

Machine Learning Engineer Nanodegree

Capstone Proposal

Carl Martin

October 22nd, 2016

Domain Background

Predicting stock markets is hard. The efficient market hypothesis states that stock prices fully reflect all available information and prices are only influenced by new information. Since news happens randomly, stock prices follow a random walk pattern. This hypothesis has been extensively studied in the past. Researchers agree, that the behavior of stock prices is approximately close to the random walk process (Qian & Rasheed, 2006).

There are two major schools of thought of stock market analysis: fundamental and technical analysis. Fundamental analysis looks to determine a company's intrinsic value, while technical analysis tries to identify market trends. There are traders from both schools of thought that are very successful, proving that traders can exploit small market inefficiencies to increase portfolio value.

Since stock returns follow approximately a random walk model, it makes sense to treat stock market prediction as a classification problem. The goal is to classify whether a given stock will go up or down. Many studies have applied machine learning to this problem and achieved accuracy up to 65% for individual stocks using multiple classifiers (Qian & Rasheed, 2006) and up to 76% for indices exploiting temporal correlations (Shen, Jian & Zhang, 2012).

Problem Statement and Evaluation Metrics

Stock market prediction in this capstone project is treated as a classification problem, where the objective is to predict whether the cumulative returns of a stock over a period of one day will be greater than zero (up) or smaller (down). The final algorithm has to return (1) the predicted label and (2) a confidence level.

The final classification algorithm will be used in a long/short equity trading strategy. Using a ranking scheme based on both the predicted label and the algorithm's confidence the strategy goes long (buys) the top x% of equities of the ranking and goes short (sells) on the bottom x% of the ranking, while maintaining equal dollar volume between the long and short position (Granizo-Mackenzie). The implementation of the machine learning algorithm is out of the scope of this capstone project, yet relevant for the evaluation metric.

Given the algorithm's future use, the evaluation metric has to penalize the algorithm heavily for false positives and false negatives, for which the algorithm assigned a high confidence level. Given the

requirements, the classification algorithm's performance will be evaluated using *logarithmic loss*, which can be calculated using the following formula:

$$-\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ji} \cdot \log(p_{ji})$$

where N is the number of samples, M the number of possible labels, y_{ji} indicates whether label j is the correct classification for instance i , and p_{ji} is the probability yielded from the algorithm for label j . The objective will be to achieve a `log_loss` value as close to zero as possible.

Datasets and Inputs

Training and Testing Dates

Data will be split sequentially (not randomly) — the first 80% of the timeseries will be used for training the algorithm, which will be evaluated on the last 20% of the available data. For stocks that existed before 01-01-2003 all available from then until now will be used.

Stock Universe

Since the machine learning algorithm will be used for trading eventually, the stocks the algorithm is trained on have to be selected carefully. I'll draw from the TradableUS Methodology (Wassermann, 2016) to select a set of 500 stocks that are tradable, liquid and financially viable at the beginning of the training time to avoid survivor bias.

Pricing Data

Using the Yahoo finance API I'll be loading adjusted close, volume, high and low as needed. The API connection code is provided to reviewer. Data will be loaded automatically as needed.

Fundamental Data

I'll be using fundamental data and index prices from Intrinio, which is provided to the reviewer with in SQLite database.

- Market Capitalization
- Earnings before Interest and Taxes
- Earnings before Interest, Taxes, Depreciation and Amortization
- Sector Information
- S&P500, DAX, DASDAQ, DIJA, FTSE HENGSHENG, ASX and NIKKEI historical pricing

Should the algorithm not perform sufficiently I might consider additional free foreign exchange rate data from Quandl, which will be provided to the reviewer in a SQLite database as well.

Solution Statement

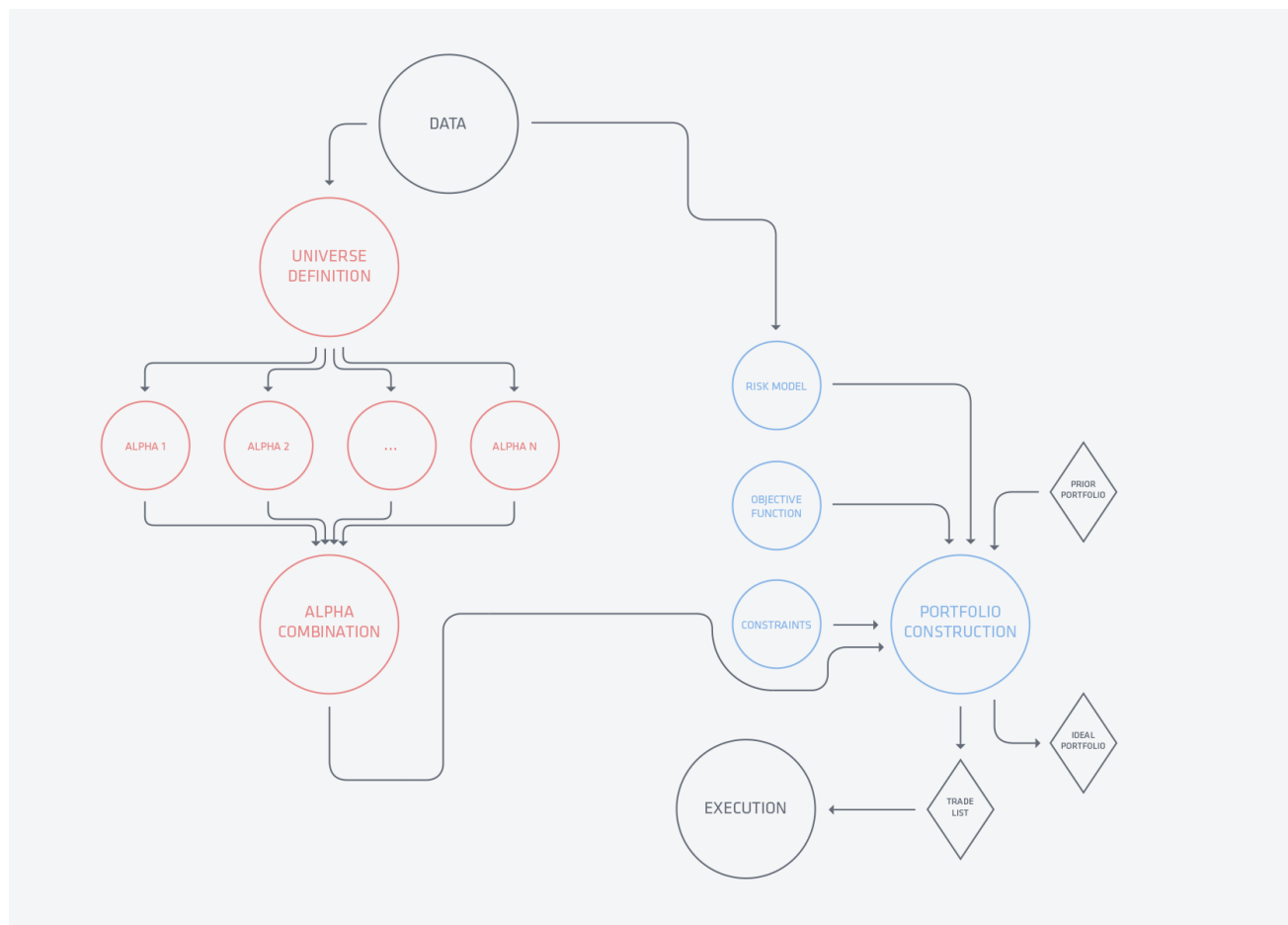
The final algorithms training interface will take a data range (`start_date`, `end_date`) and will train the algorithm on the preselected stock universe. The query interface will take `date` as an input and return a `DataFrame` with the stocks of the stock universe ranked by class (up/down) and algorithm confidence. The query date passed must be after the training `end_date`. Both, the training and querying interface, will be provided in an Jupyter Notebook as two respective Python 2.7 functions.

Benchmark Model

The benchmark model takes the scaled daily returns of the past four days for a given stock and uses the AdaBoost Classifier. The benchmark model achieved a logarithmic loss on average of 0.69. The specific implementation of the benchmark model can be found in the appendix (last section).

Project Design

My Capstone project will follow the [Professional Quant Equity Workflow](#) proposed by Jonathan Larkin, Chief Investment Officer at Quantopian. Having selected the data sources and stock universe above, I've already completed the first two steps of the workflow.



Having selected the trading 'universe', I'll move on to the feature engineering, or alpha discovery, stage of the workflow. I'll be calculating fundamental and technical indicators (alphas) with the data from the aforementioned datasets for the selected stock universe. Each alpha will be evaluated using the `alphalens` package to determine whether the alpha is predictive of future returns.

I intend to develop a system that is an ensemble of multiple machine learning classifiers. The goal of ensembles is that the performance of the entire system is better than the performance of each classifier individually. Another benefit of ensemble systems is that even if one underlying classifier fails, the system can continue to function, which makes it more robust (Qian & Rasheed, 2006).

I'll use the following algorithms as my base-learners (subject to change):

- `sklearn.svm.SVC` — Support Vector Classifier
- `sklearn.neighbors.KNeighborsClassifier` — K-Nearest Neighbors Classifier

- `sklearn.neural_network.MLPClassifier` — Neural Network Classifier
- `sklearn.ensemble.RandomForest` — Random Forest Classifier
- `sklearn.ensemble.AdaBoost` — AdaBoost Classifier

As the meta-learner, which combines the predictions of the individual classifiers, I'll be using the AdaBoost Classifier.

The classification system will be implemented using sklearn's `Pipeline` and `FeatureUnion` objects. The final Pipeline will be implemented the following way:

```
1  from sklearn.pipeline import Pipeline, FeatureUnion
2  from sklearn.base import BaseEstimator, TransformerMixin
3  from sklearn import preprocessing
4
5  # Transforms sklearn models to be used as features within pipelines
6  class ModelTransformer(BaseEstimator, TransformerMixin):
7
8      def __init__(self, model):
9          self.model = model
10
11      def fit(self, *args, **kwargs):
12          self.model.fit(*args, **kwargs)
13          return self
14
15      def transform(self, X, **transform_params):
16          return pd.DataFrame(self.model.predict(X))
17
18  pipeline = Pipeline([
19      ('features', FeatureUnion([
20          ('feature_1', Pipeline([
21              ('calculation', someCustomTransformer()),
22              ('scaling', preprocessing.RobustScaler())
23          ])),
24
25      ...
26
27      ])),
28      ('estimators', FeatureUnion([
29          ('estimator_1', ModelTransformer(someClassifier())),
```

```
30
31     ...
32
33     ])),
34     ('estimator', AdaBoost())
35 ])
```

The advantage of using sklearn Pipelines is that running GridSearch and Crossvalidation on the entire pipeline is feasible and the final model will be easier to understand.

References

Qian, Bo, and Khaled Rasheed. "Stock Market Prediction with Multiple Classifiers." *Applied Intelligence* 26.1 (2006): 25-33. Web.

Shunrong Shen, Haomiao Jiang, and Tongda Zhang., "Stock market forecasting using machine learning algorithms". *Stanford University*, 2012

Granizo-Mackenzie, Delaney. "Quantopian Lecture Series: Long-Short Equity Strategies." N.p., n.d. Web. 22 Oct. 2016.

Wassermann, Gil. "Quantopian - The Tradeable500US Is (almost) Here!" N.p., 6 July 2016. Web. 22 Oct. 2016.

Appendix

I provided the implementation of the Benchmark model below:

```
1  import pandas as pd
2  import numpy as np
3  from datetime import datetime
4  import requests
5  import json
6  from sklearn.base import BaseEstimator, TransformerMixin
7  from sklearn.pipeline import Pipeline, FeatureUnion
8  from sklearn.model_selection import GridSearchCV, RandomizedSearchCV
9  from sklearn import preprocessing
10 from sklearn.ensemble import AdaBoostClassifier
11 from sklearn.metrics import log_loss
12
13 #####
14 Necessary Classes
```

```
15 #####
16
17 class Database(object):
18
19     def __init__(self):
20
21         self.engine = create_engine('sqlite:///dataLake.db')
22
23     def save(self, tablename, data, behavior='fail', index=False):
24
25         '''behavior options:
26         fail: If table exists, do nothing.
27         replace: If table exists, drop it, recreate it, and insert data.
28         append: If table exists, insert data. Create if does not exist.'''
29
30         data.to_sql(tablename, self.engine, if_exists=behavior, index=index)
31
32     def search_ticker(self, query):
33         return pd.read_sql("SELECT security_name, ticker FROM 'Master Data'
WHERE security_name LIKE '%{}%' OR ticker LIKE '%{}%'".format(query, query),
self.engine)
34
35     def retrieve_prices(self, ticker):
36         ticker = ticker.upper()
37         prices = pd.read_sql("SELECT * from '{}'.format(ticker), self.engine)
38
39         prices = prices.set_index('date')
40         return prices
41
42 # Main Dataprovider class
43
44 class Intrinio(object):
45
46     def __init__(self):
47
48         self.user = 'b8f16afb02891a7fc1085118fd64ad47'
49         self.password = '5ad4db89b891c478803e731c8b3ff2d0'
50
51     def request(self, url, payload={}):
```

```
50
51     r = requests.get(url, params=payload, auth=(self.user, self.password))
52
53     if r.status_code == 200:
54
55         return json.loads(r.text)
56
57     else:
58         raise RuntimeError('An error has occurred while requesting data.
59 Status Code:{}'.format(r.status_code))
60
61     def extract_data(self, response):
62
63         return pd.DataFrame(response['data'])
64
65     def multiple_pages(self, url, parameters={}):
66         response = self.request(url, parameters)
67         data = self.extract_data(response)
68
69         if response['total_pages'] > 1:
70             i = 2
71
72             while i <= response['total_pages']:
73
74                 parameters['page_number'] = i
75                 data = data.append(self.request(url, parameters)['data'], ignore_index=True)
76                 i += 1
77
78             return data
79
80     def request_new_prices(self, ticker):
81
82         response = self.request('https://api.intrinio.com/prices', {'identifier': ticker, 'sort_order': 'asc'})
83         return self.extract_data(response).set_index('date')
```

```
84 # This class wraps repeatedly used functionality of assets into a Stock object
85 class Stock(object):
86
87     def __init__(self, ticker, data_provider=intrinio, database=database, validate=True):
88
89         self.data_provider = data_provider
90         self.database = database
91         self.ticker = ticker.upper()
92
93         if validate == True:
94             self.company_name = self.retrieve_company_name(ticker)
95
96         try:
97             # TODO: Check for current date
98             self.prices()
99         except:
100             self.download_prices()
101
102     def retrieve_company_name(self, ticker):
103
104         try:
105             company_name = pd.read_sql("SELECT security_name FROM 'Master Data' WHERE ticker = '{}'.format(ticker), self.database.engine).ix[0][0]
106             return company_name
107
108         except:
109             raise ValueError('The ticker {} you entered is not valid.'.format(ticker))
110
111     def compute_daily_ret(self, data):
112         data = data['adj_close']
113         daily_ret = data.copy()
114         daily_ret[1:] = (data[1:] / data[:-1].values) - 1
115         daily_ret.ix[0] = 0
116         return daily_ret
117
```



```
118     def download_prices(self):
119
120         prices = self.data_provider.request_new_prices(self.ticker)
121
122         if len(prices) < 1:
123             raise RuntimeError('An error occurred, while downloading the pri
ces.')
124
125         prices['daily_ret'] = self.compute_daily_ret(prices)
126         self.database.save(self.ticker, prices, 'replace', 'date')
127
128     def prices(self):
129         prices = self.database.retrieve_prices(self.ticker)
130         return prices
131
132 class Windowing(BaseEstimator, TransformerMixin):
133     """
134     http://stockcharts.com/school/doku.php?st=roc&id=chart_school:technical_
indicators:rate_of_change_roc_and_momentum
135     """
136     def __init__(self, t=5, feature='daily_ret'):
137
138         self.t = t
139         self.feature = feature
140
141     def fit(self, X, y=None, **fit_params):
142         return self
143
144     def transform(self, X):
145         X = X.copy()
146         columns = ['{}_{}'.format(self.feature, 0)]
147         X.loc[:, '{}_{}'.format(self.feature, 0)] = X[:, self.feature].values
148         for i in np.arange(1, self.t+1):
149             X.loc[:, '{}_{}'.format(self.feature, i)] = X[:-i][self.featur
e].values
150             X.loc[:, '{}_{}'.format(self.feature, i)] = 0
151             columns.append('{}_{}'.format(self.feature, i))
152         return X[columns]
```

```

153
154 def train_test_split(X, y, test_size=0.2, target='daily_ret'):
155     index = int(len(X) * (1 - test_size))
156     return X[:index], X[index:], y[:index][target].values, y[index:][target].values

```

```

1 #####
2 Algorithm
3 #####
4
5 intrinio = Intrinio()
6 # Set validate to false, because I didn't provide the Master Data class
7 stock = Stock('AAPL', validate=False)
8
9 X = stock.prices()[['adj_low', 'adj_high', 'adj_close', 'adj_volume', 'daily_ret']]
10
11 # Time delay between X and y - Experimenting with time delay lead to improved scores. Likely because it reduces noise. TODO: Investigate
12 time_delay = 1
13
14 # Calculate the price difference between current price and time_delay days prior
15 target_transformation = Momentum(time_delay) #
16 y = target_transformation.transform(X)
17 y.columns = ['target']
18 y['target'] = y['target'].apply(lambda x: 1 if x > 0 else -1)
19
20 # Training, testing split
21 X_train, X_test, y_train, y_test = train_test_split(X[:-time_delay], y[time_delay:], target='target')
22
23 # Pipeline for feature engineering
24 pipeline = Pipeline([
25     ('features', FeatureUnion([
26         ('daily_ret', Pipeline([
27             ('windowing', Windowing(feature='daily_ret')),
28             ('scale', preprocessing.RobustScaler())

```

```
29         ]))
30     ], n_jobs=1)),
31     ('estimator', AdaBoostClassifier(random_state=0))
32 ])
33
34 params = {
35     'estimator__n_estimators': np.arange(100,200, 10),
36     'features__daily_ret__windowing__t': np.arange(4,10)
37 }
38
39 # Basic hyperparameter optimization of algorithm and features
40 grid_search = GridSearchCV(pipeline, param_grid=params, n_jobs=12)
41 grid_search.fit(X_train, y_train)
42 classifier = grid_search.best_estimator_
43
44 # Evaluation of classifier with log_loss metric
45 y_proba = classifier.predict_proba(X_test)
46 log_loss(y_test, y_proba)
```