

Capstone Project

Machine Learning Engineer Nanodegree

Krishnan Ramaswamy

Reinforcement algorithms for trading systems

Table of Contents

Krishnan Ramaswamy.....	1
1. Problem Statement.....	2
2. Proposed Solution	3
3. Data Sets and Inputs	3
4. Benchmark	4
5. Design.....	4
5.1. Random Action Agent	5
5.2. Random Q-Learning Agent	5
5.3. Flexible Q-learning Agent.....	6
6. Evaluation metrics	10
7. Input and running the algorithm	11
7.1. Indicators.....	11
7.2. Strategies.....	12
7.3. Instances.....	13
8. Automation tools, logs and embedded runtime statistics.....	14
9. Implementation & Code snippets.....	15
9.1. Framework	15
9.1.1. Agent Factory	15

9.2.	Playback Simulator	16
9.2.1.	Indicators.....	17
9.3.	Agents	18
9.3.1.	Random Action Agent	18
9.3.2.	Simple/Random Q-Learning Agent.....	19
9.3.3.	Flexible Q-learning Agent	19
10.	Run & Observation.....	21
10.1.	Random Action Agent	21
10.2.	Simple/Random Reward Q-Learning Agent	24
10.3.	Flexible Q-learning Agent	24
11.	Improve the Q-Learning Driving Agent	27
12.	Conclusion	32
13.	References	33
5.	Market-Neutral Trading: Combining Technical and Fundamental Analysis Into 7 Long-Short Trading Systems	34

1. Problem Statement

The definition of automated or algorithmic trading can be classified as following groups.

1. Market making. Mostly to control and drive the direction of market movement, which derived from the order books' information, to provide liquidity and profit from bid/ask spread from market makers' perspectives. Stochastic control in high frequency trading plays important role here.
2. Algorithmic trading. The motivation is to identify patterns and intelligent decision making like professional/experienced trader using combination of technical, fundamental and market sentiments.

Traders/investors uses many technical indicators and trading strategies which have been developed by researchers and hedge fund managers/companies. Also statistical and other computer aided machine learning approaches are prevalent in hedge fund industry. More recently, there are new cloud based services are being available – like <https://stocksneural.net/> which helps in predicting the direction of stock

prices using machine learning/neural network techniques.

The objective of this project is to utilize machine learning/reinforcement techniques to identify market directors and corresponding decisions to BUY and SELL and provide better performance/risk than S&P

2. Proposed Solution

The objective of reinforcement learning based approach is not the minimization of the sum-of-squares error which is one of the objective of conventional supervised learning but the acquisition of an optimal policy under which the learning agent achieves the maximal average reward from the environment.

In this project, I am proposing an improvised reinforcement learning framework with many(flexible) indicators involved in prediction criteria with trading policies more effectively. The value approximator is trained using a regularizing technique for the prevention of divergence of the parameters with S&P. I am hoping to demonstrate a stock trading system implemented using the proposed framework will outperform the market average or individual stock being compared.

3. Data Sets and Inputs

There are various sources to get stock trading. We will use Yahoo finance to get

- S&P daily data for 15 years (involving both bear and bull markets)
- One can specify the ETF, stocks in instances.dict (explained in later sessions) and the framework dynamically get the data
- Daily data for about 10 ETFs and several large/med/small companies involving various industries

Note that the inputs will contain multiple metrics, such as opening price (Open), highest price the stock traded at (High), how many stocks were traded (Volume) and closing price adjusted for stock splits and dividends (Adjusted Close); your system only needs to predict the Adjusted Close price.

4. Benchmark

S&P index fund performance used as benchmark to compare the performance of various learners and strategies. In this project, we have implemented three learning agents and 3 types of strategies that utilize different combinations of learning agents.

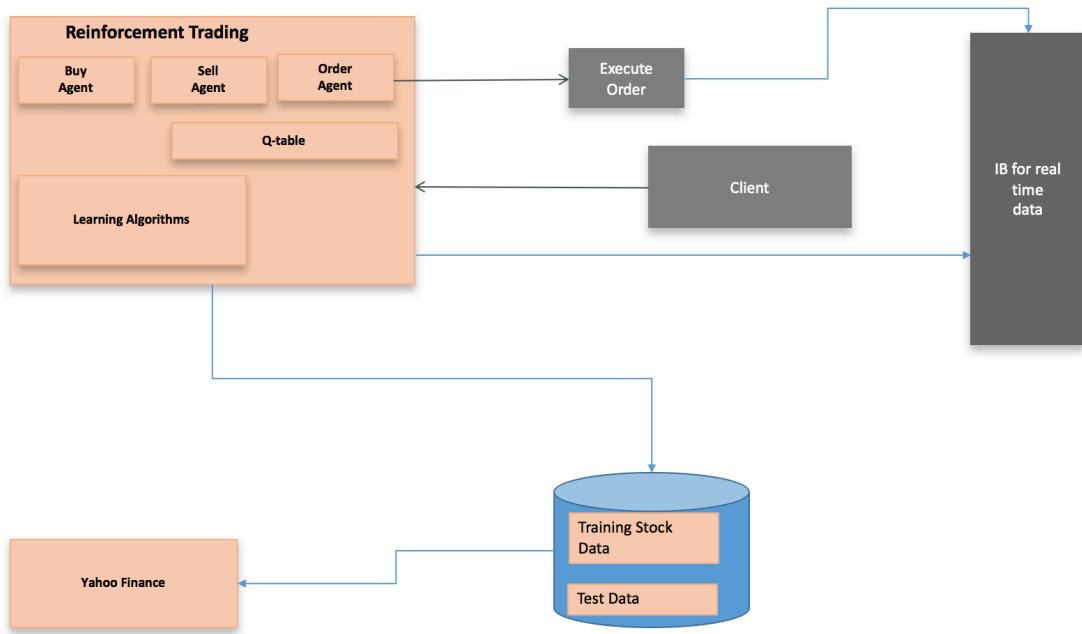
One can specify the benchmark data set as part of instance.dict. For example, following is an example in which SPY is specified as benchmark.

```
"instance2" : {
    "experiment"      : "experiment1",
    "tickerWeights"  : {
        "SPY": 0.20,
        "IJS": 0.05,
        "EFA": 0.15,
        "EEM": 0.05,
        "AGG": 0.20,
        "JNK": 0.05,
        "DJP": 0.10,
        "GLD": 0.20
    },
    "benchmark"       : "SPY",
    "startDate"       : "Dec 2010 04:00AM",
    "endDate"         : "Oct 2016 12:00AM",
    "strategy"        : "monthlyTradeStrategy_randomActionAgent",
    "equity"          : 50000.00,
    "status"          : "A",
    "alpha_params"    : [0.2, 0.5, 0.9],
    "gamma_params"   : [0.2, 0.8],
    "epsilon_params" : [0.2, 0.8],
    "total_trials"   : 1
},
```

5. Design

Our goal is to achieve maximum gain in a year and beat the market average. This framework aims to maximize the profits of investments by considering not only monthly trends of stock prices but also day price movements of stocks. I have instrumented weekly, monthly and quarterly strategies which will use price action based indicators coming from weekly, monthly and quarterly price change indicators.

There are 3 types of agents and each agent has its own goal to achieve and interacts with others to share episodes in the learning process:



5.1. Random Action Agent

Buy and SELL decision is performed by random action. In this project, we have used Random Action agent for 3 different strategies – weekly, monthly and quarterly. You can configure and change [add more strategies] this using instances, strategies.dict files in data/strategies directory

5.2. Random Q-Learning Agent

Buy and SELL decision is performed by using Q-table. However, the reward function uses random approach to decide the reward.

In this project, we have used Random Action agent for 3 different strategies – weekly, monthly and quarterly. You can configure and change [add more strategies] this using instances, strategies.dict files in data/strategies directory

5.3. Flexible Q-learning Agent

Buy and SELL decision is performed by Agent that uses reinforcement technique using Q-table. The reward function uses several indicators and associated signals to derive the appropriate reward.

The {state} variable used can have any number of signals based on the strategy used. I have developed flexible, dynamic way of adding more state member which then can instructed via the strategy class. With this approach, I can add a new strategy class which can opt to use as many indicators for state member and as many indicators for reward calculation.

States, Actions, and Rewards

One of the most important keys to achieve reasonable performance in machine learning is the representation of the input space.

The calculation of reward is as follows.

The agent is given zero reward for the first BUY signal calculation of its reward is postponed until the agent sells the stock (or sell signal is received)

The rate of daily changes of stock prices are given while the agent takes HOLD. But when it takes SELL, zero reward is given because the signal means that it exits the market.

When the sell price is determined by the agent, the buy signal agent receives the profit rate, considering the transaction cost (*transaction cost is not yet instrumented in my code yet*), as the reward.

Learning Algorithm

An episode starts with a stock at a certain day chosen by the environment (input configuration). If the agent takes NOT-BUY as a response to the state of the given stock, the episode ends and a new episode starts.

If the stock cannot be purchased, the episode ends and invokes the buy signal agent with 0 as the reward.

Each agent has equations for reward, delta and Q-value, but all of them are expressed in the same notation for the brevity. Time index is also

abbreviated for the same reason. If the agent fails to sell the stock with its offer price, it sells the stock at close, the last traded price of the day. The update rules and function approximation are based on the ordinary Q-learning.

TD Sequential vs conventional momentum indicators

Conventional momentum indicators such as RSI are typically calculated between 0 to 100 have constraint of overbought and oversold zones which makes them less reliable when price action switches between ranges and trends. TD Setup indicators, on the other hand adjust dynamically in line with the price action.

An Q-learning example with ϵ -greedy policy:

Use the Q-learning update function:

$$Q(s,a) = Q(s,a) + \alpha * (R + \gamma * \max(Q(a,a')) - Q(s,a))$$

where R is defined as:

$$R = r_t * \text{state} * \text{wealth}_t * |\text{state-action}|$$

and r is the percentage return at time t.

Action:

{BUY, SELL, HOLD}

No of States in Q-Learning Strategy 1

- Moving average is used as a primary technique to decide BUY/SELL/HOLD.
- If the short window moving average > long window average, then there is a possibility of a trend reversal
- This is one of the basic and well known technical indicators

Following are the fields used to define the state

```
State = { prev_action, sw_gt_lw_moving_average_signal, position_change_per }
```

No of states in the state definition I have selected is

$$\begin{aligned} &= \text{No of positions *} \\ &\quad [\#\text{of possible prev action}] * \\ &\quad [\#\text{sw_gt_lw_moving_average_signals}] * \\ &\quad [\text{Position \% increase (round to digit)}] \\ \\ &= \text{No of positions * [BUY, SELL, NONE]} * [\text{YES, NO}] * [1 \dots 100] \end{aligned}$$

If we assume that No of positions for our experiment is around 10 for diversified EFT based portfolio, then total no of states is

$$\begin{aligned} &= 10 * 3 * 2 * 100 \\ &= 6000 \end{aligned}$$

The total (max) number of states is 6000 and is reasonable to hold in a memory. Q-table data structure is compact and for each run look up time to get the Q info will be efficient.

No of States in Q-Learning Strategy 2

- In this strategy, two crossover signals are used. It is very similar to Q-Learning strategy 1 but the idea is to combine multiple indicators which helps in deciding the trade.
- The indicators used here could RSI, MACD and any other momentum indicators. I have so far developed only RSI based indicators (14_20 and 50_200) but one can add more those indicators easily with this framework. Please note that when you have more indicators, of course, the State table will become huge.

State:

```
State = { prev_action, rsa_14_50_days_crossover_signal,
rsa_14_50_days_crossover_signal, position_change_per }
```

No of states in the state definition I have selected is

$$\begin{aligned}
&= \text{No of positions} * \\
&\quad [\text{#of possible prev action}] * \\
&\quad [\# rsa_14_50_days_crossover_signal] * \\
&\quad [\# rsa_14_50_days_crossover_signal] * \\
&\quad [\text{Position \% increase (round to digit)}] \\
\\
&= \text{No of positions} * [\text{BUY, SELL, NONE}] * [\text{YES, NO}] * [\text{YES, NO}] * [\text{1 ... 100}]
\end{aligned}$$

If we assume that No of positions for our experiment is around 10 for diversified EFT based portfolio, then total no of states is

$$\begin{aligned}
&= 10 * 3 * 2 * 2 * 100 \\
&= 12000
\end{aligned}$$

No of States in Q-Learning Strategy 3

In this approach, new kind of indicator is added. It is called de-mark TD Sequence indicator.

```
State = { prev_action, rsa_14_50_days_crossover_signal,
rsa_14_50_days_crossover_signal, TD_seq_indicator, position_change_per }
```

No of states in the state definition I have selected is

$$\begin{aligned}
&= \text{No of positions} * \\
&\quad [\text{#of possible prev action}] * \\
&\quad [\# rsa_14_50_days_crossover_signal] * \\
&\quad [\# rsa_14_50_days_crossover_signal] * \\
&\quad [\# TD_seq_indicator] * \\
&\quad [\text{Position \% increase (round to digit)}] \\
\\
&= \text{No of positions} * [\text{BUY, SELL, NONE}] * [\text{YES, NO}] * [\text{YES, NO}] * \\
&\quad [\text{YES, NO}] * [\text{1 ... 100}]
\end{aligned}$$

If we assume that No of positions for our experiment is around 10 for diversified EFT based portfolio, then total no of states is

$$\begin{aligned}
 &= 10 * 3 * 2 * 2 * 2 * 100 \\
 &= 24000
 \end{aligned}$$

6. Evaluation metrics

The objective of the above reinforcement algorithm based approach is to integrate various optimal policies and learn for maximum performance gain.

Each run will produce a report which includes

- Total Return,
- CAGR
- Sharp Ratio, Annual Volatility, R-Squared and Maximum Daily Draw Down



In our experiment, you will see several of these reports involving 3 different strategies trying 3 different learning algorithms with several possible combinations of alpha, gamma and epsilon. You can find the reports in data/reports directory

7. Input and running the algorithm

7.1. Indicators

- data:strategies/indicators.dict
- One can add new indicators under smartrader/indicator directory. For example, you will see rsi_indicators.py and demark-indicators.py
- You can add new indicators in these files or create a new file and add indicators
- The new indicators created should be defined in indicatros.dict for to get processed

```
{  
    "rsi_indicators.RSI14" : {  
        "indicator"      : "rsi_indicators.RSI14",  
        "window"         : 14,  
        "status"         : "A",  
        "dependencies"   : []  
    },  
    "rsi_indicators.RSI50" : {  
        "indicator"      : "rsi_indicators.RSI50",  
        "window"         : 50,  
        "status"         : "A",  
        "dependencies"   : []  
    },  
    "rsi_indicators.RSICrossOver" : {  
        "indicator"      : "rsi_indicators.RSICrossOver",  
        "status"         : "A",  
        "dependencies"   : ["rsi_indicators.RSI14", "rsi_indicators.RSI50"]  
    },  
    "rsi_indicators.DailyValueChangeCalculator" : {  
        "indicator"      : "rsi_indicators.DailyValueChangeCalculator",  
        "window"         : 50,  
        "status"         : "A",  
        "dependencies"   : []  
    },  
    "rsi_indicators.MonthlyValueChangeCalculator" : {  
        "indicator"      : "rsi_indicators.MonthlyValueChangeCalculator",  
        "window"         : 50,  
        "status"         : "A",  
        "dependencies"   : []  
    },  
    "rsi_indicators.QuarterlyValueChangeCalculator" : {  
        "indicator"      : "rsi_indicators.QuarterlyValueChangeCalculator",  
        "window"         : 50,  
        "status"         : "A",  
        "dependencies"   : []  
    },  
    "rsi_indicators.movingAverage" : {  
        "indicator"      : "rsi_indicators.movingAverage",  
        "window"         : 50,  
        "status"         : "A",  
        "dependencies"   : []  
    }  
}
```

7.2. Strategies

- data/strategies/strategies.dict
- One can add new strategy under smartrader/strategies directory. For example, you will see following strategies now.
 - o monthlyTradeStrategy.py
 - o quaterlyTradeStrategy.py
 - o simpleBuyAndHold.py
 - o weeklyTradeStrategy.py
- You can add new strategies in these files or create a new file and add strategies
- The new indicators created should be defined in strategies.dict for to get processed

```

"weeklyTradeStrategy_learningAgent1": {
    "positionSizer" : "rebalance.LiquidateRebalancePositionSizer",
    "riskManager" : "example.ExampleRiskManager",
    "portfolioHandler" : "portfolio_handler.PortfolioHandler",
    "complianceHandler" : "example.ExampleCompliance",
    "executionHandler" : "ib_simulated.IBSimulatedExecutionHandler",
    "title" : "weeklyTradeStrategy_learningAgent1",
    "tearsheetStatistics" : "tearsheet.TearsheetStatistics",
    "displayStrategy" : "DisplayStrategy",
    "strategyHandler" : "weeklyTradeStrategy",
    "learningHandler" : "learningAgent1"
},
"monthlyTradeStrategy_learningAgent1": {
    "positionSizer" : "rebalance.LiquidateRebalancePositionSizer",
    "riskManager" : "example.ExampleRiskManager",
    "portfolioHandler" : "portfolio_handler.PortfolioHandler",
    "complianceHandler" : "example.ExampleCompliance",
    "executionHandler" : "ib_simulated.IBSimulatedExecutionHandler",
    "title" : "monthlyTradeStrategy_learningAgent1",
    "tearsheetStatistics" : "tearsheet.TearsheetStatistics",
    "displayStrategy" : "DisplayStrategy",
    "strategyHandler" : "monthlyTradeStrategy",
    "learningHandler" : "learningAgent1"
},
"quaterlyTradeStrategy_learningAgent1": {
    "positionSizer" : "rebalance.LiquidateRebalancePositionSizer",
    "riskManager" : "example.ExampleRiskManager",
    "portfolioHandler" : "portfolio_handler.PortfolioHandler",
    "complianceHandler" : "example.ExampleCompliance",
    "executionHandler" : "ib_simulated.IBSimulatedExecutionHandler",
    "title" : "quaterlyTradeStrategy_learningAgent1",
    "tearsheetStatistics" : "tearsheet.TearsheetStatistics",
    "displayStrategy" : "DisplayStrategy",
    "strategyHandler" : "quaterlyTradeStrategy",
    "learningHandler" : "learningAgent1"
},

```

7.3. Instances

- dataestrategiesinstances.dict
- instances configuration determines the scope of actual run.
- runMain and framework parses this and run the experiment for the instances and input data set.
- You can add any number of instances and the framework takes care of running all instances one by one

```

"instance4" : {
    "experiment"      : "experiment1",
    "tickerWeights" : {
        "SPY": 0.7,
        "GDX": 0.1,
        "GLD": 0.2
    },
    "benchmark"       : "SPY",
    "startDate"       : "Dec 2010 04:00AM",
    "endDate"         : "Oct 2016 12:00AM",
    "strategy"        : "weeklyTradeStrategy_learningAgent1",
    "equity"          : 50000.00,
    "status"          : "D",
    "alpha_params"   : [0.2],
    "gamma_params"  : [0.2],
    "epsilon_params" : [0.2],
    "total_trials"   : 10
},
"instance5" : {
    "experiment"      : "experiment1",
    "tickerWeights" : {
        "SPY": 0.20,
        "IJS": 0.05,
        "EFA": 0.15,
        "EEM": 0.05,
        "AGG": 0.20,
        "JNK": 0.05,
        "DJP": 0.10, // DJP
        "GLD": 0.20
    },
    "benchmark"       : "SPY",
    "startDate"       : "Dec 2010 04:00AM",
    "endDate"         : "Oct 2016 12:00AM",
    "strategy"        : "monthlyTradeStrategy_learningAgent1",
    "equity"          : 50000.00,
    "status"          : "A",
    "alpha_params"   : [0.2, 0.5, 0.9],
    "gamma_params"  : [0.2, 0.8],
    "epsilon_params" : [0.2, 0.8],
    "total_trials"   : 10
},

```

8. Automation tools, logs and embedded runtime statistics

As part of this project, I have developed several automation scripts that enables one to run experiment with several different alpha,

gamma and epsilon values in a flexible way that can be initiated from command line. This automaton enables one to initiate an experiment with different trials, hyperparameters without changing the code.

Moreover, I have instrumented the code base (agent.py, environment.py) etc to record run time statistics about the trial and its progress. Upon completion of experiment, an analysis document is created (by analysis the results program) which summarizes the top trials/experiments having high rewards and completion rate.

9. Implementation & Code snippets

9.1. Framework

- Framework consists of the main program and many sub components
- Each sub component can be considered as a plugin
- All the plugins for each run are assembled in Agent Factory
- Framework takes care of persisting the results in data/results directory
- Framework can be started by app/runMain.py

9.1.1. Agent Factory

- Reads the plugins and configuration and assemble/wire the plugin components
- As you see below, there are various plugins such as position_sizer, risk_manager, indicators, strategies, display, statistics etc.
- It is the responsibility of Agent Factory to assemble this based on strategy.dict and make it ready for the run

```

74     tmpArray=strategyHandler.split(".")
75     # Use the monthly liquidate and rebalance strategy
76     strategyClass = getattr(importlib.import_module("smarttrader.strategy."+tmpArray[0]), "Strategy")
77     displayClass = getattr(importlib.import_module("smarttrader.strategy.display"), displayStrategy)
78
79     self.strategy = strategyClass(self.tickers, self.events_queue)
80     strategy = Strategies(self.strategy, displayClass())
81
82     # Use the liquidate and rebalance position sizer
83     # with prespecified ticker weights
84     tmpArray = positionSizer.split(".")
85     rebalancingClass = getattr(importlib.import_module("smarttrader.position_sizer."+tmpArray[0]), tmpArray[1])
86     self.positionSizer = rebalancingClass(ticketWeights)
87
88     # Use an example Risk Manager
89     tmpArray = riskManager.split(".")
90     riskManagerClass = getattr(importlib.import_module("smarttrader.risk_manager."+tmpArray[0]), tmpArray[1])
91     self.risk_manager = riskManagerClass()
92
93     # Use the default Portfolio Handler
94     tmpArray = portfolioHandler.split(".")
95     portHandlerClass = getattr(importlib.import_module("smarttrader.portfolio."+tmpArray[0]), tmpArray[1])
96     self.portfolio_handler = portHandlerClass(self.initial_equity, self.events_queue, self.price_handler,
97     | self.positionSizer, self.risk_manager)
98
99     # Use the ExampleCompliance component
100    tmpArray = complianceHandler.split(".")
101    complianceClass = getattr(importlib.import_module("smarttrader.compliance."+tmpArray[0]), tmpArray[1])
102    self.compliance = complianceClass(config)
103
104    # Use a Indicator component
105    self.indicator_list = list()
106    self.indicator_dict = dict()
107    for indicator, value in indicators.iteritems():
108        print indicator, value
109        tmpArray = value['indicator'].split(".")
110        indicatorClass = getattr(importlib.import_module("smarttrader.indicators."+tmpArray[0]), tmpArray[1])
111        instance = indicatorClass(self.tickers, self.events_queue,value)
112        self.indicator_list.append(instance)
113        self.indicator_dict[value['indicator']] = instance
114
115    # Use a simulated TB Execution Handler
116

```

9.2. Playback Simulator

- Very important part of the application
- It basically invokes the pricing/market data provider and for each data tick
- Starting from the start date of the trading, it calls the indicators, process the indicators for the tick and call appropriate strategy class.

```

print("Running Backtest...")
while self.price_handler.continue_backtest:
    try:
        event = self.events_queue.get(False)
    except queue.Empty:
        self.price_handler.stream_next()
    else:
        if event is not None:

            print("----- Price EVENT Received from market backtest-----")
            if event.type == EventType.TICK or event.type == EventType.BAR:
                self.cur_time = event.time

            graph = dict()
            indicatorPath = getPath('strategies.data.path')
            indicators = json.load(open(indicatorPath + "indicators.dict"))
            for indicator, value in indicators.items():
                dependencies = value['dependencies']
                graph[indicator] = dependencies

            order = topoSort.topological(graph)
            while True:
                try:
                    ind = order.pop()
                    print(self.indicator_dict)
                    indClass = self.indicator_dict[ind]
                    indClass.updateIndicators(event, self.agent)

                except IndexError:
                    break

                print("      **** Processing signals -----")
                self.strategy.calculate_signals(event, self.agent)
                print("      **** Update portfolio -----")
                self.portfolio_handler.update_portfolio_value()
                print("      **** Update Stats -----")
                self.statistics.update(event.time, self.portfolio_handler)

            elif event.type == EventType.SIGNAL:
                self.portfolio_handler.on_signal(event)
            elif event.type == EventType.ORDER:
                self.execution_handler.execute_order(event)
            elif event.type == EventType.FILL:
                self.portfolio_handler.on_fill(event)

```

9.2.1. Indicators

Framework invokes Indicator classes for two reasons. First one is to process the market data and identify the signals (such as SELL or BUY). Second one is to provide the signal up on the request.

```

06     self.invested = False
07
08     self.sw_bars = dict()
09     self.lw_bars = dict()
10     self.signal = dict()
11
12     self.tickers = tickers
13     self.events_queue = events_queue
14     self.ticks = 0
15
16     self.period = indicator_details['window']
17
18     for ticker in tickers :
19         self.sw_bars[ticker] = deque(maxlen=self.short_window)
20         self.lw_bars[ticker] = deque(maxlen=self.long_window)
21
22     for ticker in tickers :
23         self.signal[ticker] = "NONE"
24
25
26     def getSignal(self, ticker):
27         print "TIC", ticker, self.signal
28         return SignalEvent(ticker, self.signal[ticker])
29
30     def updateIndicators(self, event, agent=None):
31
32         if event.type in [EventType.BAR, EventType.TICK] and event.ticker in self.tickers :
33
34             self.lw_bars[event.ticker].append(event.adj_close_price)
35             self.sw_bars[event.ticker].append(event.adj_close_price)
36
37
38             #Enough_bars_are_present_for_trading
39             if len(self.lw_bars[event.ticker]) == self.lw_bars[event.ticker].maxlen:
40                 # Calculate the simple moving averages
41                 short_sma = np.mean(self.sw_bars[event.ticker])
42                 long_sma = np.mean(self.lw_bars[event.ticker])
43                 # Trading signals based on moving average cross
44                 if short_sma > long_sma and not self.invested:
45                     print("LONG: %s" % event.time)
46                     self.signal[event.ticker] = "BOT"
47
48             elif short_sma < long_sma and self.invested:
49                 print("SHORT: %s" % event.time)

```

9.3. Agents

9.3.1. Random Action Agent

```

b
7     class Action():
8
9         def __init__(self, context, instanceDict):
0             self.state = State()
1
2         def reset(self, destination=None):
3             print "Inside reset"
4
5         def getTotalRewards(self):
6             return 0
7
8         def getLegalActions(self):
9             """
0                 returns the legal action from the current state
1             """
2             return ['BOT', 'SLD', 'NONE']
3
4         def update(self, event, learning_indicators, reward_indicators, indicator_dict, portfolio):
5
6             pri Parameter 'learning_indicators' value is not used more... (⌘F1) dict
7
8             # Select action according to your policy
9             action = random.choice(list(self.getLegalActions()))
0
1
2             try:
3                 print " Action ", action
4                 if DEBUG:
5                     input("")
6
7             except SyntaxError:
8                 pass
9
0
1
2             return action
3
4
5
6
7
8
9
0
1

```

9.3.2. Simple/Random Q-Learning Agent

```

def getTotalRewards(self):
    return self.totalRewards

def act(self, event, action, reward_indicators, indicator_dict):
    # random reward. You would see different reward structure based on portfolio/position increase in learningAgent1
    return random.choice(list([0.5, 0.8, 2, 1, 0.2, 0.7, 2, 3]))

```

9.3.3. Flexible Q-learning Agent

```

11     class State():
12
13         def __init__(self, ticker, signals=None, lastAction=None):
14             self.signals = signals
15             self.ticker = ticker
16             self.lastAction = lastAction
17
18         def getStateString(self):
19
20             print "Inside learningAgent1:getStateString()"
21             returnStr = self.ticker + ":" +
22             for key, value in self.signals.iteritems():
23                 returnStr += " " + str(key) + ":" + str(value.action)
24             returnStr += "—" + self.lastAction
25
26             return returnStr

```

```

71
72     def update(self, event, learning_indictors, reward_indicators, indicator_dict, portfolio):
73
74         print "Inside learningAgent1:Action:Update:", event, indicator_dict
75
76         signals = dict()
77         for ln in learning_indictors :
78             signals[ln] = indicator_dict[ln].getSignal(event.ticker)
79
80         # Update state from env. In this case, backtest/replay will call the update function
81         tempState = State(event.ticker, signals,"NONE")
82         print "current STATE:", tempState.getStateString()
83
84         # Select action according to your policy
85         action = self.qTable.chooseAction(tempState)
86
87         # Execute action and get reward
88         reward = self.act(event, action, reward_indicators, indicator_dict)
89
90         try:
91             print " Action ", action
92             print " Reward for action ", reward
93             if DEBUG:
94                 input("")
95
96         except SyntaxError:
97             pass
98
99         newstate = State(event.ticker,signals,action)
100
101        # Learn policy based on state, action, reward
102        self.qTable.updateQTable1(self.state,
103                                action=action,
104                                reward=reward,
105                                newstate=newstate)
106        self.state = newstate
107        self.totalRewards += reward
108
109
110

```

```

def getTotalRewards(self):
    return self.totalRewards

def act(self,event, action,reward_indicators, indicator_dict):
    # As of now, only one reward indicator can be passed.
    value = indicator_dict[reward_indicators[0]].getSignal(event.ticker)
    print "% chaned value:", value

    ## Since, I am passing weekly performance/value_change as reward indicator
    ## the change is not expected to be more than 10%. However, if I start using monthly or Qty reward indicator
    ## the value likely to be higher. In that case, we need to normalize. The idea here is tha
    ## if you get 10% up in a week, the reward given is maximum -which is 1

    if value > 10:
        return 1
    elif value > 5 :
        return 0.5
    elif value > 2 :
        return 0.1

```

```

163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195 def chooseAction(self, state):
196     print "\n"
197     if True or random.random() < self._epsilon:
198         legalActions = self.getLegalActions()
199         bestAction = random.choice(legalActions)
200         print " Going for random action :", bestAction
201         return bestAction
202     else:
203
204         maxQ = 0.0
205         bestAction = "NONE"
206         for action in self.getLegalActions():
207             q = self.getQValue(state,action=action)
208             if q >= maxQ :
209                 bestAction = action
210                 maxQ = q
211             print " Going for best action for the state:", state.getStateString(), "=>, action,"\\t\\t", q, ":" , \
212                 maxQ
213
214         if maxQ == 0.0 :
215             legalActions = self.getLegalActions()
216             bestAction = random.choice(legalActions)
217             print " No prev state-action found. Going for random action for the state:", bestAction
218
219         else :
220             print " +++++++ Self Learned action for the state:", bestAction, "having Q:", maxQ
221
222     return bestAction
223
224
225 def updateQTable1(self,state,action,reward,newstate):
226
227     oldQvalue = self.getQValue(state,action=action)
228     maxValue = self.getMaxQValue(state)
229
230     newQvalue = oldQvalue + self._alpha * ( (reward + self._gamma * maxValue) - oldQvalue)
231     self.setQValue(state,action,new_value=newQvalue)
232
233     print " QTable : "
234     print json.dumps(self._table, indent=4)

```

10. Run & Observation

10.1. Random Action Agent

Weekly Rebalancing

BuyAndHoldStrategy



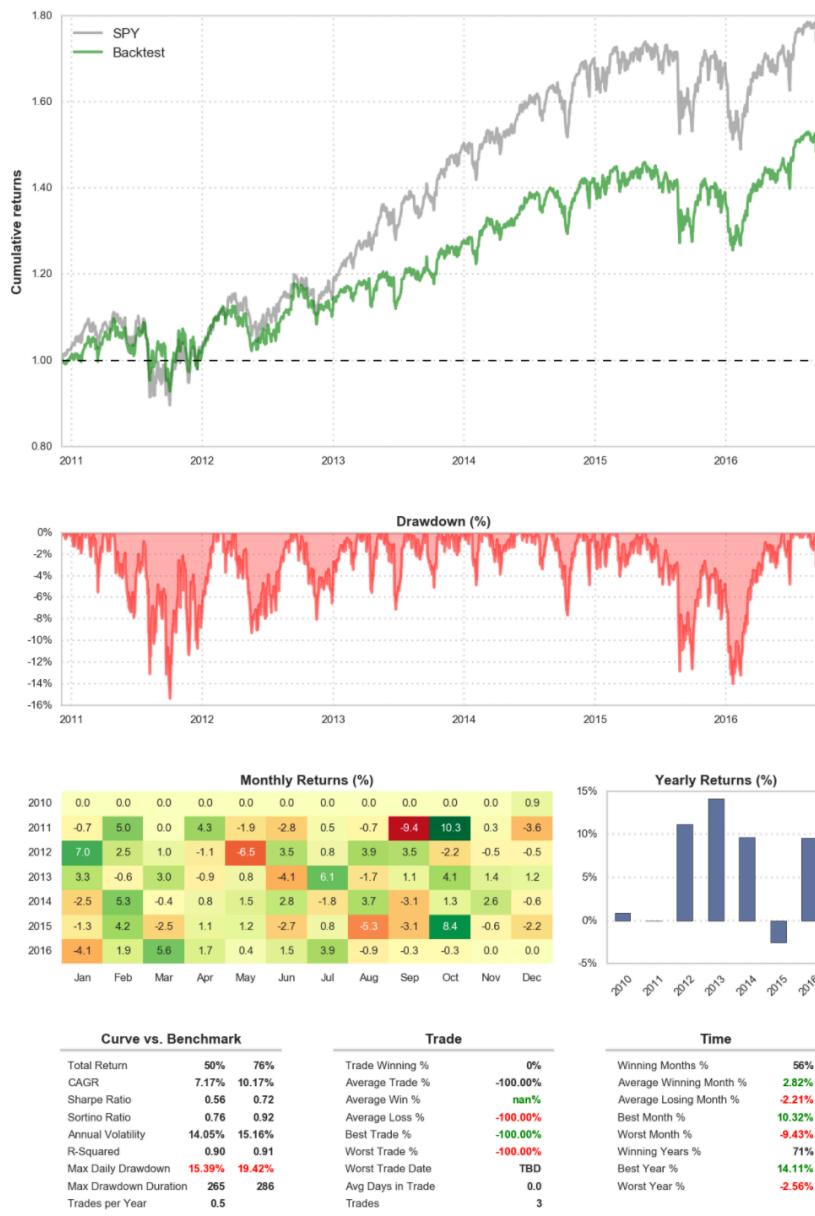
Monthly Returns (%)												
2010	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec
0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.6
2011	1.6	2.4	-0.5	2.6	-0.4	-1.8	-0.8	-3.2	-4.7	7.0	-0.6	0.7
2012	3.5	2.8	1.5	-0.2	-3.8	2.3	1.3	1.6	1.4	-1.4	0.4	-0.3
2013	3.3	1.1	2.3	1.7	0.8	-2.0	3.8	-2.6	2.3	3.7	2.1	1.2
2014	-2.1	3.5	0.2	0.7	2.1	1.1	-1.2	3.3	-1.7	2.1	2.3	-0.8
2015	-1.6	3.9	-1.5	0.6	0.8	-2.4	2.0	-4.9	-2.1	6.4	0.1	-2.0
2016	-3.4	0.2	4.9	0.3	1.3	0.4	2.9	-0.0	-1.0	0.0	0.0	0.0



Curve vs. Benchmark		Trade		Time		
Total Return	54%	75%	Trade Winning %	0%	Winning Months %	63%
CAGR	7.69%	10.11%	Average Trade %	-100.00%	Average Winning Month %	2.09%
Sharpe Ratio	0.76	0.71	Average Win %	nan%	Average Losing Month %	-1.80%
Sortino Ratio	0.99	0.92	Average Loss %	-100.00%	Best Month %	7.05%
Annual Volatility	10.54%	15.16%	Best Trade %	-100.00%	Worst Month %	-4.90%
R-Squared	0.92	0.91	Worst Trade %	-100.00%	Winning Years %	86%
Max Daily Drawdown	12.01%	19.42%	Worst Trade Date	TBD	Best Year %	18.86%
Max Drawdown Duration	284	286	Avg Days in Trade	0.0	Worst Year %	-1.26%
Trades per Year	0.3		Trades	2		

Monthly Rebalancing

weeklyTradeStrategy_randomActionAgent



Quarterly rebalancing

Not included in the word doc to reduce space. Please refer data/results if you want to see the details

10.2. Simple/Random Reward Q-Learning Agent

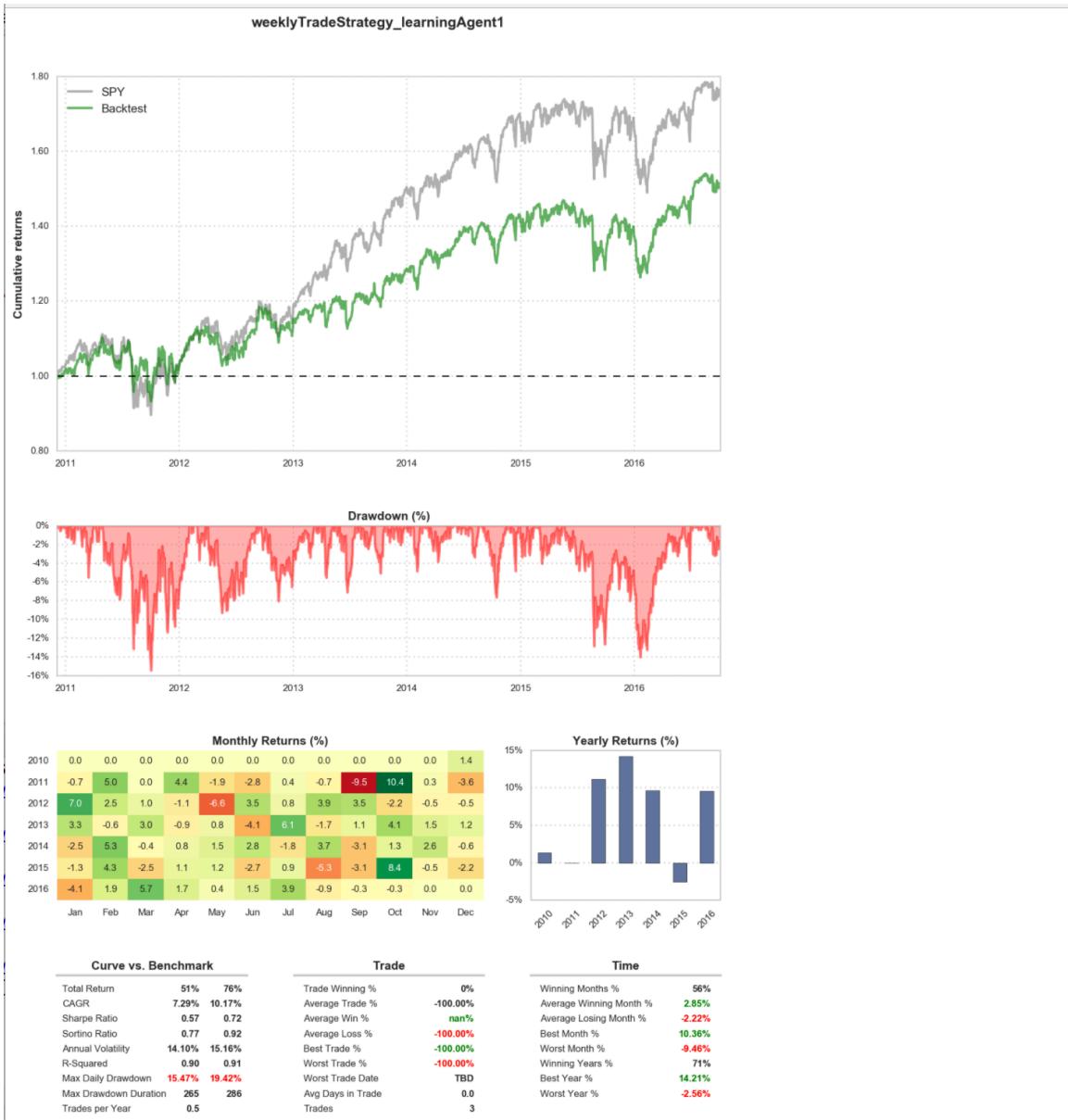
Not included in the word doc to reduce space. Please refer data/results if you want to see the details

10.3. Flexible Q-learning Agent

This is the most interesting experiment. First, I ran with low alpha, gamma and Epsilon for many trials. In most cases, I could not beat the market.

Weekly Rebalancing

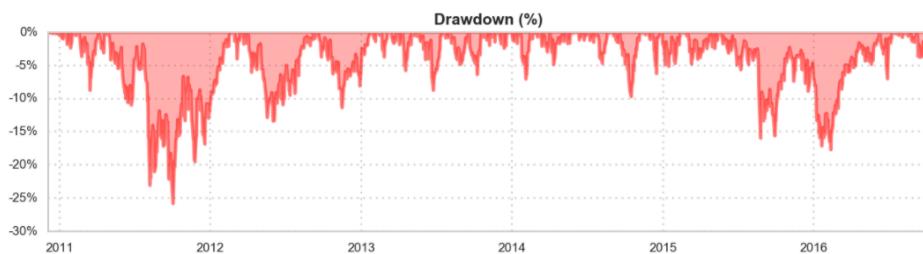
Alpha 0.2, Gamma, 0.2 Epsilon 0.2



I observed that as I increase Gamma and Epsilon, the results started to looking better. After many trials, following is the best performance I observed.

Alpha 0.5, Gamma, 0.2 Epsilon 0.2 and Iteration 10

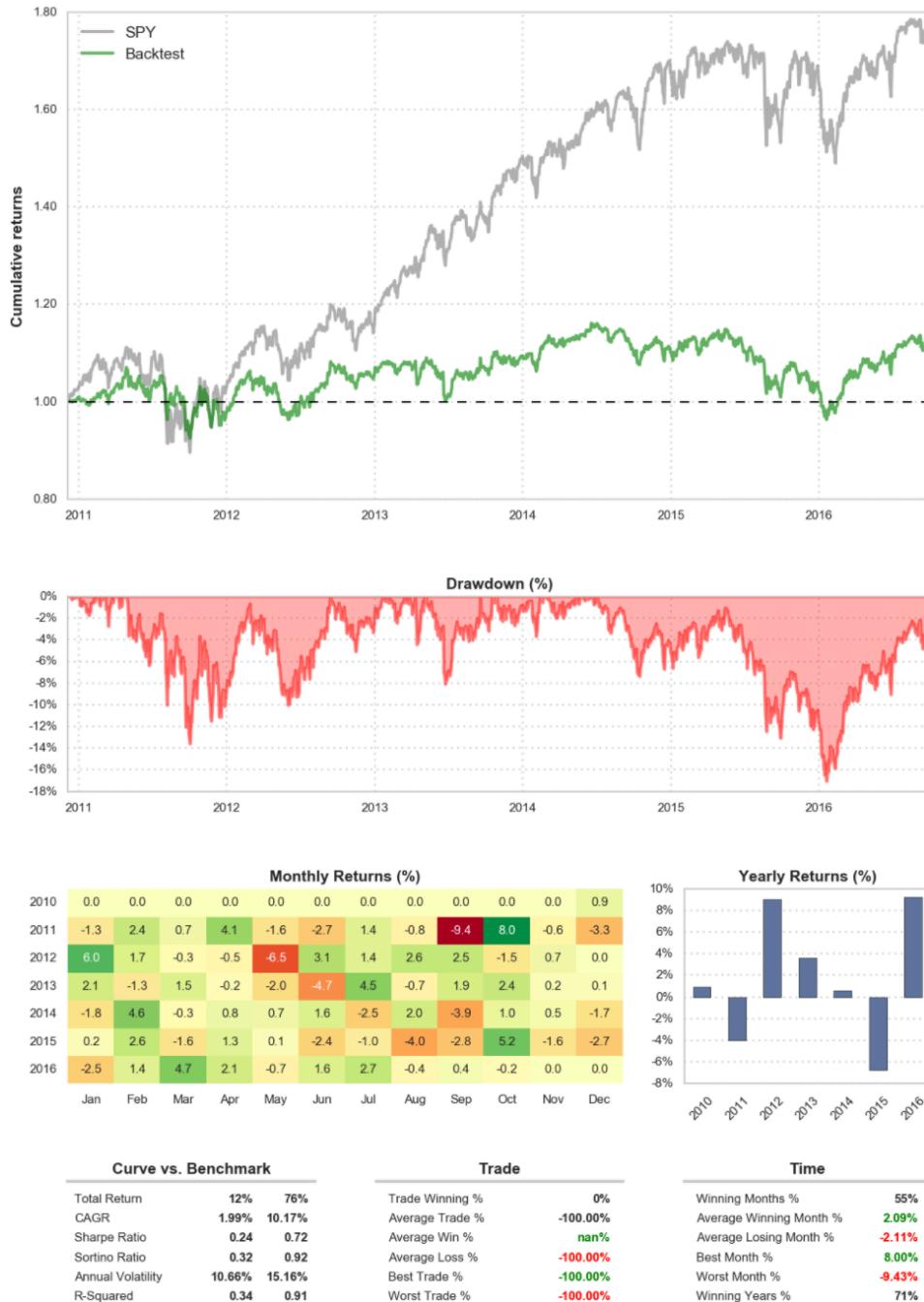
weeklyTradeStrategy_learningAgent1



Curve vs. Benchmark		Trade		Time		
Total Return	92%	76%	Trade Winning %	0%	Winning Months %	59%
CAGR	11.90%	10.17%	Average Trade %	-100.00%	Average Winning Month %	3.86%
Sharpe Ratio	0.65	0.72	Average Win %	nan%	Average Losing Month %	-3.07%
Sortino Ratio	0.85	0.92	Average Loss %	-100.00%	Best Month %	16.71%
Annual Volatility	20.56%	15.16%	Best Trade %	-100.00%	Worst Month %	-12.50%
R-Squared	0.91	0.91	Worst Trade %	-100.00%	Winning Years %	71%

Monthly Rebalancing

Monthly rebalancing is not helping much. I am yet to try with different params.



Quarterly rebalancing

11. Improve the Q-Learning Driving Agent

- Enhance Q-learning agent so that, after sufficient training, the agent is able to beat the benchmark or perform better.
- Parameters in the Q-Learning algorithm, such as the learning rate (`alpha`), the discount factor (`gamma`) and the exploration rate (`epsilon`) all contribute to the agent's ability to learn the best action for each state.

You can specify the params in `instances.dict` file. For example, for following instance, different values of alpha, gamma and epsilon is provided

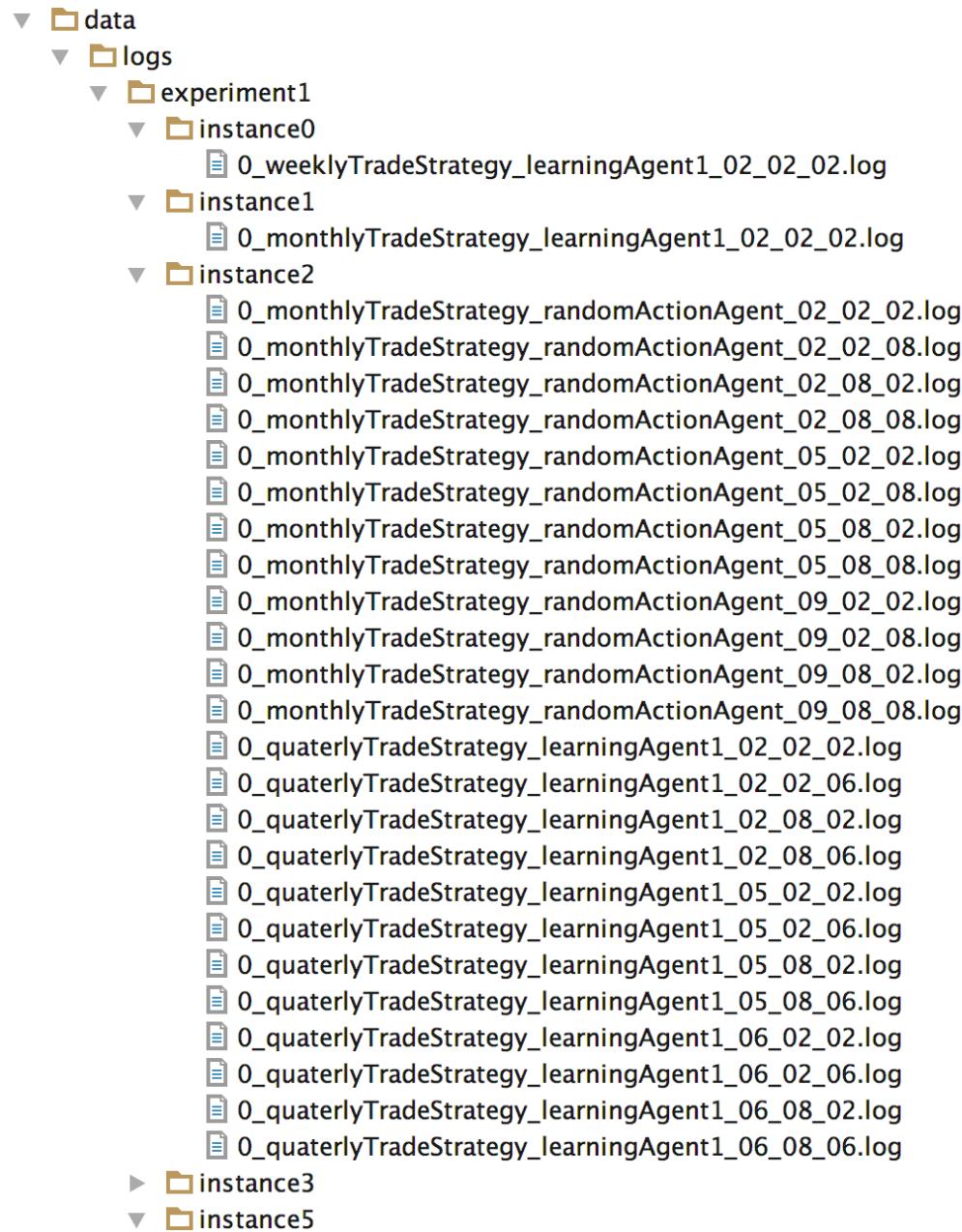
```
"instance5" : {
    "experiment"      : "experiment1",
    "tickerWeights"  : {
        "SPY": 0.20,
        "IJS": 0.05,
        "EFA": 0.15,
        "EEM": 0.05,
        "AGG": 0.20,
        "JNK": 0.05,
        "DJP": 0.10,
        "GLD": 0.20
    },
    "benchmark"       : "SPY",
    "startDate"       : "Dec 2010 04:00AM",
    "endDate"         : "Oct 2016 12:00AM",
    "strategy"        : "monthlyTradeStrategy_learningAgent1",
    "equity"          : 50000.00,
    "status"          : "A",
    "alpha_params"    : [0.2, 0.5, 0.9],
    "gamma_params"   : [0.2, 0.8],
    "epsilon_params" : [0.2, 0.8],
    "total_trials"   : 10
},
```

Run the simulation. Adjust one or several of the above parameters and iterate this process.

With automation tools, logs and runtime stats and analysis over runtime stats, we will run two experiments. Each experiments consist of several trials of various combination of alpha, gamma and epsilon.

Now, when you run mainDriver.py it will automatically run for various combination of alpha, gamma and epsilon from the above sets.

For each run, the corresponding log file is created under data/logs

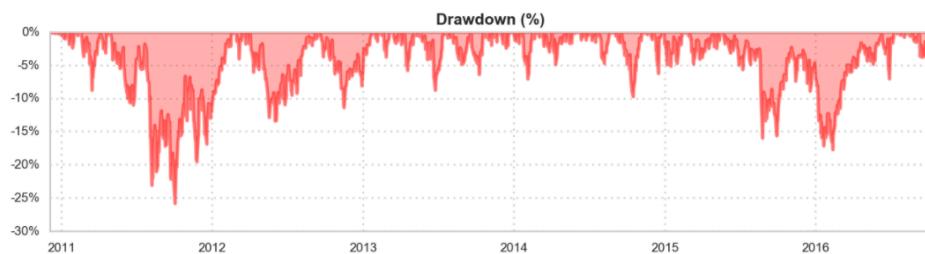


After completing the experiment, grep on Total rewards under logs directory.

Best RUN from many trials...

Alpha 0.5, Gamma, 0.2, Epsilon 0.2 and Iteration 10

weeklyTradeStrategy_learningAgent1



Curve vs. Benchmark		Trade		Time		
Total Return	92%	76%	Trade Winning %	0%	Winning Months %	59%
CAGR	11.90%	10.17%	Average Trade %	-100.00%	Average Winning Month %	3.86%
Sharpe Ratio	0.65	0.72	Average Win %	nan%	Average Losing Month %	-3.07%
Sortino Ratio	0.85	0.92	Average Loss %	-100.00%	Best Month %	16.71%
Annual Volatility	20.56%	15.16%	Best Trade %	-100.00%	Worst Month %	-12.50%
R-Squared	0.91	0.91	Worst Trade %	-100.00%	Winning Years %	71%

12. Conclusion

By looking at experiment1 logs results, we can realize that

- Beating stock market is not an easy task
- There are some strategies which beats the market. However please note that we have not considered transaction costs/tax etc.
- In my experiment, weekly rebalancing with multiple indicators (RSI, d-mark TD as signal indicator and weekly performance as reward indicator) with alpha=0.5, gamma=0.2 and epsilon=0.2 give the optimal result.
- Agent with alpha as 0.5 or 0.9 , gamma as 0.2 and epsilon as 0.2 runs with high reward and best completion rate with minimal penalty and best run time.
- Higher alpha means that agent uses higher learning rate. Learning rate controls the learning step size, that is, how fast learning takes place.
- In our experiment, higher alpha does not always result in best completion runs. Seems like 0.5 alpha is optimal. You can review the detail logs to confirm this.
- Future reinforcements are weights controlled by a value gamma between 0 and 1. A higher value of gamma means that the future matters more for the Q-value of a given action in a given state.
- In our experiment, higher gamma does not help much. Gamma with 0.2 value provides the best results.
- The reason for higher gamma is not helpful because the information calculated for reward may be faulty for one reason or another.
- It is better to update more gradually, to use the new information to move in a particular direction, but not to make too strong a commitment.
- Higher epsilon is counterproductive in our experiment. Usually higher epsilon helps when there is more randomness in the behavior of env response.
- Seems like, there is not much of variation of responses from env and

higher epsilon does not help much – hence low epsilon results in better success rate.

Behavior of Agents and reaching the optimal policy

- *Agent with alpha as 0.5 or 0.9 , gamma as 0.2 ad epsilon as 0.2 runs with high reward and best completion rate with minimal penalty and best run time.*
- The optimal policy for this kind of agent is to balance between exploration and exploitation. Alpha with 0.5 and gamma 0.2 achieves the optimal exploitation and 0.2 epsilon is optimal for exploration.

In summary, further study is needed to confirm the results – specifically across multiple bull and bear market.

With this framework, I am planning to further try different indicators and combinations and see whether we can consistently beat the market with a margin that will include transaction costs/tax in beating the market.

13. References

1. Mnih, Volodymyr, et al. “Playing atari with deep reinforcement learning.” *arXiv preprint arXiv:1312.5602* (2013).
2. <http://cs229.stanford.edu/proj2014/David%20Montague,%20Algorithmic%20Trading%20of%20Futures%20via%20Machine%20Learning.pdf>
3. Algorithmic and High-Frequency Trading, **Print ISBN-13:** 978-1-107-09114-6 **By:** Álvaro Cartea; Sebastian Jaimungal; José Penalva
 - a. **Publisher:** Cambridge University Press
4. Analysis of Financial Time Series
 - b. **By:** [RUEY S. TSAY](#)

- c. **Publisher:** John Wiley & Sons
 - d. **Pub. Date:** August 30, 2010
 - e. **Print ISBN:** 978-0-470-41435-4
 - f. **Web ISBN:** 0-470414-35-9
5. Market-Neutral Trading: Combining Technical and Fundamental Analysis Into 7 Long-Short Trading Systems
- g. **By:** Thomas K. Carr
 - h. **Publisher:** McGraw-Hill
 - i. **Pub. Date:** December 17, 2013
 - j. **Print ISBN-13:** 978-0-07-181310-5
6. quantstart.com for python libraries for yahoo data reader and portfolio management