

# REINFORCEMENT LEARNING FOR AUTOMATED TRADING

Reid Falconer<sup>a</sup>, Sam MacIntyre<sup>a</sup>, Hector Cano<sup>a</sup>, Maximilian Zebhauser<sup>a</sup>

*<sup>a</sup>Barcelona Graduate School of Economics, Barcelona, Spain*

---

*Keywords:* Reinforcement Learning, Deep Q-Learning, Automated Trading, Neural Networks

---

## 1. Introduction

Algorithmic trading for stocks is attractive for both researchers and market practitioners. Existing approaches for algorithmic trading can be categorised into knowledge-based methods and machine learning (ML) based methods. Knowledge-based methods design trading strategies based on either financial research or trading experience; ML-based methods, in contrast, learn trading strategies from the historical market data. A distinct advantage of the ML-based methods is that they can discover profitable patterns that are not yet known to people.

Among various ML methods, reinforcement learning (RL) is particularly exciting and is considered a third ML paradigm alongside unsupervised and supervised learning. Nevertheless, unlike the other approaches, RL considers the whole problem of a goal-directed agent that interacts in an uncertain environment. This approach involves learning what actions are necessary to take in order to maximise a numerical reward signal.

The most important distinguishing features of a RL problem are:

1. They behave as closed-loop problems given that its learned actions influence later inputs

---

*Email addresses:* `reid.falconer@barcelonagse.eu` (Reid Falconer),  
`sam.macintyre@barcelonagse.eu` (Sam MacIntyre), `hector.cano@barcelonagse.eu` (Hector Cano),  
`maximilian.zebhauser@barcelonagse.eu` (Maximilian Zebhauser)

2. Learners must try different operations to discover which strategy yields the most reward
3. Actions may affect next situations and all subsequent rewards (Sutton et al., 1998)

Among its main applications are: resources management in computer clusters, games, traffic light control and robotics Mnih et al. (2013). This project aims to apply RL techniques to make decisions in the stock market given that it involves the interaction of an active agent that has to make decisions based on an imperfect information environment while also interacting with other market participants. Some previous findings indicate that RL can be successfully applied to the portfolio problem and its performance exceeds the supervised learning approach (Neuneier, 1996) and Q-learning algorithm operates better than kernel-based methods (Bertoluzzo and Corazza, 2012).

In this paper, we apply the deep Q-learning approach to algorithmic trading. Our goal is to build a deep Q-learning system that determines when to buy, sell or hold based on the current and historical market data. Our experiments on the both Apple (AAPL) and Wawel (WWL) stocks demonstrate that the deep Q-learning system is highly effective and that the deep Q-learning model outperforms the benchmarks such as a random decision policy and a buy and hold strategy.

The paper is organised as follows: Firstly, a conceptual framework is presented, highlighting the underlying principles of reinforcement learning. Section 3 describes the Q-learning approach, and Section 4 presents the implementation details of the deep Q-learning system. Section 5 presents the data, settings and results. Finally, Section 6 concludes.

## 2. Reinforcement Learning

Before delving into the specifics of employing reinforcement learning to the problem of automated trading, it will be informative to discuss the general theory and its underlying principles.

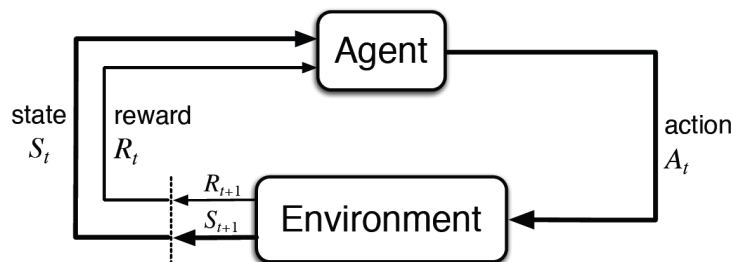
Reinforcement learning aims to maximise a given reward signal by undertaking certain actions (in a restricted space). In this framework, an agent must take the state of the environment as input and take actions to alter the future state. A measurable goal

related to the environment is also necessary for the problem formulation. Beyond this, each reinforcement learning problem contains four sub-elements: a *policy*, a *reward signal* and a *value function* (Sutton et al., 1998).

The policy defines the agent's actions in different environment states. The reward signal defines the goal and should be maximised throughout the learning process. The value function maps the current state to a value so the agent can make optimal longer run decisions. It can be seen as the expected total future reward that can be obtained beginning from that state. Most of the challenges associated with the implementation of reinforcement learning derive from the estimation of the value function.

Formally, we construct a Markov Decision Process (MDP). In an ideal situation, we would have access to the value function directly in tabular form when we have a tractable action and state space.

Figure 1: Agent and Environment Interaction



Source: Sutton et al. (1998)

At a sequence of discrete time steps  $t = 0, 1, 2, 3, \dots$ , the agent interacts with the environment. At each step  $t$ , the agent receives state information  $S_t \in \mathcal{S}$  and performs an action  $A_t \in \mathcal{A}(s)$ . As a consequence of the action, the agent receives a reward  $R_{t+1} \in \mathcal{R} \subset \mathbb{R}$  and transitions to a new state  $S_{t+1}$ .

In the context of a *Markov* decision process, the future rewards ( $R_t$ ) and states ( $S_t$ ) only depend on the previous state and action.

The general reinforcement learning paradigm involves finding an optimal policy  $\pi$  to maximise the expected discounted return. The discount factor is required to ensure that rewards in the distant future are less valuable than current rewards.

$$G_t \doteq R_{t+1} + \gamma R_t + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

Value and action-value functions allow the actions of the agent to be assessed under the implementation of a particular policy. The value function and action-value functions respectively are defined below:

$$v_{\pi}(s) \doteq \mathbb{E}_{\pi} [G_t | S_t = s] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right], \text{ for all } s \in \mathcal{S}$$

$$q_{\pi}(s, a) \doteq \mathbb{E}_{\pi} [G_t | S_t = s, A_t = a] = \mathbb{E}_{\pi} \left[ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

Ideally, the value function is decomposed into the following (known as the *Bellman's Equation*):

$$\begin{aligned} v_{\pi}(s) &\doteq \mathbb{E}_{\pi} [G_t | S_t = s] \\ &= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_{\pi} [G_{t+1} | S_{t+1} = s']] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')], \quad \text{for all } s \in \mathcal{S} \end{aligned}$$

Both expressions relate to a specific state and action taken at any time  $t$ .

A reinforcement learning problem principally involves pursuing the optimal policy  $\pi$  which is said to maximise the value and action-value functions:

$$v_{*}(s) \doteq \max_{\pi} v_{\pi}(s)$$

$$q_{*}(s, a) \doteq \max_{\pi} q_{\pi}(s, a)$$

In the most simple cases where the value and action-value functions are specified, dynamic programming can be used to derive the optimal policy  $\pi$ .

In the algorithmic trading problem, the value function cannot be ascertained easily in this

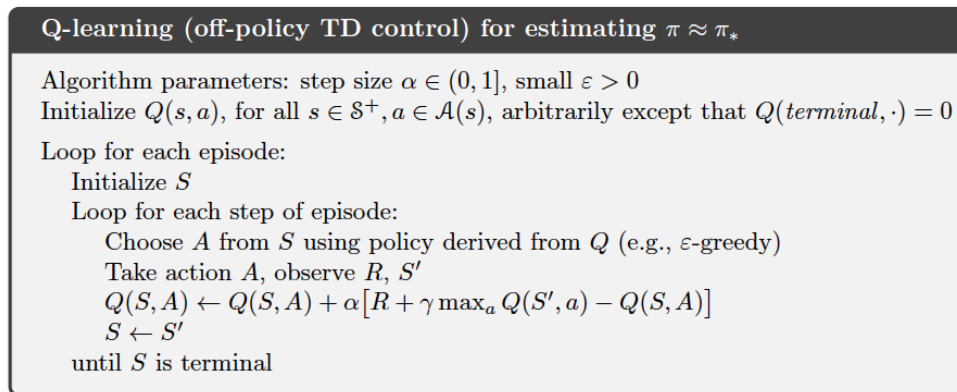
way. To deal with these situations and arbitrarily large state space, approximate solution methods must be used. This is known as a partially observable Markov decision process as the state is only observed indirectly and cannot be fully known (we cannot know the trading behaviour of other agents for example) and we do not have access to the transition probabilities between states.

Q-Learning is a technique whereby the value functions are repeatedly estimated based on the rewards of our actions and assumes no prior model specification.

### 3. Q-Learning

Q-learning attempts to estimate  $q_*$  (optimal action-value function) without any regard for the policy followed. From a high-level perspective, the Q-Learning algorithm proceeds by randomly initialising  $Q$ , perform actions, measure reward and update  $Q$  accordingly (see Figure 2). The final output after a training period should be a stable approximation of the  $q_*$ .

Figure 2: Q-Learning Algorithm



Source: Sutton et al. (1998)

Notice the Bellman equation appearing in the update phase of the algorithm.

### 4. Deep Q Learning

An extension of this idea (which we aim to employ) is to use neural networks to approximate the Q-function ( $Q(s, a)$ ). In situations where the state space is enumerable, we

can produce a Q table which specifies the action-value function for each possible state. However, if the state space is intractable or very large, this is generally not possible or computationally intensive. A neural network is an appropriate tool for our use case due to the infinite nature of the state space. Estimating a table corresponding to every possible state would be excessive in regards to memory requirements.

To train the neural network on the state space, we must define a loss function. The Bellman Equation defines the optimal result; thus we can use this to calculate our loss as follows:

$$\hat{Q}(s, a) = R(s, a) + \gamma \max_{a' \in A} Q(s, a)$$

$$\text{Loss} = \|Q - \hat{Q}\|_2$$

In general, a partition of the data is used to train the neural network and approximate the  $q_*$  function, and this is then used as our action-value function for deciding the optimal policy.

A neural network with two hidden layers is implemented using Tensorflow<sup>1</sup> in our case.

#### 4.1. Problem Formulation - Algorithmic Trading

Now we must formulate the trading problem as a Markov Decision Process and define the states, actions and rewards (Xiong et al., 2018).

- State:  $\mathcal{S} = \{prices, holdings, balance\}$  where *prices* refers to the current prices of all the stocks in our portfolio, *holdings* the quantity of each stock held and *balance* as the total portfolio value. In our simple variant, we are only trading one stock but include a history of prices also (default 200)
- Actions:  $\mathcal{A} = \{buy, sell, hold\}$  For simplicity and computational tractability, we restrict our action space to buy 1 stock, sell 1 stock or hold.

---

<sup>1</sup>TensorFlow is a free and open-source software library for dataflow and differentiable programming across a range of tasks.

- Rewards:  $R_t \in \mathcal{R}$  can be defined as the change in the portfolio value due to an action  $A_t$ . Our Reward was defined as  $R_t = balance_t - balance_{t-1}$
- Policy:  $\pi$  which governs the trading strategy at state  $S$ . We converge on the optimal policy by approximating the  $q_*$  function.
- Action-value function:  $q_\pi(s, a)$  as defined above. The expected reward we obtain by following policy  $\pi$ , choosing action  $A$  while in state  $S$ . The action-value function is approximated by a neural network.

No transaction cost is considered in this study, and the algorithmic trader can only trade with a single stock, Apple (AAPL) or Wawel (WWL). Furthermore, a negative portfolio is not permissible ( $balance_t \geq 0$  for all  $t$ ) and the trader can only sell owned stock ( $holdings_t \geq 1$ ). In our primary model, the initial  $balance_0$  is set to \$1000 and  $price_0$  to the price of the stock at our chosen start time  $t_0$

## 5. Data, Settings and Results

The proposed deep Q-learning system is evaluated on two stocks: Apple (AAPL) and Wawel (WWL). We first describe the data and configurations of the experiment and then present the performance results.

### 5.1. Data and Settings

The databases used in the experiment involve nine years of daily data on Apple and Wawel stocks<sup>2</sup>, ranging from 2010-01-01 through 2019-03-01. The databases are divided into training data sets (2010-01-01 — 2012-04-13) and test sets (2012-04-16 — 2019-02-28). Only the daily closing price is used in this study, though other features can be easily incorporated into the model.

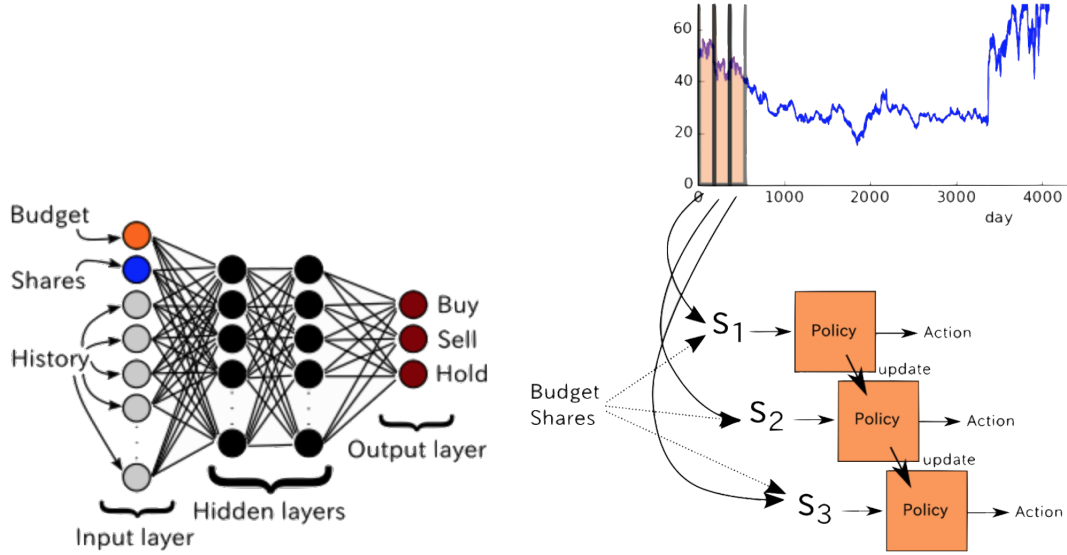
We compare the deep Q-learning system with two benchmark strategies: a random decisions policy (RAND) and a buy-and-hold (BH) strategy. For deep Q-learning, the training datasets are used to initialise the deep Q-network, and then the system runs in an online

---

<sup>2</sup>All data was downloaded using Yahoo finance

fashion where trading decision making and model adaptation are conducted simultaneously. For the Q-learning part, the discount factor is set to  $\gamma = 0.5$ . For the deep-learning part, the architecture of the deep Q-network involves four layers in total (two are hidden). The input units (features) are composed by the current budget ( $balance_t$ ), the current number of stocks held ( $holdings_t$ ) and 200 days of stock history while the output units correspond to the three actions in trading (see Figure 3). The learning rate for Q-network is 0.01, and the training stops after 5 epochs. Furthermore, we introduce a new hyper-parameter  $\epsilon$  (set at  $\epsilon = 0.9$ ) to keep our solution from getting “stuck” when applying the same action over and over. Thus the algorithm exploits the best option with probability  $\epsilon$  and explores a random option with probability  $1 - \epsilon$ .

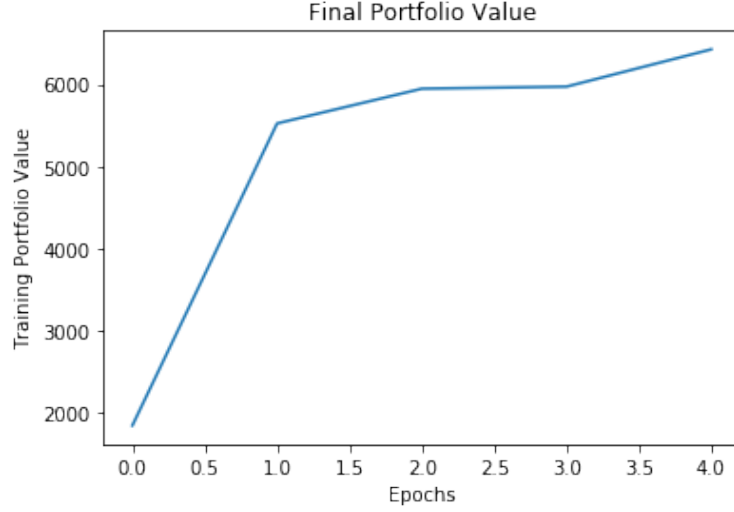
Figure 3: The architecture of the deep Q-network and the rolling window scheme



As mentioned above, to ensure that the Q function is being learned correctly and trending towards an optimal policy, we train the neural network over five epochs. From the plot below (Figure 4), the final portfolio value achieved continues to increase, signalling convergence towards an optimal policy. The hyper-parameters were tuned to derive maximum gain.



Figure 4: Illustrative convergence of the deep Q-network



## 5.2. Results

The results of the three trading approaches (BH, RAND, and Deep Q-learning) on the AAPL test set are presented in Figure 5, and the results on the WWL test set are shown in Figure 6. In each figure, the green line illustrates the Deep Q-learning decision policy while the red and blue lines depict the BH and RAND strategies respectively. It can be seen that on both of the two test sets, the deep Q-learning system accumulates more value than the other two systems. A more detailed numerical comparison is shown in Table 1 where we report two widely used measures for stock trading: accumulated return and the Sharpe ratio. The Sharpe ratio is defined as:

$$\text{Sharpe Ratio} = \frac{R_p - R_f}{\sigma_p}$$

where  $R_f$  is the risk-free rate (assumed to be 8%),  $R_p$  the return of the portfolio and  $\sigma_p$  is the standard deviation of the portfolio. A larger Sharpe Ratio indicates that the portfolio achieves a better return for its volatility. In a practical setting, this would be an attractive feature for a potential investor. Thus, it can be seen that our deep Q-learning system outperforms the other two methods in terms of total return and the Sharpe Ratio.

Therefore, the results show our deep Q-learning model performs well on both the two stocks (AAPL being a strong stock with a stable upward trend and WWL being a more volatile stock with no apparent trend). This advantage of deep Q-learning can be at-

Table 1: Comparison of Trading Performance

	AAPL			WLL		
	DQL	BH	RAND	DQL	BH	RAND
Accumulated Return(%)	1.67	1.33	0.59	0.64	0.41	0.18
Sharp Ratio	0.29	0.22	0.26	0.25	0.20	0.09

tributed to two advantages. One of them is the ability to detect the status of the market from the raw and noisy data, and the other is the online nature that adapts itself to new market status quickly. More interesting observations can be found in the portfolio action plots (see Figure 7 and Figure 8). From the positions held by the Q-learning system, it seems that it has learned how to take different actions in different market situations which can largely be attributed to the power of the deep Q-network in discovering the status of the market from the noisy historical price signals.

Figure 5: The performance of various trading strategies on AAPL

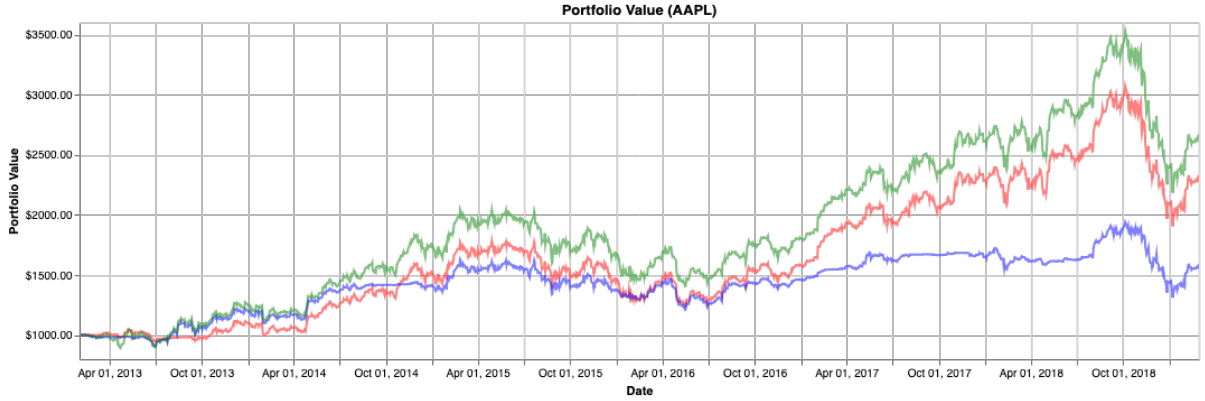


Figure 6: The performance of various trading strategies on WLL

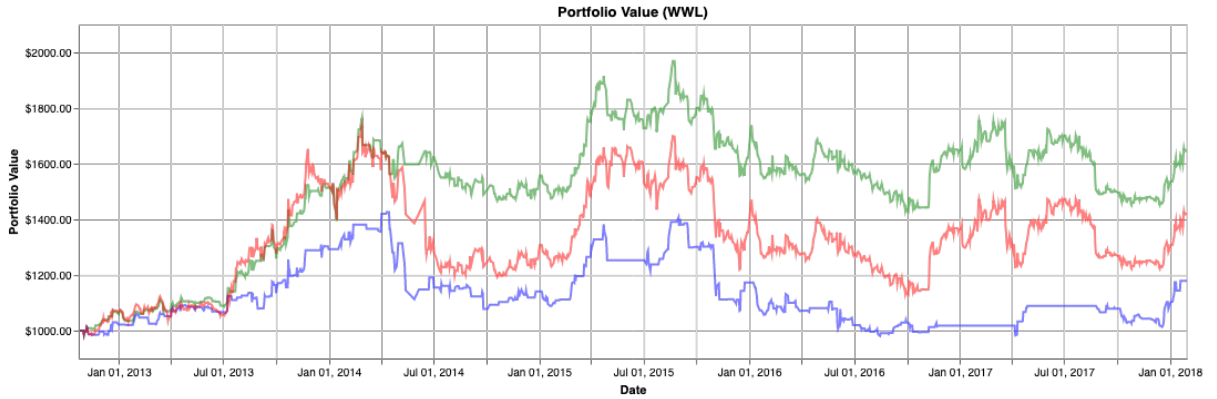


Figure 7: The positions held by various trading strategies on AAPL

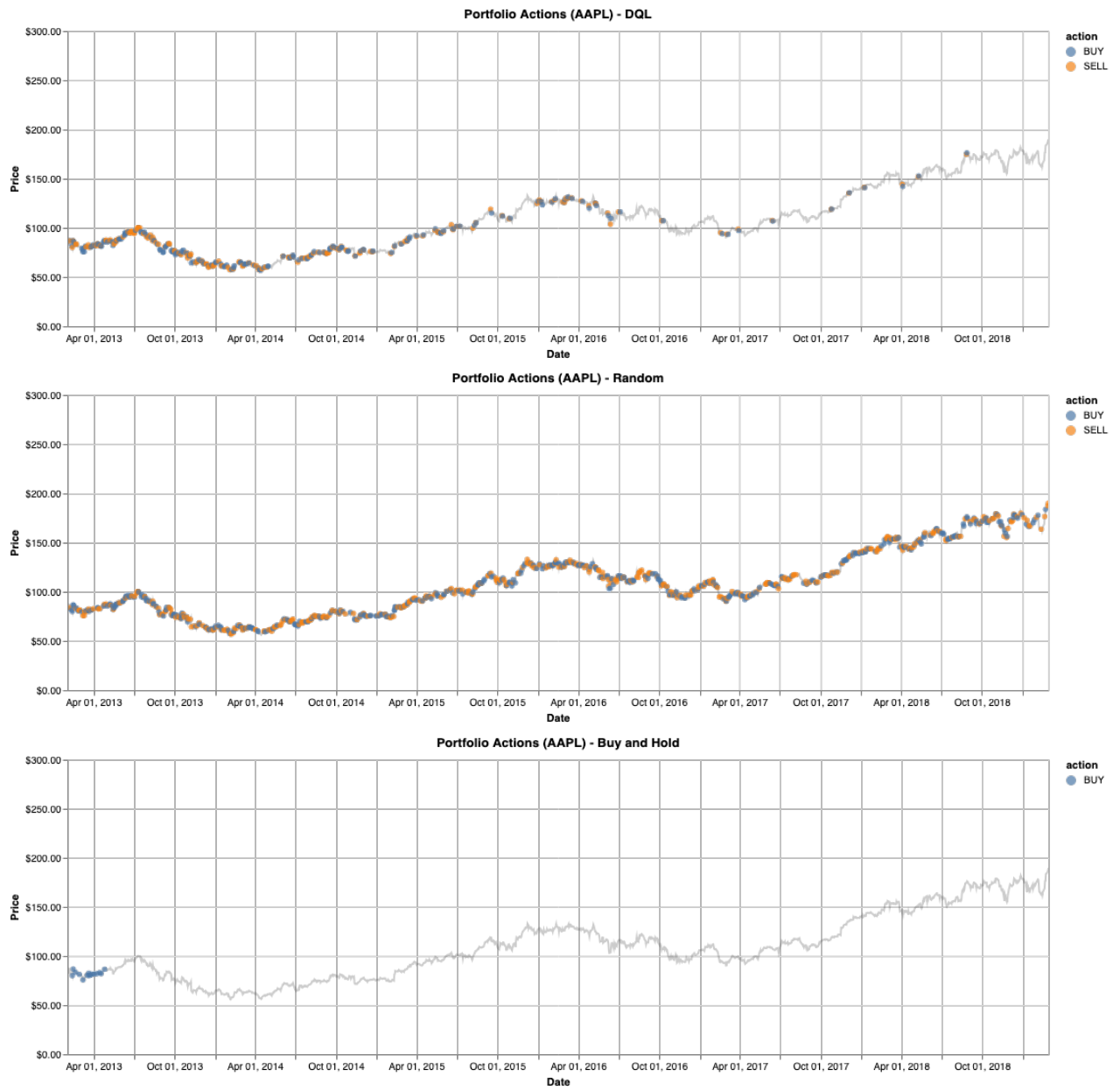
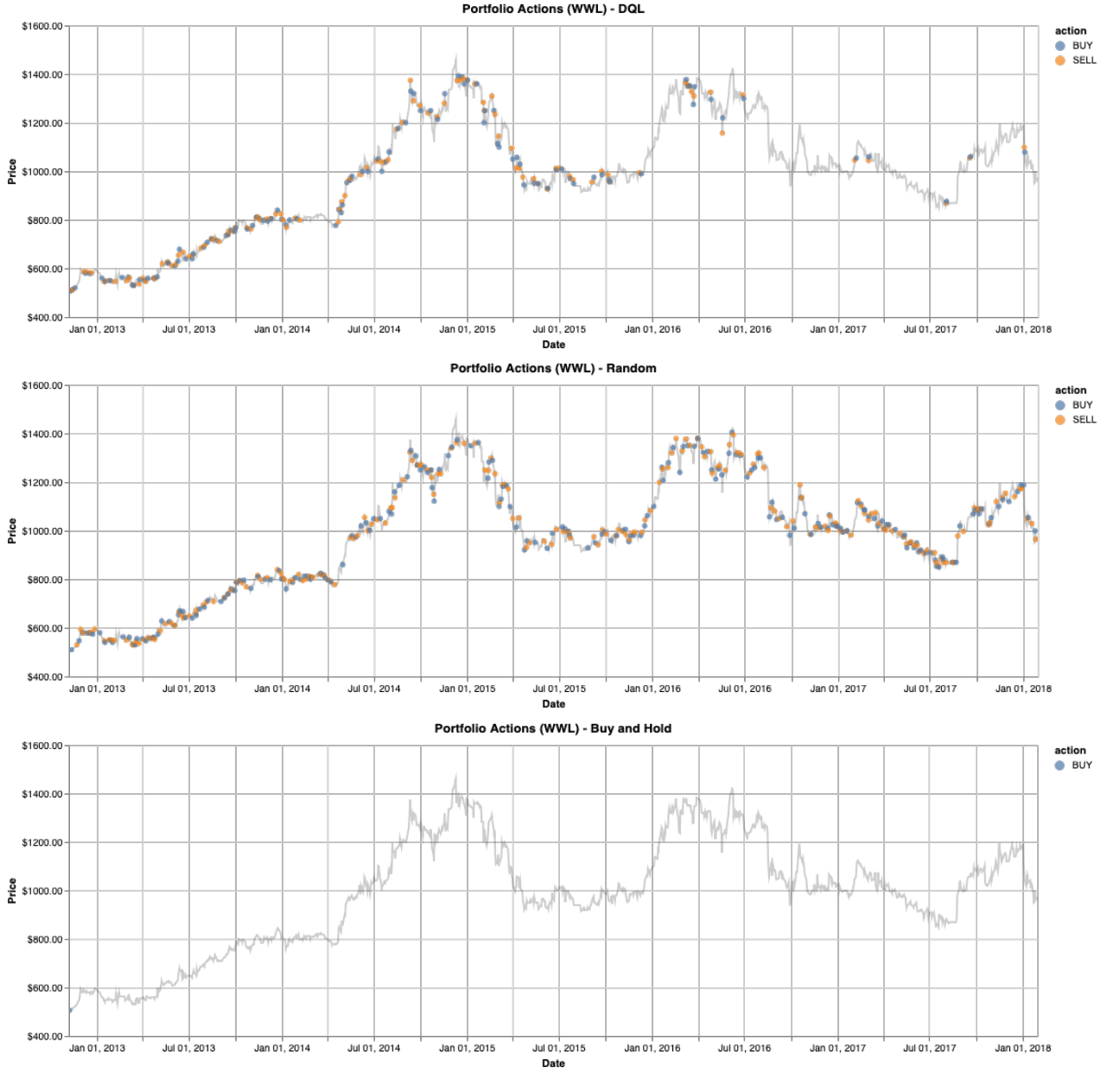


Figure 8: The position held by various trading strategies on WWL



## 6. Conclusion

The investigation aimed to implement reinforcement learning in the context of stock trading. A deep Q-Learning approach was used to approximate the Q-function (action-value function) and develop a trading policy. Through our experiments with AAPL and WLL stocks, we have managed to show that with minimal training and hyper-parameter tuning, our agent can outperform random actions and a Buy and Hold strategy. We have also demonstrated the agent's ability to manage the volatility of the portfolio effectively.

There are many possible future developments that could be explored such as using multiple stocks, modifying the reward function and doing further hyper-parameter tuning. Furthermore, the initial budget appeared to be a strong determinant of the behaviour of our agent. A more thorough investigation into this and its implications for algorithmic trading would be an interesting extension.

## 7. References

- Bertoluzzo, F. and Corazza, M. (2012), ‘Testing different reinforcement learning configurations for financial trading: Introduction and applications’, *Procedia Economics and Finance* **3**, 68–77.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D. and Riedmiller, M. (2013), ‘Playing atari with deep reinforcement learning’, *arXiv preprint arXiv:1312.5602* .
- Neuneier, R. (1996), Optimal asset allocation using adaptive dynamic programming, in ‘Advances in Neural Information Processing Systems’, pp. 952–958.
- Sutton, R. S., Barto, A. G. et al. (1998), *Introduction to reinforcement learning*, Vol. 135, MIT press Cambridge.
- Xiong, Z., Liu, X.-Y., Zhong, S., Walid, A. et al. (2018), ‘Practical deep reinforcement learning approach for stock trading’, *arXiv preprint arXiv:1811.07522* .