

---

# MATHEMATICS OF REINFORCEMENT LEARNING WITH APPLICATIONS TO QUANTITATIVE FINANCE

---

A PREPRINT

**Yuanhao JIANG\***  
School of Mathematics  
The University of Edinburgh  
Edinburgh, UK  
s2132254@ed.ac.uk

August 29, 2022

## ABSTRACT

This paper provides a new way of pricing in the car insurance industry, which is, using reinforcement learning algorithms to train pricing agents to promote proper price, such that our customer is most likely to purchase our product, and we also get the highest profit and a better portfolio. We first introduce general reinforcement learning knowledge, including how we model the interaction between agent and environment using Markov decision process, and the mathematics behind those methods to solve the reinforcement learning problem, especially the Actor-Critic algorithm. Then we start building the environment, and use artificial neural networks as our pricing models, and writing algorithms that are specific to our financial problem. In the end, we will compare the Actor-Critic algorithms with other popular algorithms including REINFORCE and PPO in terms of their rate of convergence to a stable reward.

**Keywords** Reinforcement Learning · Artificial Neural Networks · Quantitative Finance

## 1 Introduction

Modern way of calculating car insurance premium is by estimating how much risk the customer poses to the insurance provider, and more risk results in higher premium. The estimation of the risk for a customer is usually done by comparing various risk factors with the statistics of the database collected so far. Those factors might include customer's age, address, occupation, car's mileage, brand and model, and so on. For example, younger drivers pose higher risks since they are more likely to have an accident due to the lack of driving experience compare to elder drivers.

In this paper, instead of using the above method to promote car insurance, we use RL methods to train agents that will interact with the market and response with price. More specifically, we consider the insurance promoting scenario as markov decision process, with the current state being the current customer features and current portfolio. At each time step, the agent promote a price according to the current customer features and the current portfolio, after that the customer resopnse to the promotion. The reward is then given based on the customer response and the portfolio. Then the portfolio updates and a new customer arrives. This set up will help us to solve the pblem in the context of reinforcement learning, further details are discussed in the remaining sections.

## 2 Reinforcement Learning

The reinforcement learning method we want to use is based on the Markov decision process (MDP), as illustrated in Figure 1. However, in RL, the agent is not told which actions to take, but instead must discover which actions yield the

---

\*Current student at the University of Edinburgh.

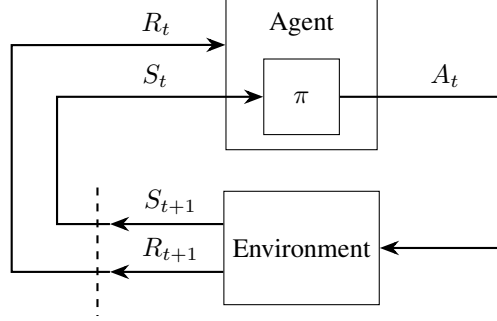


Figure 1: MDP Illustration [2]

most reward by trying them. Our goal at each time step is to find a policy  $\pi$  to maximize the expected reward along the process afterwards, which is given by the value function

$$V(s) = \mathbb{E} \left[ \sum_{i=0}^{\infty} \gamma^i R_{t+i+1} \mid S_t = s \right]$$

To achieve this, we need to evaluate the current policy, and improve it accordingly.

## 2.1 Policy Evaluation

We want to evaluate the value function under a given policy  $\pi$ . According to Bellman equation,

$$v_{\pi}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

We can use it to iteratively compute  $v_{\pi}$ . The iteratively update is

$$v_{k+1}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s]$$

which is proven that it will converge to the true value given any initial value function.

## 2.2 Policy Improvement

Now we have the evaluated value function for the policy, we want to improve it accordingly. To do this, we introduce the q-function under policy  $\pi$

$$q_{\pi}(s, a) = \mathbb{E} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s, A_t = a]$$

The policy improvement theorem says that, let  $\pi$  and  $\pi'$  be any pair of policy, for all  $s$  we have

$$q_{\pi}(s, \pi'(s)) \geq v_{\pi}(s) \Rightarrow v_{\pi'}(s) \geq v_{\pi}(s)$$

Now we can improve the policy by

$$\pi'(s) = \operatorname{argmax}_a q_{\pi}(s, a)$$

for all  $s$ .

## 2.3 Policy Iteration

By iteratively applying policy evaluation and policy improvement, we can obtain a sequence of monotonically improving policies and value functions [2] :

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi_* \xrightarrow{E} v_*$$

where  $\xrightarrow{E}$  denotes a policy evaluation and  $\xrightarrow{I}$  denotes a policy improvement. This way of finding an optimal policy is called policy iteration.

## 2.4 Value Iteration

One drawback to policy iteration is that each of its iterations involves policy evaluation, which itself is a protracted iterative computation requiring multiple sweeps through the state set.

We now truncate the policy evaluation such that it is stopped after just one sweep [2] (one update of each state):

$$\begin{aligned} v_{k+1}(s) &= \max_a q(s, a) \\ &= \max_a \mathbb{E} [R_{t+1} + \gamma v_k(S_{t+1}) \mid S_t = s, A_t = a] \end{aligned}$$

We keep iterate until  $v$  converges (to  $v_*$ ), then we record the last policy  $\pi(s) = \operatorname{argmax}_a q(s, a)$ .

## 2.5 Generalized Policy Iteration

As long as both processes, policy evaluation and policy improvement, continue to update all states, the ultimate result is typically the same — converge to the optimal value function and an optimal policy [2].

Generalized policy iteration (GPI) is used to refer to the general idea of letting policy evaluation and policy improvement interact, independent of the granularity of the two processes [2] (instead of letting each complete before the other begins).

## 2.6 Temporal-Difference Learning

When we don't have the complete model of the environment, we need to estimate the value function from only experience. There are two methods to do so: Monte Carlo Methods and Temporal-Difference. We will use Temporal-Difference method in this project.

To estimate  $v_\pi(S_t)$  using only the experience, we take the average on all experienced trajectory' values. Temporal-Difference method uses  $G_i = [R_{t+1}]_i + \gamma[v_\pi(S_{t+1})]_i$  to estimate the value of the  $i$ th trajectory starting at  $S_t$ , with  $v_\pi$  being the estimated value function at the  $i$ th time  $S_t$  is visited, which can be initialized to anything at start. Now we have the incremental formula to update our value function:

$$\begin{aligned} [v_\pi(S_t)]_{k+1} &= \frac{1}{k} \sum_{i=1}^k G_i \\ &= \frac{1}{k} \left( G_k + \sum_{i=1}^{k-1} G_i \right) \\ &= \frac{1}{k} (G_k + (k-1)[v_\pi(S_t)]_k) \\ &= [v_\pi(S_t)]_k + \frac{1}{k} (G_k - [v_\pi(S_t)]_k) \\ &= [v_\pi(S_t)]_k + \frac{1}{k} ([R_{t+1}]_k + \gamma[v_\pi(S_{t+1})]_k - [v_\pi(S_t)]_k) \end{aligned}$$

To simplify:

$$v_\pi(S_t) \leftarrow v_\pi(S_t) + \alpha [R_{t+1} + \gamma v_\pi(S_{t+1}) - v_\pi(S_t)]$$

with  $\alpha$  being the stepsize. And this also applies to the  $q$  function:

$$\begin{aligned} q_\pi(S_t, A_t) &\leftarrow q_\pi(S_t, A_t) \\ &\quad + \alpha [R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) - q_\pi(S_t, A_t)] \end{aligned}$$

## 2.7 Policy Gradient Methods

When the state space is arbitrarily large, we cannot find an optimal policy or the optimal value function given the limited resources and time. Instead of keeping track of the value and the action(s) to select for each state (tabular methods), we can use parameterized functions, for example, artificial neural networks (ANNs), to approximate the value function and the policy function.

Previously methods learned the values of actions and then select actions based on their estimated action values (q function). Now we consider methods that instead learn a parameterized policy that select actions without consulting a value function. A value function may still be used to learn the policy parameter, but is not required for action selection [2].

We use  $\theta$  and  $w$  to denote the policy's parameter vector and value function's weight vector, respectively. Then we have our policy  $\pi(a | s, \theta) = \Pr\{A_t = a | S_t = s, \theta_t = \theta\}$  and value function  $v_\pi(s, w)$ . We also define the scalar performance measure (objective) to be

$$J(\theta) = \begin{cases} v_{\pi_\theta}(s_0); & \text{episodic case} \\ \lim_{t \rightarrow \infty} \mathbb{E}[R_t | S_0, A_{0:t-1} \sim \pi]; & \text{continuing case} \end{cases}$$

Policy gradient methods seek to maximize this performance measure, so their updates approximate gradient ascent in  $J$ :

$$\theta_{t+1} = \theta_t + \alpha \widehat{\nabla J(\theta_t)}$$

where  $\widehat{\nabla J(\theta_t)}$  is a stochastic estimate whose expectation approximates the gradient of the performance measure w.r.t  $\theta_t$ .

### 2.7.1 The Policy Gradient Theorem

The policy gradient theorem says

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a | s, \theta)$$

where  $\mu$  is the stationary distribution of the state. Thus we can derive

$$\begin{aligned} \nabla J(\theta) &\propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi(a | s, \theta) \\ &= \mathbb{E}_\pi \left[ \sum_a q_\pi(S_t, a) \nabla \pi(a | S_t, \theta) \right] \\ &= \mathbb{E}_\pi \left[ \sum_a \pi(a | S_t, \theta) q_\pi(S_t, a) \frac{\nabla \pi(a | S_t, \theta)}{\pi(a | S_t, \theta)} \right] \\ &= \mathbb{E}_\pi \left[ q_\pi(S_t, A_t) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \right] \end{aligned}$$

The policy gradient theorem can be generalized to include a comparison of the action value to an arbitrary baseline  $b(s)$ :

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a (q_\pi(s, a) - b(s)) \nabla \pi(a | s, \theta)$$

with the same derivation we have

$$\nabla J(\theta) \propto \mathbb{E}_\pi \left[ (q_\pi(S_t, A_t) - b(S_t)) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \right]$$

And the discounted version of above is:

$$\nabla J(\theta) \propto \mathbb{E}_\pi \left[ \gamma^t (q_\pi(S_t, A_t) - b(S_t)) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \right]$$

One natural choice for baseline is an estimate of the state value,  $\hat{v}(S_t, w)$ , where  $w$  is the weight vector to be learned.

### 2.7.2 Actor-Critic Methods

Recall the TD method, we make use of  $q_\pi(s, a) = \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a]$ , by taking one sample estimate, and choosing  $\hat{v}(S_t, w)$  as baseline, we have

$$\nabla J(\theta) \propto \mathbb{E}_\pi \left[ (R_{t+1} + \gamma \hat{v}(S_{t+1}, w) - \hat{v}(S_t, w)) \frac{\nabla \pi(A_t | S_t, \theta)}{\pi(A_t | S_t, \theta)} \right]$$

Again, by taking one sample estimate we have the update formula for  $\theta$ :

$$\begin{aligned}\theta_{t+1} &= \theta_t + \alpha_\theta \left( R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}) \right) \frac{\nabla \pi(A_t | S_t, \theta_t)}{\pi(A_t | S_t, \theta_t)} \\ &= \theta_t + \alpha_\theta \delta_t \nabla \ln \pi(A_t | S_t, \theta_t)\end{aligned}$$

As for the weight vector  $\mathbf{w}$  for the estimated value function, at each time step  $t$ , we want to update it by minimizing the squared error,  $[v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)]^2$ :

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \alpha'_w \nabla [v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t)]^2 \\ &= \mathbf{w}_t - \alpha_w \left( v_\pi(S_t) - \hat{v}(S_t, \mathbf{w}_t) \right) \nabla \hat{v}(S_t, \mathbf{w}_t)\end{aligned}$$

where  $v_\pi$  is the true value function. Again, since  $v_\pi(s) = \mathbb{E}[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s]$ , we take one sample estimate, which gives us the update formula for  $\mathbf{w}$ :

$$\begin{aligned}\mathbf{w}_{t+1} &= \mathbf{w}_t - \alpha_w \left( R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}_t) \right) \nabla \hat{v}(S_t, \mathbf{w}_t) \\ &= \mathbf{w}_t - \alpha_w \delta_t \nabla \hat{v}(S_t, \mathbf{w}_t)\end{aligned}$$

The actor-critic method is based on the above updates formulas, details are shown in Algorithm 1.

---

**Algorithm 1** Actor-Critic

---

- 1: Input: a differentiable policy parameterization  $\pi(a|s, \theta)$
  - 2: Input: a differentiable state-value function parameterization  $\hat{v}(s, \mathbf{w})$
  - 3: Step sizes  $\alpha_\theta > 0, \alpha_w > 0$
  - 4:  $\theta \leftarrow \mathbf{0}, \mathbf{w} \leftarrow \mathbf{0}$
  - 5: **loop** (for each episode)
  - 6:   Initialize  $S$  (first state of the episode)
  - 7:   **loop** (for each time step  $i$ )
  - 8:      $A \sim \pi(\cdot | S, \theta)$
  - 9:     Take action  $A$ , observe  $S', R$
  - 10:     $\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$
  - 11:     $\mathbf{w} \leftarrow \mathbf{w} + \alpha_w \delta \nabla \hat{v}(S, \mathbf{w})$
  - 12:     $\theta \leftarrow \theta + \alpha_\theta \gamma^i \delta \nabla \ln \pi(A | S, \theta)$
  - 13:     $S \leftarrow S'$
  - 14:   **end loop**
  - 15: **end loop**
- 

### 3 Build the Environment

The environment does the following things:

1. record the current state that can be queried anytime
2. when an action, i.e. a price, is given according to the current state, it outputs the reward, and goes to next state

At each time step, the state consists of a randomly generated customer's features and the current portfolio. After the agent performs an action, i.e. promotes a price to the customer, the reward is then computed based on the profit on this customer and the portfolio up to now.

#### 3.1 Customer Features

A customer's features is represented by an array (or vector), denoted by  $x$  with 16 entries with each generated randomly from a distribution, including normal distribution, binomial distribution, and so on. For example, a feature vector  $x$  is described as Table 1

Table 1: Customer Feature Vector Description

gender	age	car cost	miles	brand	random feature 0	random feature 1	...
1.0	42.723	94032.096	30086.311	29.264	8.0	68.740	...

After a price, denoted by  $c$ , is promoted to a customer, the customer will give a response, denoted by  $r$ , which is 0 if he does not buy the product, and 1 otherwise. And the profit we have for that customer, denoted by  $pf$ , can be computed by

$$pf = r * (c - g(x))$$

where  $g(x)$  is the cost on that customer.

### 3.2 Response

The response is generated by a generalized linear model, whose input is a vector of a customer features vector  $x$  appended with a price  $c$ , i.e.  $r = glm(x[0], \dots, x[15], c)$ .

The GLM used in the project is fitted using a randomly generated data set with each observation is a customer's features followed by a price, then followed by a response, as illustrated by Table 2.

Table 2: Customer Response Examples

	customer feature 0	...	customer feature 15	price	response
1	1.0	...	2.0	1109.504	1.0
2	0.0	...	1.0	665.8	0.0
...	...	...	...	...	...

And the link function for the GLM is chosen to be the logistic function.

### 3.3 State

The state at start of each time step is the feature vector of the new customer to promote (might be modified a bit), appended by the portfolio, which is also a data vector consists of the following entries:

- average profit over all customers so far, expect for the new customer at current time step, denoted by  $avg\_pf$  (initialized to 0), which is computed by

$$avg\_pf_t = \frac{1}{t} \sum_{i=0}^{t-1} pf_i, t > 0$$

- the portion of the buyers over all customers, expect for the new one, denoted by  $p$  (initialized to 0), which is computed by

$$p_t = \frac{1}{t} \sum_{i=0}^{t-1} r_i, t > 0$$

- the variance of the portion of buyers in their categories, denoted by  $var$ : we divide customers (expect for the new one at current time step) into 4 categories, according to their features, and for each category, we compute the portion of the customers who buy our product, denoted by  $pt_i$  with  $i = 1, 2, 3, 4$ , then we compute the variance of the 4 portions:

$$var = Var(pt) = \frac{\sum_{i=1}^4 (pt_i - \bar{pt})^2}{3}$$

Later we use this to compute reward to make sure the agent will promote to all types of customers instead of only focusing on one type. Note that  $var \leq 1$ .

Thus, the state  $s$ , or  $s_t$ , end up with the form

$$s_t = (x'_t[0], \dots, x'_t[-1], avg\_pf_t, p_t, var_t)$$

where  $x'_t$  is the modified  $x_t$ , including preprocess of converting some categorical data to meaningful data to the model.

### 3.4 Reward

The reward at each time step  $t$ , after the price  $c_t$  is promoted to the customer  $x_t$ , denoted by  $R_t$ , is calculated by

$$R_t = \text{pf}_t (1 - h(\text{var}_t))$$

where  $h$  is a function that indicates how much we care about the variance, i.e. how important it is to try to promote to customers of all categories. In this project I choose  $h$  to be such that  $h(\text{var}) = \sqrt{\text{var}}/2$ .

## 4 Build the model

With the environment well built, we now build the model for our policy and the value function. In this project we use artificial neuro networks (ANNs) as the parameterized policy and value function, each with three linear layers, and with the ReLU function as the activation function. And the input of both function is the state, as illustrated in Figure 2.

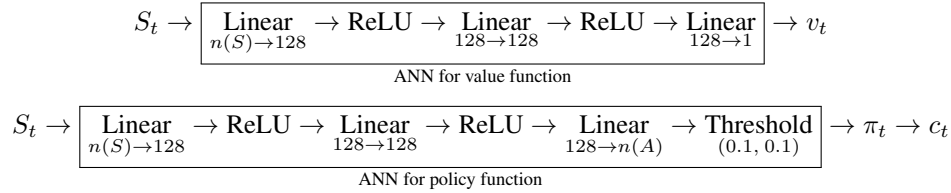


Figure 2: ANN model for policy and value function

where  $n(S)$  denotes the vector length of the state and  $n(A)$  denotes the number of actions.

We generate the policy distribution by scoring each action. E.g. denoted the policy function output by  $y$ , then the probability of choosing action  $i$  ( $i = 0, 1, \dots, n(A) - 1$ ) is computed by

$$P(i) = \frac{y_i}{\sum_{j=0}^{n(A)-1} y_j}$$

Note that we add a Threshold layer at the end of policy function, which is defined by

$$\text{Threshold}(0.1, 0.1) = \max(0.1, x)$$

which makes sure that

1. the probabilities are valid (no negative values, greater score results in greater probability);
2. all actions have the probability to be chosen at early stage, which is beneficial for training since it increase the exploration for different actions.

## 5 RL with Actor-Critic

In this project, we consider the episodic learning procedure, where each episode has finite time steps,  $T$ , to be specified as a constant, and we will also reset the environment at the start of each episode.

### 5.1 The Algorithm

Algorithm 2 uses RL with actor-critic method, and with the help of what we have built in previous sections, to solve the problem. And Figure 3 illustrates the algorithm.

### 5.2 Training Result

Figure 4 shows the training result of Algorithm 1 (300 moving average reward vs iteration), with 600 iterations (5 episodes per iteration),  $T = 300$ ,  $\gamma = 0.99$ ,  $\alpha_\theta = \alpha_w = 3 \times 10^{-4}$ . The graph shows that the moving average reward converges within only 50 iterations, that is, 250 episode, 75000 parameter updates for each network, and stays quite stable afterwards.

**Algorithm 2** Actor-Critic for financial problem

---

```

1: Initialize ANN policy parameterization  $\pi(c|s, \theta)$ , with any  $\theta$ 
2: Initialize ANN state-value function parameterization  $\hat{v}(s, w)$  with any  $w$ .
3: Step sizes  $\alpha_\theta > 0, \alpha_w > 0$ 
4: loop (for each episode)
5:   Reset the environment
6:   Generate an  $x$  (first customer of the episode)
7:   Initialize  $S$  (first state of the episode)
8:   for  $t = 0, 1, \dots, T$  do
9:      $c \sim \pi(\cdot|S, \theta)$ 
10:    Take action  $c$ , observe  $S', R$  from environment
11:     $\delta \leftarrow R + \gamma \hat{v}(S', w) - \hat{v}(S, w)$ 
12:     $w \leftarrow w + \alpha_w \delta \nabla \hat{v}(S, w)$ 
13:     $\theta \leftarrow \theta + \alpha_\theta \gamma^t \delta \nabla \ln \pi(A|S, \theta)$ 
14:     $S \leftarrow S'$ 
15:   end for
16: end loop

```

---

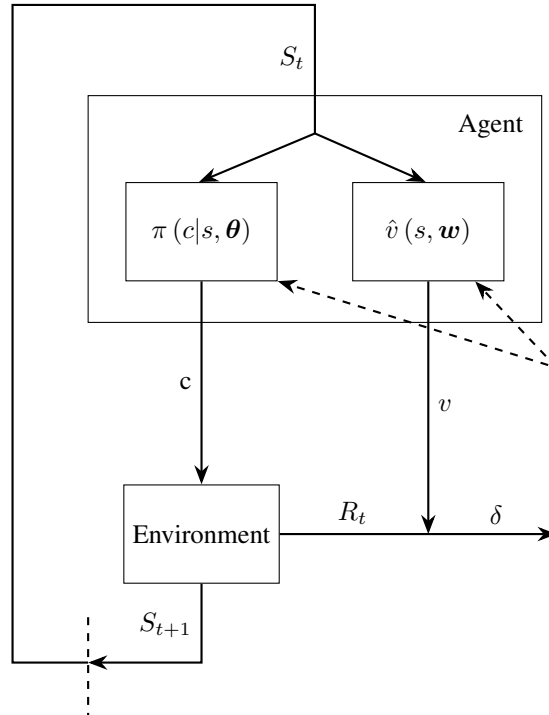


Figure 3: Actor-Critic for Finance Algorithm Illustration



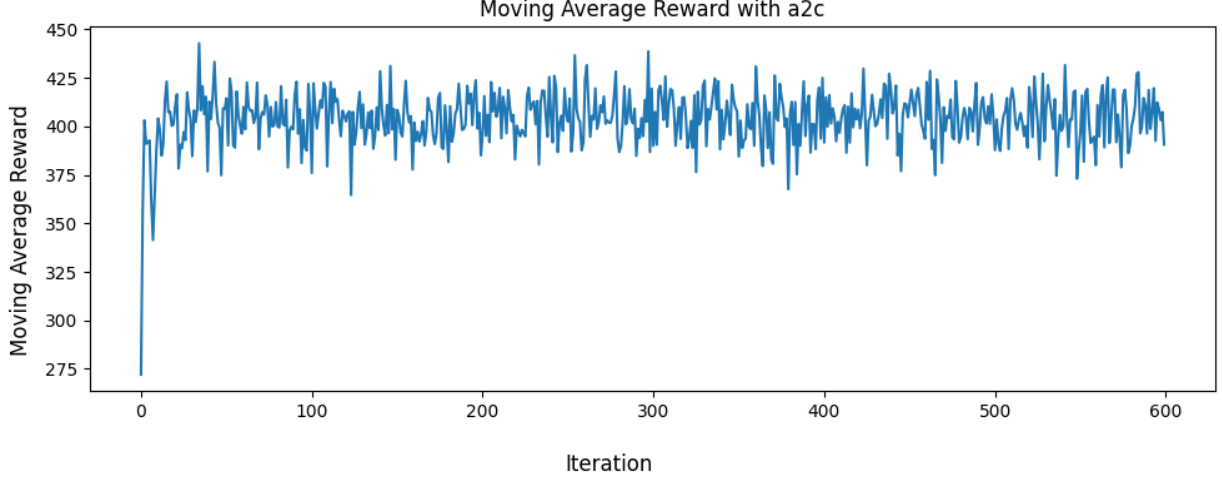


Figure 4: Training Result for Algorithm a2c

## 6 Actor-Critic, Reinforce, and PPO

### 6.1 Algorithms

To further show how Actor-Critic algorithm performs, we compare it to another two popular algorithms: REINFORCE (Algorithm 3) and PPO (Algorithm 4). The policy function of the two algorithms is implemented the same way as we did in Actor-Critic algorithm, and the same is true for the value function used in PPO method.

---

#### Algorithm 3 Reinforce for financial problem [2]

---

- 1: Initialize ANN policy parameterization  $\pi(c|s, \theta)$ , with any  $\theta$
  - 2: Step sizes  $\alpha_\theta > 0$
  - 3: **loop** (for each episode)
  - 4:   Reset the environment
  - 5:   Generate an episode  $S_0, C_0, R_1, \dots, S_{T-1}, C_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$
  - 6:   **for**  $t = 0, 1, \dots, T$  **do**
  - 7:      $G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$
  - 8:      $\theta \leftarrow \theta + \alpha_\theta \gamma^t G \nabla \ln \pi(C_t|S_t, \theta)$
  - 9:   **end for**
  - 10: **end loop**
- 

### 6.2 Training result

Figure 5 shows the training result of these three algorithms (300 moving average reward vs iteration), with 600 iterations (5 episodes per iteration),  $T = 300$ ,  $\gamma = 0.99$ ,  $\alpha_\theta = 3 \times 10^{-4}$ ,  $\alpha_w = 3 \times 10^{-10}$  (if applicable),  $\epsilon = 0.5$  (if applicable). As we can see, the Actor-Critic algorithm converges the fastest, within 50 iterations, followed by REINFORCE algorithm, which converges within 250 iterations, despite a small drop of moving average around iteration 300, which is quickly corrected by the algorithm within 50 iterations. PPO algorithm converges the slowest, with about 300 iterations.

## 7 Conclusion

To summarize, the moving average rewards after convergence of these three algorithms are nearly the same (around 400), which is as expected. The Actor-Critic algorithm performs really well compared to another two algorithms, it has highest converge rate, the reward is also very stable after convergence. In the mean time, the slow convergence for PPO algorithm is expected since it uses a clamp function to prevent model from changing too fast in a single update,

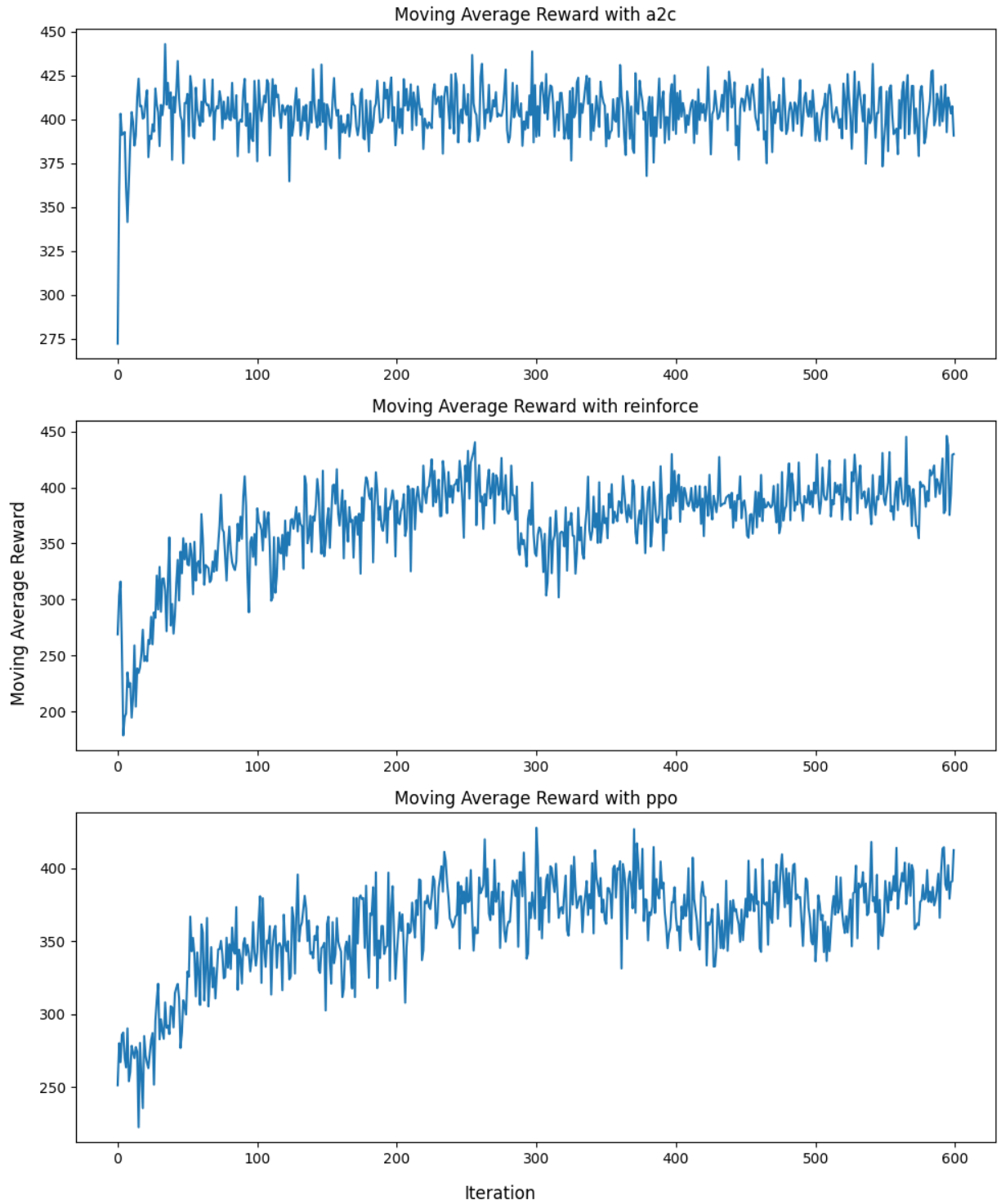


Figure 5: Training Result Comparison for Three Algorithm

**Algorithm 4** PPO for financial problem [1]

---

```

1: Initialize ANN policy parameterization  $\pi(c|s, \theta)$ , with any  $\theta$ 
2: Initialize ANN state-value function parameterization  $\hat{v}(s, w)$  with any  $w$ .
3: Step sizes  $\alpha_\theta > 0, \alpha_w > 0$ 
4: for iteration  $k = 0, 1, 2, \dots$  do
5:   Collect set of trajectories  $\mathcal{D}_k = \{\tau_i\}$ , following  $\pi(\cdot|\cdot, \theta_k)$ , where  $\tau = S_0, C_0, R_1, \dots, S_{T-1}, C_{T-1}, R_T$ 
6:   for  $t = 0, 1, \dots, T$  in each trajectory  $\tau$  in  $\mathcal{D}_k$  do
7:      $G_t \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$ 
8:      $A_t \leftarrow G_t - \hat{v}(S_t)$ 
9:   end for
10:  loop (update multiple times)
11:     $\theta \leftarrow \arg \max_{\theta} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \min \left( \frac{\pi(c_t|S_t, \theta)}{\pi(c_t|S_t, \theta_k)} A_t, \text{clip} \left( \frac{\pi(c_t|S_t, \theta)}{\pi(c_t|S_t, \theta_k)}, 1 - \epsilon, 1 + \epsilon \right) A_t \right)$ 
12:     $w \leftarrow \arg \min_w \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left( \hat{v}(S_t, w) - G_t \right)^2$ 
13:  end loop
14: end for

```

---

although at current stage it seems to be unnecessary for our specific problem setup. Further investigation needs to be done for the unexpected small drop of the reward for the REINFORCE algorithm.

## References

- [1] Joshua Achiam. Spinning Up in Deep Reinforcement Learning. 2018.
- [2] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018.