# Machine Learning Engineer Nanodegree

## Capstone Proposal

Arya Roy
23rd March, 2018

## 1. Domain Background

The emergence of structured market data and improvement in the computational power in the last decade has led to machine learning getting successfully applied in quantitative trading. Finance practitioners often derisively refer to this methodology as data mining as financial data are not only quite limited (unless we use tick data), they are also not very stationary in the statistical sense as discussed in [1]. That is, the probability distribution of returns does not stay constant forever. If we use Machine Learning algorithms on these data, it is very easy to come up with trading rules that worked extremely well in certain past periods, but fail terribly going forward. However, this problem also plagues handcrafted models, but remedial action is possible. No such luck with machine-learned rules.

Despite the reluctance to use AI techniques in trading, this is a rapidly advancing field, where researchers are working towards developing techniques to solve the overfitting bias as mentioned in the previous paragraph. One such method is the use of a reinforcement learning framework to generate profitable trading signals. One key advantage of this method over supervised learning approaches is does not require learning from examples provided by a knowledgeable external supervisor[2].

## 2. Problem Statement

The goal of the project is to develop an adaptive learning agent that trains itself based on historical price data on how to maximize the risk-adjusted return (including net transaction cost) i.e. Shape ratio. The environment used here is the historical closing prices of 50 large-cap stocks listed in the National Stock Exchange (Nifty 50 - ^NSEI) ranging from 4th Jan, 2010 to 21t Feb, 2018. The states of the agent are defined based on whether the returns are above, within or below

a dynamic price momentum band around the moving average of return (window of 12 trading days). On each trading day, the agent must generate a valid action – buy stock, do nothing or sell stocks. Based on its actions, the agent should receive a reward or penalty equal to the risk-adjusted stock returns (profitable or loss making) and the transaction cost as a function of the closing price for every new buying or selling of a stock. Based on the rewards or penalties it gets or each stock, the agent should learn an optimal trading strategy to maximize the risk-adjusted returns of a portfolio of stocks. In order to learn from actual trading data can be costly as it involves the possibility of making losses while the agent is learning. Hence I have used a novel method called Dyna Q as given in [2] in order to solve the problem of insufficient market data. Using Dyna Q algorithm, the real experiences, passing back and forth between the environment and the policy, affect the policy and the value functions in much the same way as simulated experiences generated by the model of the environment.

## 3. Datasets and Inputs

The dataset consists of market data (open, high, low, close, volume and open interest) for all the stocks listed In the Nifty 50 for every trading day ranging from 4th January, 2010 till 21st February, 2018 (a total of 2012 trading days) as shown in Figure 1.

```
>>> import os
>>> from fileinput import filename
>>> def explore() :
...     base='E:/Bhavcopy/data/NSE-EOD'
...     filenames=os.listdir(base)
...     print ("Number of files:",len(filenames))
...
>>> explore()
Number of files: 2012
```

*Figure 1: Code snippet showing number of files*

The start date is chosen such that the data excludes the 2008-09 recession. This data is sourced directly from the National Stock Exchange database using an open-source application [3]. There are files corresponding to each trading day with a list of tickers. The sample of 1 file is shown in Table 1 below:

2

| <ticker> | <date> | <open> | <high> | <low> | <close> | <volume> | <o/i> |
|---|---|---|---|---|---|---|---|
| 20MICRONS | 2/21/2018 | 51.75 | 52 | 50.05 | 50.35 | 46224 | 24176 |
| 21STCENMGM | 2/21/2018 | 45.15 | 45.15 | 43.45 | 43.45 | 10217 | 0 |
| 3IINFOTECH | 2/21/2018 | 5.75 | 5.85 | 5.6 | 5.7 | 2910036 | 1550093 |
| 3MINDIA | 2/21/2018 | 20799 | 20906 | 20430 | 20543.45 | 1446 | 682 |
| 5PAISA | 2/21/2018 | 333.15 | 350.6 | 333.15 | 350.6 | 8535 | 7001 |

*Table 1: Trading data in the imported file fpr 2/21/2018*

The <ticker> column corresponds to the tickers that are traded in the market on the particular day. The <date> column gives the date on which the trade is executed. <open> column gives the price at which the ticker commenced trading on the exchange on the given day. <high> column gives the highest price of the ticker on the given day. The <low> column gives the lowest price of the ticker on the given day. <close> column gives the last price at which the stock traded on the particular day. The <volume> column gives the total number of shares traded on that day. The <o/i> column gives the total number of outstanding contracts at the end of the day.

For training and testing our learning algorithm, the data is preprocessed to create a single table with the closing price data with the ticker names (all tickers as of 16[th] January, 2018) as the columns and the trading dates as the rows. During this time the index has shown a net upward trend, starting from 5232.2 and ending at 1039.75. Since there is a significant change in the index in the period under consideration, the daily returns will be calculated and used for the purpose of learning instead of the closing price data to avoid bias.

Additionally, the following additional columns are calculated based on the existing data in order to compute the agent's state:

*risk_adjusted* : (Daily return – Bond rate)

*risk_adjusted_moving* : Moving average of risk_adjusted with a moving window of 12 trading days

*risk_adjusted_stdev* : Standard deviation of the risk_adjusted data in the moving window of 12 trading days

*risk_adjusted_high* : It is calculated for each stock in the portfolio as given below :

*risk_adjusted_moving + 1.5*risk_adjusted_stdev*

*risk_adjusted_low* : It is calculated for each stock in the portfolio as given below :

*risk_adjusted_moving – 1.5*risk_adjusted_stdev*


Now based on the above calculated data, the state of the individual stock is decided:

-1 : stock return is below *risk_adjusted_low*

0 : stock return is between *risk_adjusted_low* and *risk_adjusted_high*

1 : stock return is above *risk_adjusted_high*

```
>>> df1 = pd.read_csv('E:/Bhavcopy/data/final.csv')
>>> df1['Date']= pd.to_datetime(df1['Date'])
>>> df1['Date'] = df1['Date'].dt.strftime('%Y/%m/%d')
>>> df1 = df1.set_index('Date')
>>> df1 = df1.filter(items=['HDFCBANK','NSENIFTY'])
>>> l=len(df1.columns)
>>> for i in range(l):
...     df1['_'.join(['return',str(i)])]=df1.iloc[:,i].rolling(window=2).apply(lambda x: x[1]/x[0]-1)
...     df1['_'.join(['return_moving',str(i)])]=df1.iloc[:,-1].rolling(window=12).apply(lambda x: x.mean())
...     df1['_'.join(['return_stdev',str(i)])]=df1.iloc[:,-2].rolling(window=12).apply(lambda x: x.std())
...     df1['_'.join(['return_high',str(i)])]=df1.iloc[:,-2]+0.5*df1.iloc[:,-1]
...     df1['_'.join(['return_low',str(i)])]=df1.iloc[:,-3]-0.5*df1.iloc[:,-2]
...
>>> df1.iloc[:,[2,4,7,9]].loc['2017/10/26':'2018/02/21'].plot()
```

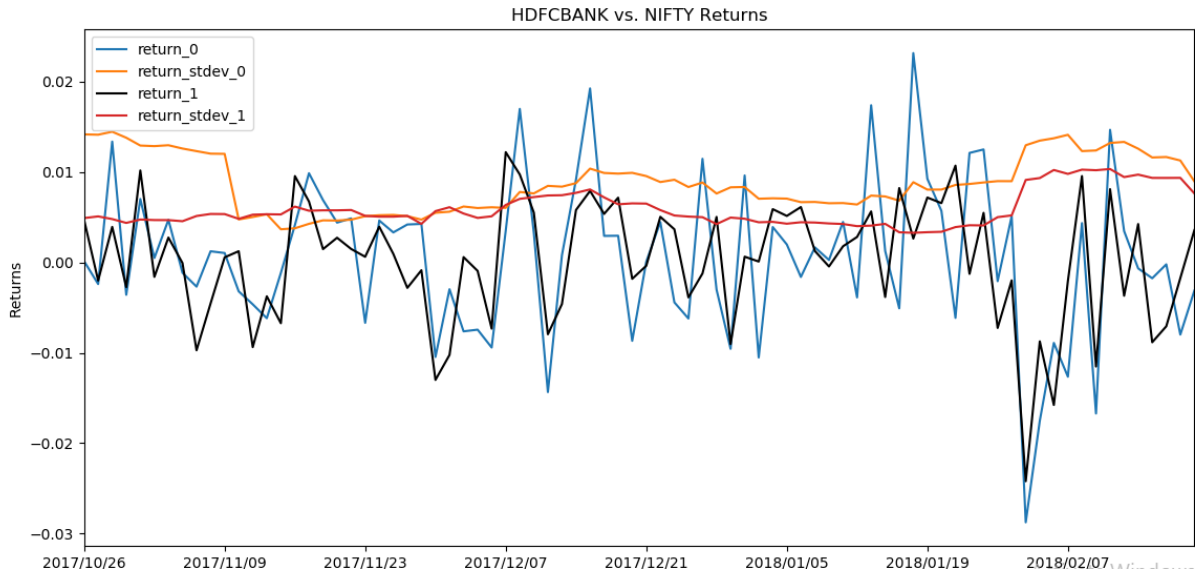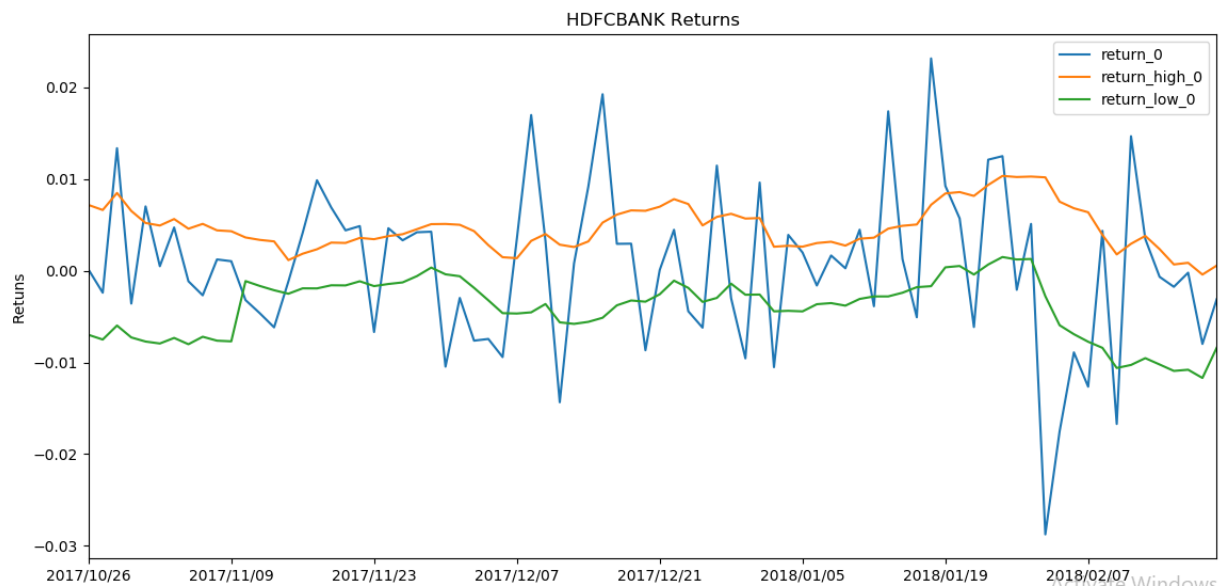*Figure 2: Code snippet showing calculation of additional fields*

4

*Figure 3: Plot of returns and volatility of 'HDFCBANK' and 'NIFTY'*

The above figure shows the relation between the returns on a stock 'HDFCBANK' (*return_0*) and that of the exchange 'NIFTY' [1](*return_1*). The standard deviation of the 'HDFCBANK' returns is in general higher than that of 'NIFTY'. Also, it is observable from the above data that there is a strong correlation between the returns of the stock compared to the exchange.

The above figure shows the returns of the 'HDFBANK' stock as well the bandwidth of 1 standard deviation on either side of the moving average of the returns on the stock. For the 'HDFCBANK' stock, it is in state 1 when the returns are above the orange line (*return_high* : 1 standard deviation above the moving average), state 0 when the returns are between the orange line and the green line(*return_low* : 1 standard deviation below the moving average of stock returns) and in state -1 when returns are below the green line.

Furthermore, I use the 10 year sovereign bond rates as the risk-free rate. The annual rates are converted to daily rates in order to calculate the risk-adjusted daily rates. The returns on the sovereign bonds are considered the benchmark for our model. In this project it will also be demonstrated that this agent works better than a model taking random action on the same set of stocks,

---

[1] Source : https://en.wikipedia.org/wiki/NIFTY_50

## 4. Solution Statement

The goal in using reinforcement learning to adjust the parameters of a system is to maximize the expected payoff or reward that is generated due to the actions of the system according to [2]. This is accomplished through trial and error exploration of the environment. The system receives a reinforcement signal from its environment (a *reward)* that provides information on whether its actions are good or bad. The performance function at time *T* can be expressed as a function of the cumulative return (product of daily returns).

$$U_T = U(R_1, R_2, \ldots, R_T)$$

Here the Markov Decision Process (MDP) is used to implement Reinforcement Learning according to [4]. The MDP model consists of a discrete set of states S and a discrete set of actions A. The set of actions are as follows:

$$a_t \in (\textit{None, Buy, Sell})$$

*None* indicates that the agent takes no action in the current time step. *Buy* indicates that the agent makes a *Buy* order of 1 share. *Sell* indicates that the agent makes a *Sell* order of 1 share.

At each discrete time step *t*, the agent computes its current state and takes an action $a_t$. Based on this, the environment provides the reward $r_t = r(s_t, a_t)$ that is the payoff due to the action taken. The environment also provides the next state $s_{t+1} = f(s_t, a_t)$.

The reward and next state functions depend only on the current state and action. Thus it does not have a memory.

The task of the agent is to learn a policy $\Omega$ that maps each state to an action ($\Omega: S \rightarrow A$), selecting its next action $a_t$ based solely on the current observed state $s_t$ that is $\Omega(s_t) = a_t$. The optimal policy, or control strategy, is the one that produces the greatest possible cumulative reward over time. So, we can state that:

$$V^{\Omega}(s_t) = \mu_0 r_t + \mu_1 r_{t+1} + \mu_2 r_{t+2} + \ldots\ldots = \sum_{i=0}^{\infty} \mu_i r_{t+i}$$

$V^{\Omega}(s_t)$ is the discounted cumulative reward with the discount factor $\mu_i$ E (0,1). V represents the cumulative value achieved by following the policy $\Omega$ from an initial state $s_t$. $\mu_i$ represents the relative value of delayed versus immediate reward.

If we set $\mu_i = 0$, then only immediate reward is considered. As $\mu_i \rightarrow 1$, the future rewards are given greater emphasis compared to the immediate reward. In our case the value of $\mu_i$ is given below:

$$\mu_i = \sin(\exp(-i))$$

The optimal policy $\Omega^*$ that will maximize $V^{\Omega}(s)$ for all states s can be written as:

$$\Omega^* = \text{argmax}_{\Omega} V^{\Omega}(s) \quad, \forall s$$

There is no direct method to learn $\Omega^*: S \rightarrow A$ since there is no available training data of the form (s,a). The optimal policy needs to be learnt from the available immediate reward data $r(s_i, a_i)$ where $I = 1,2,3\ldots$

We need to maximize $V^*(s_t)$ such that for all states s, the agent should prefer the next state $s_1$ over $s_2$ wherever $V^*(s_1) > V^*(s_2)$. We also need to learn $V^*(s)$ independently since the agent can only choose the action and not the state and it cannot perfectly predict the cumulative reward and immediate next state for every possible state-action transition.

To learn $V^*(s)$, we define a function $Q(s,a)$ such that its value is the maximum discounted cumulative reward that can be achieved starting from state s and applying action a as the first action. So we can write:

$$Q(s,a) = r(s,a) + \mu V^*(f(s,a))$$

$F(s,a)$ is the state resulting from applying action a to state s that is chosen by the optimal policy. Hence what we are trying to achieve is:

$$\Omega^*(s) = \arg \max_a Q(s,a)$$

The above equation implies that the optimal policy can be obtained with the current state a and current action a only if a maximizes $Q(s,a)$. Thus even if the agent has no knowledge of the functions r and f, it can find the optimal policy by choosing the maximum Q value.

In order to incorporate transaction cost, a fixed percentage of 1.5% is deducted from daily returns. This is because it is found that the sum of trading costs, including tax, brokerage etc. , is approximately 1.5% of daily returns.

The Q Learning algorithm requires many interactions with the real world. But this can be expensive especially in Algorithmic Trading as noted in [5]. Thus the Dyna Q method is used to hallucinate the real world experiences by randomly choosing a previously observed state s and

action a taken in state s to update the Q table. This process is iterated 100 times in our model for every real time step. This ensures that the Q table converges to that corresponding to the optimum policy with limited real world data.

Initialize Q(s, a), for all s, a

Do forever:

(a) s $\leftarrow$ current (nonterminal) state (direct update in every timestep)

(b) a $\leftarrow$ ε-greedy(s, Q)

(c) Execute action a; observe resultant state, s', and reward, r

(d) Q(s, a) $\leftarrow$ Q(s, a) + μ(r + arg max $_{s'}$ Q(s', a') − Q(s, a))

(e) State Action Pair (s, a) $\leftarrow$ (s', r)

(f) Repeat N time, (Dyna Q update):

s $\leftarrow$ random previously observed state

a $\leftarrow$ random action previously taken in s

State Action Pair (s', r) $\leftarrow$ (s, a)

Q(s, a) $\leftarrow$ Q(s, a) + μ(r + arg max $_{s'}$ Q(s', a') − Q(s, a))

Figure 4: Dyna Q algorithm

must converge to an optimal policy. Since we are using a deterministic MDP, this will be possible when the agent chooses actions in such a way that it visits every possible state-action pair infinitely often. We will also be using a stochastic component ω with a probability of 1% in order to randomly choose an action a from the available subset of actions A.

## 5. Benchmark Model

The benchmark will be a random agent that randomly chooses actions from the set A at every time step t. The goal of the learner agent is to outperform the random agent in terms of the risk-adjusted returns (Sharpe Ratio).

The available time series data is divided into 2 sections – training dataset and testing dataset. The training index consists of 1723 time series points from 2010/1/4 to 2016/12/22. The testing index consists of 283 time series points from 2016/12/23 to 2018/02/21. First I will analyze if the learning agent was able to improve its performance on the training dataset after different trials. After that, I will use the policy learned to simulate the learning agent behavior in different behavior in the testing dataset. The performance will be compared in the testing dataset for both the learning agent and the random agent.

We will also compare the performance of the learning agent against that of buying and holding 5-year sovereign bonds during the testing period. The 5-year sovereign bond is considered to have the risk-free rate of interest. The risk-adjusted return of the learning agent must be greater than the risk –free rate of return in order for this Q Learning strategy to be feasible.

## 6. Evaluation Metrics

The 3 metrics used to evaluate the performance of the proposed model are: Sharpe Ratio, Maximum Drawdown and Maximum Drawdown Duration.

According to [8], Sharpe Ratio is defined as the excess return over the risk-free rate per unit of underlying risk. Here the proxy for portfolio risk is volatility or standard deviation of returns. Here the daily volatility of returns is calculated as opposed to the weekly or monthly one to give more conservative Sharpe ratio.

$$S = \frac{Portfolio\ Return - (Risk-free\ Rate)}{Standard\ Deviation\ of\ Portfolio\ Return}$$

In the above equation, S is the Sharpe Ratio. The higher the Sharpe Ratio, the higher the expected risk-adjusted return from the portfolio. The next metric used is Maximum Drawdown. According to [9], Maximum Drawdown is the indicator of the downside risk associated with a particular stock or portfolio. It is defined as the maximum loss from peak to trough of a portfolio before the portfolio returns are positive again. The third metric used is the Maximum Drawdown Duration. This is simply the duration in days when the portfolio value falls for Maximum Drawdown. Hence for the Maximum Drawdown period,

$$\text{Maximum Drawdown} = \frac{(Peak\ Vale - Trough\ Value)}{Peak\ Value}$$

We will find out if the proposed model provides healthy returns (daily Sharpe ratio above 0.2) in different scenarios. We will test the robustness of the model by testing the optimal policy learned over 5 different testing datasets with 50 trading days included in each.

We will observe the daily Sharpe ratio to decide whether the proposed model gives a significant risk-adjusted return in the 5 different testing periods. In the period with the least risk-adjusted return, we will check the maximum drawdown and the maximum drawdown duration to ascertain how capital-intensive the trading strategy is.

## 7. Project Design

In this project we have used the Reinforcement Learning Algorithm (Q Learning Algorithm) to build a trading agent that learns how to trade based on market conditions and its own states. This is a model-free approach which learns the optimal trading policy based solely on the market parameters during the testing period. The agent learns the optimal policy to maximize profits based on the market statistics part of the training dataset. It populates the Q table that lists the expected returns for a particular stock in the portfolio based on the state the agent is in. In the testing period, the learning agent references the learned Q table in order to find the next state it

should be in to maximize returns. The learning agent shows a significant improvement in performance from the random agent that takes random trading decisions. The random agent is seen to give better returns than the risk-free rate of return of 5 year sovereign bonds. But this is seen when the average return of the random agent is taken after 20 trials. The learning agent gives significantly better returns than the risk-free rate since it outperforms the random agent as shown in the previous section.

One issue faced in training the learning agent was that there is a considerable chance that the agent will not converge to an optimal policy. In order to counter this problem, the agent is configured to pick a random action in 1 out of 10 training instances. This ensures that the policy is not biased towards 1 state. Also the decay parameter used falls as we go further back in time. This ensures that the expected returns approximation gives more weight to recent observations compared to older observations.

The reinforcement learning agent goes through 100 iterations in order to converge to an optimum policy. The policy learnt is used by the agent to take trading decisions during the testing period. One of the key challenges was to wrangle data to create a training dataset with the parameters required in order to represent distinct states for the learning agent. The next challenge was to represent the states in a manner such that the optimal policy can be learnt in 100 iterations or less. If the number of states is greater than 3 then learning the optimal policy will take many more iterations. As the number of states becomes significantly higher, the chances of the Q learning algorithm to converge to an optimal policy within a finite number of trials reduces. Establishing the environment was also difficult. At every time step, the agent has to evaluate which state it is in based on the environment parameters. During the training period, the agent must calculate the immediate reward for taking an action to move to the next state. Based on this, the agent is able to estimate iteratively the Q value or expected value of taking a trading action to buy, sell or hold.

The next challenge was implementing the reinforcement learning algorithm. The dataset consists of data corresponding to individual tickers. However the agent chooses a portfolio of stocks to

optimize the portfolio returns. This is more complicated than using the same algorithm on a single stock.

Dyna Q implementation was necessary since the number of data points was insufficient for the convergence to the optimal policy. Although it increased the processing time by a few minutes, this was only a small price to pay for the right approximation of the Q table corresponding to the optimal policy.

## 8. References

[1] Chan, Ernest P. Machine Trading: Deploying Computer Algorithms to Conquer the Markets. John Wiley &Sons, Inc., 2017.

[2] Sutton, Richard S., and Andrew G. Barto. Reinforcement Learning: an Introduction. The MIT Press, 2012.

[3] Kapadia, H. (2016). GETBHAVCOPY. Mumbai: Hemen Kapadia.

[4] Xin Du, Jinjian Zhai, and Koupin Lv. Algorithm trading using q-learning and recurrent reinforcement learning

[5] Classroom.udacity.com. (2018). Udacity. [online] Available at:
https://classroom.udacity.com/courses/ud501

[6] John Moody and Matthew Saffell. 1999. Reinforcement learning for trading. In Proceedings of the 1998 conference on Advances in neural information processing systems II, David A. Cohn (Ed.). MIT Press, Cambridge, MA, USA, 917-923.

[7] Lee J.W., O J. (2002) A Multi-agent Q-learning Framework for Optimizing Stock Trading Systems. In: Hameurlain A., Cicchetti R., Traunmüller R. (eds) Database and Expert Systems Applications. DEXA 2002. Lecture Notes in Computer Science, vol 2453. Springer, Berlin, Heidelberg

[8] Adam Hayes, C. (2018). Sharpe Ratio. [online] Investopedia. Available at:
https://www.investopedia.com/terms/s/sharperatio.asp

[9] Elvis Picardo, C. (2018). Maximum Drawdown (MDD). [online] Investopedia. Available at:

https://www.investopedia.com/terms/m/maximum-drawdown-mdd.asp