔ ↔ ↔ ↔ ↔ ↔

9. "fix-yahoo-finance 0.0.21," [PyPI - the Python Package Index](#) ↔ ↔

10. "Standard & Poor's 500," *Wikipedia* ↔ ↔ ↔
11. "Utilities for Developers," *scikit-learn.org* ↔ ↔ ↔
12. udacity, "Project: Finding Donors for CharityML," *github.com* ↔ ↔ ↔
13. "F1 score," *Wikipedia* ↔ ↔ ↔ ↔ ↔
14. "sklearn.metrics.fbeta\_score," *scikit-learn.org* ↔ ↔
15. "pandas.DataFrame," *Python Data Analysis Library*, *pandas.pydata.org* ↔
16. "Imputation of missing values," *scikit-learn.org* ↔
17. "pandas.DataFrame.describe," *pandas.pydata.org* ↔
18. NumPy, *SciPy.org* ↔
19. matplotlib.org ↔
20. seaborn: statistical data visualization, *seaborn.pydata.org* ↔
21. "Outlier detection with several methods," *scikit-learn.org* ↔
22. "Comparing anomaly detection algorithms for outlier detection on toy datasets," *scikit-learn.org* ↔
23. "Standardization, or mean removal and variance scaling," *scikit-learn.org* ↔
24. "Compare the effect of different scalers on data with outliers," *scikit-learn.org* ↔ ↔
25. "sklearn.model\_selection.train\_test\_split," *scikit-learn.org* ↔
26. Rohit Walimbe, "Avoiding Look Ahead Bias in Time Series Modelling," *datasciencecentral.com* ↔
27. "Pipeline: chaining estimators," *scikit-learn.org* ↔
28. "Tuning the hyper-parameters of an estimator," *scikit-learn.org* ↔ ↔
29. "Cross-validation: evaluating estimator performance," *scikit-learn.org* ↔ ↔

30. "sklearn.model\_selection.TimeSeriesSplit," *scikit-learn.org* ↔ ↔

31. "Feature importances with forests of trees," *scikit-learn.org* ↔

32. "Feature selection," *scikit-learn.org* ↔