

Reinforcement Learning: An Introduction

Attempted Solutions

Chapter 4

Scott Brownlie & Rafael Rui

1 Exercise 5.1

Consider the diagrams on the right in Figure 5.1. Why does the estimated value function jump up for the last two rows in the rear? Why does it drop off for the whole last row on the left? Why are the frontmost values higher in the upper diagrams than in the lower?

The estimated value function jumps up for the last two rows in the rear because these rows correspond to the player having 20 or 21, which are the only two values on which the player sticks. The player cannot lose when he sticks on 21 and is unlikely to lose when sticking on 20. For values less than 20 the player always hits, whereas the dealer only hits if he has less than 17, which appears a better strategy, hence why the value function is low for values less than 20.

The value function drops off for the whole last row on the left because this row corresponds to the dealer showing an ace. The ace is the best card to have because it can either be counted as 1 or 11 and thus increases the dealer's chances of winning.

The frontmost values are higher in the upper diagrams because they correspond to the player having a usable ace. In this case the player has less chance of going bust. For example, if he has an ace and a 1, giving a total of 12, and he hits and receives a 10, the ace will then be counted as 1 and the player will still be in the game.

2 Exercise 5.2

Suppose every-visit MC was used instead of first-visit MC on the blackjack task. Would you expect the results to be very different? Why or why not?

The only way that the player can have the same sum more than once in a single episode is if he has a usable ace. For example, if he has an ace and a 5 then his sum will be 16. If he then receives a 6 then, to avoid going bust, the ace will be counted as 1 and his sum will be $1+5+6 = 12$. There is now the possibility that he once again ends up with a sum of 16 (by receiving a 4, for example). However, the same sum with and without a usable ace are classed as different states, hence it is impossible for the player to visit the same state more than once in a single episode and the estimated value functions will be exactly the same.

3 Exercise 5.3

What is the backup diagram for Monte Carlo estimation of q_π ?

4 Exercise 5.4

The pseudocode for Monte Carlo ES is inefficient because, for each state–action pair, it maintains a list of all returns and repeatedly calculates their mean. It would be more efficient to use techniques similar to those explained in Section 2.4 to maintain just the mean and a count (for each state–action pair) and update them incrementally. Describe how the pseudocode would be altered to achieve this.

We would initialise $N(s, a) \leftarrow 0$ for all $s \in \mathcal{S}, a \in \mathcal{A}(s)$ and then, given that S_t, A_t is the first visit to some state-action pair in some episode, we would increment $N(S_t, A_t)$ by 1 and update $Q(S_t, A_t)$ as follows:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{1}{N(S_t, A_t)} [G - Q(S_t, A_t)].$$

5 Exercise 5.5

Consider an MDP with a single nonterminal state and a single action that transitions back to the nonterminal state with probability p and transitions to the terminal state with probability $1 - p$. Let the reward be +1 on all transitions, and let $\gamma = 1$. Suppose you observe one episode that lasts 10 steps, with a return of 10. What are the first-visit and every-visit estimators of the value of the nonterminal state?

The first-visit estimate of the value of the nonterminal state is 10. The every-visit estimate is

$$\frac{1}{10} (10 + 9 + 8 + \dots + 1) = 5.5.$$

6 Exercise 5.6

What is the equation analogous to (5.6) for action values $Q(s, a)$ instead of state values $V(s)$, again given returns generated using b ?

It is

$$Q(s, a) = \frac{\sum_{t \in \mathcal{T}(s, a)} \rho_{t:T(t)-1} G_t}{\sum_{t \in \mathcal{T}(s, a)} \rho_{t:T(t)-1}}$$

where $\mathcal{T}(s, a)$ denotes the set of all time steps in which action a is taken in state s .

7 Exercise 5.7

In learning curves such as those shown in Figure 5.3 error generally decreases with training, as indeed happened for the ordinary importance-sampling method. But for the weighted importance-sampling method error first increased and then decreased. Why do you think this happened?

As mentioned in the chapter, given a single observed return the numerator and denominator of the weighted-average estimate cancel so that the estimate is equal to the observed return, whose expectation is $v_b(s)$ rather than $v_\pi(s)$. The results in Figure 5.3 were averaged over 100 runs. However, the state is not necessarily visited on every episode and, in fact, it will not be visited during the majority of episodes. Suppose that the state is visited n times during the first episode over 100 runs, where n is almost certainly significantly smaller than 100. As the initial estimate of the state's value is 0, in n out of 100 runs the estimate of the state's value will remain 0 after a single episode, thus pulling the average over all runs towards 0. The mean square error over 100 runs will therefore be pulled towards the square of the true value of the state, that is, $-0.27726^2 = 0.07687$.

Now after two episodes it is more likely that the state has been visited at least once, however, it is unlikely that the state has been visited during both episodes. Therefore, over 100 runs, after two episodes we will have less estimates of 0 and more biased estimates with expectation $v_b(s)$. Hence why the mean square error increases. Similarly, after three episodes we will have even more biased estimates with expectation $v_b(s)$ and it is still unlikely that the state will be visited more than once during the first three episodes in any single run. And so on.

However, after a certain number of episodes m , which appears to be around eight or nine in this example, the state will begin to be visited more than once during the first m episodes of each run. As a result the bias and hence also the mean square error will begin to decrease. As mentioned in the chapter, the bias converges asymptotically to zero, as does the square error in this example.

8 Exercise 5.8

The results with Example 5.5 and shown in Figure 5.4 used a first-visit MC method. Suppose that instead an every-visit MC method was used on the same problem. Would the variance of the estimator still be infinite? Why or why not?

For the every-visit MC method the mean of the square of the importance-sampling ratio in an episode of length n is

$$\frac{1}{n} \left(\frac{1}{0.5^2} + \frac{1}{0.5^4} + \frac{1}{0.5^6} + \dots + \frac{1}{0.5^{2n}} \right) = \frac{2^2}{n} (1 + 2^2 + 2^4 + \dots + 2^{2n-2})$$

So the expected value of the square of the importance-sampling ratio is

$$\begin{aligned} & \frac{1}{2} \cdot 0.1 \cdot 2^2 \\ & + \frac{1}{2} \cdot 0.9 \cdot \frac{1}{2} \cdot 0.1 \cdot \frac{2^2}{2} (1 + 2^2) \\ & + \frac{1}{2} \cdot 0.9 \cdot \frac{1}{2} \cdot 0.9 \cdot \frac{1}{2} \cdot 0.1 \cdot \frac{2^2}{3} (1 + 2^2 + 2^4) \\ & + \frac{1}{2} \cdot 0.9 + \frac{1}{2} \cdot 0.9 \cdot \frac{1}{2} \cdot 0.9 \cdot \frac{1}{2} \cdot 0.1 \cdot \frac{2^2}{4} (1 + 2^2 + 2^4 + 2^6) + \dots \\ & = 0.1 \sum_{k=0}^{\infty} \frac{0.9^k}{2^k} \frac{2}{k+1} (1 + 2^2 + 2^4 + \dots + 2^{2k}) \\ & = 0.2 \sum_{k=0}^{\infty} \frac{0.45^k}{k+1} (1 + 2^2 + 2^4 + \dots + 2^{2k}) \\ & > 0.2 \sum_{k=0}^{\infty} \frac{0.45^k \cdot 2^{2k}}{k+1} \\ & = 0.2 \sum_{k=0}^{\infty} \frac{1.8^k}{k+1}. \end{aligned}$$

Now for $k \geq 8$,

$$\begin{aligned} 1.8^k &> (k+1)^2 \\ \iff \frac{1.8^k}{k+1} &> k+1 \end{aligned}$$

Hence, the above sum, and therefore the variance of the estimator, is infinite.

9 Exercise 5.9

Modify the algorithm for first-visit MC policy evaluation (Section 5.1) to use the incremental implementation for sample averages described in Section 2.4.

We would initialise $N(s) \leftarrow 0$ for all $s \in \mathcal{S}$ and then, given the S_t is the first visit to some state in some episode, we would increment $N(S_t)$ by 1 and update $V(S_t)$ as follows:

$$V(S_t) \leftarrow V(S_t) + \frac{1}{N(S_t)} [G - V(S_t)].$$

10 Exercise 5.10

Derive the weighted-average update rule (5.8) from (5.7). Follow the pattern of the derivation of the unweighted rule (2.3).

We have

$$\begin{aligned} V_{n+1} &= \frac{\sum_{k=1}^n W_k G_k}{\sum_{k=1}^n W_k} \\ &= \frac{1}{\sum_{k=1}^n W_k} \left(W_n G_n + \sum_{k=1}^{n-1} W_k G_k \right) \\ &= \frac{1}{\sum_{k=1}^n W_k} \left(W_n G_n + \sum_{k=1}^{n-1} W_k \frac{\sum_{k=1}^{n-1} W_k G_k}{\sum_{k=1}^{n-1} W_k} \right) \\ &= \frac{1}{\sum_{k=1}^n W_k} \left(W_n G_n + V_n \sum_{k=1}^{n-1} W_k \right) \\ &= \frac{1}{\sum_{k=1}^n W_k} \left(W_n G_n + V_n \left(\sum_{k=1}^n W_k - W_n \right) \right) \\ &= \frac{1}{\sum_{k=1}^n W_k} \left(W_n G_n + V_n \sum_{k=1}^n W_k - V_n W_n \right) \\ &= V_n + \frac{W_n}{\sum_{k=1}^n W_k} [G_n - V_n] \\ &= V_n + \frac{W_n}{C_n} [G_n - V_n] \end{aligned}$$

where $C_n = \sum_{k=1}^n W_k$.

11 Exercise 5.11

In the boxed algorithm for off-policy MC control, you may have been expecting the W update to have involved the importance-sampling ratio $\frac{\pi(A_t|S_t)}{b(A_t|S_t)}$, but instead it involves $\frac{1}{b(A_t|S_t)}$. Why is this nevertheless correct?

Because π is a deterministic greedy policy, it always chooses the current best action. On the third last line of the algorithm $\pi(S_t)$ is set to the current best action and on the second last line we exit the inner loop if $A_t \neq \pi(S_t)$. Therefore, the W update on the last line is only executed if $A_t = \pi(S_t)$, in which case $\pi(A_t|S_t) = 1$.

12 Exercise 5.13

Show the steps to derive (5.14) from (5.12).

