

RBE 595 — Reinforcement Learning
Week #5 Assignment

Arjan Gupta

Problem 1

When is it suited to apply Monte-Carlo to a problem?

Answer

Monte-Carlo methods are best suited to be applied to problems where we do not have a model of the environment (i.e., the dynamics of the environment are unknown). For example, sometimes it is simply not practical to model the complexity of the environment. In such cases, the agent must learn about the environment by interacting with it and using the obtained rewards to update its policy via the action-value function.

Problem 2

When does the Monte-carlo prediction performs the first update?

Answer

The Monte-Carlo prediction performs the first update after an episode terminates. This is because the Monte-Carlo method is an episodic method, i.e., it learns from a series of state, action, and reward tuples that occur in an episode.

Problem 3

What is off-policy learning and why it is useful?

Answer

Off-policy learning is a method of reinforcement learning where the agent learns about the environment by observing the behavior of another agent, called the *behavior policy*, which is the policy responsible for exploration and interaction. However, the agent performs evaluation and optimization using a different policy, called the *target policy*.

Off-policy learning is useful for the following reasons,

- It avoids the unlikely assumption of exploring starts. In some Monte Carlo algorithms, the exploratory behavior comes from random starting states, however this is not always possible. Off-policy learning allows the agent to learn about the environment without this assumption.
- Existing knowledge can be leveraged by learning from the behavior of other agents. The behavior policy can be a simple random policy, or it can be a policy that has been learned from past experience.
- It avoids the situation where the agent is stuck in a suboptimal policy because it is not exploring enough. This would happen if the agent is using a greedy and deterministic policy to learn about the environment.
- It avoids unexpected actions that may occur during exploration. Instead, the behavior policy can continuously explore while the target policy learns. This is particularly useful in cases where the environment is dangerous or expensive to explore, or if humans are involved.

Problem 4

(Exercise 5.5, page 105) 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?

Answer

Since this problem does not involve a behavior and target policy, we will not use the importance sampling ratio. Instead, we can manually calculate the first-visit and every-visit estimators of the value of the non-terminal state.

Given episode

Let the nonterminal state be s (and let s_i denote the i^{th} time s was visited) and the terminal state be s' . Let the action be a (and let a_i denote the i^{th} time s was visited). The reward is $r = 1$ for all transitions. Also, $\gamma = 1$.

The given episode is as follows,

$$s_0 \xrightarrow{a_1} s_1 \xrightarrow{a_2} s_2 \xrightarrow{a_3} s_3 \xrightarrow{a_4} s_4 \xrightarrow{a_5} s_5 \xrightarrow{a_6} s_6 \xrightarrow{a_7} s_7 \xrightarrow{a_8} s_8 \xrightarrow{a_9} s_9 \xrightarrow{a_{10}} s'$$

As we can see by the subscript of a , there are 10 rewards of 1 each. Therefore, the total return is 10.

First-Visit Estimator

The first-visit estimator of the value of the nonterminal state is calculated as follows,

$$\begin{aligned} V(s) &= 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 + 1 \\ &= 1(10) = 10 \end{aligned}$$

Every-Visit Estimator

The every-visit estimator of the value of the nonterminal state is calculated as follows,

$$\begin{aligned} V(s) &= \frac{1}{10}(1 + 2 + 3 + 4 + 5 + 6 + 7 + 8 + 9 + 10) \\ &= \frac{1}{10}(55) = 5.5 \end{aligned}$$

Problem 5

(Exercise 5.7, page 108) 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?

Answer

As shown in the lectures by Dr. Navid Dadkhah Tehrani, here is a table showing the bias and variance comparison between the ordinary importance-sampling method and the weighted importance-sampling method.

	Ordinary Importance-Sampling	Weighted Importance-Sampling
Bias	Un-biased	Biased (eventually unbiased)
Variance	Large	Low

For weighted importance-sampling, the bias is initially high, but it eventually becomes unbiased. That initial bias is what causes the error to increase initially. However, after a large number of episodes, the bias decreases and the error decreases as well.

Problem 6

(Exercise 5.8, page 108) 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?

Answer

No, the variance of the estimator would not be infinite. In fact, after the first episode that ends with the left action, the value of the state $V(s)$ would be 1 and would remain so for all subsequent episodes.

Problem 7

(Exercise 3.17) What is the Bellman equation for action values, that is, for q_π ? It must give the action value $q_\pi(s, a)$ in terms of the action values, $q_\pi(s', a')$, of possible successors to the state-action pair (s, a) . Hint: the backup diagram below corresponds to this equation. Show the sequence of equations analogous to (3.14), but for action values.

Answer

From the textbook, the action-value function for a policy π is defined as,

$$\begin{aligned}
 q_\pi(s, a) &\doteq \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a] \\
 &= \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right] \\
 &= \mathbb{E}_\pi \left[R_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k R_{t+k+2} \mid S_t = s, A_t = a \right] \\
 &= \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \\
 &= \mathbb{E}_\pi [R_{t+1} \mid S_t = s, A_t = a] + \gamma \mathbb{E}_\pi [G_{t+1} \mid S_t = s, A_t = a]
 \end{aligned}$$

Now, let us consider the first and second terms of the above equation separately.

First Term

$$\mathbb{E}_\pi [R_{t+1} \mid S_t = s, A_t = a] = \sum_{r \in \mathcal{R}} r \cdot p(r \mid s, a) = \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} r \cdot p(s', r \mid s, a)$$

Second Term

$$\begin{aligned}
 \gamma \mathbb{E}_\pi [G_{t+1} \mid S_t = s, A_t = a] &= \gamma \sum_{g \in \mathcal{G}} g \cdot p(g \mid s, a) \\
 &= \gamma \sum_{g \in \mathcal{G}} \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} g \cdot p(g \mid s', a') \cdot p(s', r \mid s, a) \cdot \pi(a' \mid s')
 \end{aligned}$$

Where, $\sum_{g \in \mathcal{G}} g \cdot p(g \mid s', a') = \mathbb{E}_\pi [G_{t+1} \mid S_{t+1} = s', A_{t+1} = a'] = q_\pi(s', a')$

Therefore the second term is,

$$\gamma \mathbb{E}_\pi [G_{t+1} \mid S_t = s, A_t = a] = \gamma \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} q_\pi(s', a') \cdot p(s', r \mid s, a) \cdot \pi(a' \mid s')$$

Now, combining the first and second terms, we get,

$$q_\pi(s, a) = \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} r \cdot p(s', r | s, a) + \gamma \sum_{r \in \mathcal{R}} \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} q_\pi(s', a') \cdot p(s', r | s, a) \cdot \pi(a' | s')$$

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) \left[r + \gamma \sum_{a'} \pi(a' | s') q_\pi(s', a') \right]$$

Which is the Bellman equation for action values, i.e., for q_π .

Backup Diagram Confirmation

This equation can be verified by looking at the backup diagram given in the prompt. The backup diagram shows that we start with the state-action pair (s, a) . To get to the next state, we are subjected to the environment $p(s', r | s, a)$. The reward r is added to the discounted return G_{t+1} . This brings us to our new state, s' . At this point, the equation would look as follows,

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]$$

However we still need to eliminate the $v_\pi(s')$ term. To do this, we go through our policy, π , to get the action a' that we would take in the state s' . Now the equation becomes,

$$q_\pi(s, a) = \sum_{s', r} p(s', r | s, a) \left[r + \gamma \sum_{a'} \pi(a' | s') q_\pi(s', a') \right]$$

So, the Bellman equation for action values, i.e., for q_π , is confirmed by the backup diagram.

Problem 8

(Exercise 3.22) Consider the continuing MDP shown below. The only decision to be made is that in the top state, where two actions are available, left and right. The numbers show the rewards that are received deterministically after each action. There are exactly two deterministic policies, π_{left} and π_{right} . What policy is optimal if $\gamma = 0$? If $\gamma = 0.9$? If $\gamma = 0.5$?

Answer

The discounted return is defined as,

$$G_t \doteq R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \cdots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (3.8)$$

Case 1: $\gamma = 0$

When $\gamma = 0$, the left policy rewards are calculated as follows,

$$G_{\text{left}} = 1 + 0 + 0 + \cdots = 1$$

Similarly, the right policy rewards are calculated as follows,

$$G_{\text{right}} = 0 + 0 + \cdots = 0$$

In this case, the **left** policy is optimal.

Case 2: $\gamma = 0.9$

When $\gamma = 0.9$, the left policy rewards are calculated as follows,

$$\begin{aligned} G_{\text{left}} &= 1 + 0.9 \cdot 0 + 0.9^2 \cdot 1 + \cdots \\ &= 1 + 0.9^2 + 0.9^4 + \cdots \\ &= \sum_{k=0}^{\infty} 0.9^{2k} \\ &= \sum_{k=0}^{\infty} 0.81^k \\ &= \frac{1}{1 - 0.81} = \frac{1}{0.19} \\ &= 5.263 \end{aligned}$$

Similarly, the right policy rewards are calculated as follows,

$$\begin{aligned} G_{\text{right}} &= 0 + 0.9 \cdot 2 + 0 + 0.9^3 \cdot 2 + \cdots \\ &= 0.9 \cdot 2 + 0.9^3 \cdot 2 + \cdots \\ &= 2 \cdot \sum_{k=0}^{\infty} 0.9^{2k+1} = 2 \cdot \sum_{k=0}^{\infty} (0.9)(0.81)^k = 2 \cdot \frac{0.9}{1 - 0.81} \\ &= \frac{1.8}{0.19} = 9.474 \end{aligned}$$

In this case, the **right** policy is optimal.

Case 3: $\gamma = 0.5$

When $\gamma = 0.5$, the left policy rewards are calculated as follows,

$$\begin{aligned}
 G_{\text{left}} &= 1 + 0.5 \cdot 0 + 0.5^2 \cdot 1 + \dots \\
 &= 1 + 0.5^2 + 0.5^4 + \dots \\
 &= \sum_{k=0}^{\infty} 0.5^{2k} = \sum_{k=0}^{\infty} 0.25^k \\
 &= \frac{1}{1 - 0.25} = \frac{1}{0.75} \\
 &= 1.333
 \end{aligned}$$

Similarly, the right policy rewards are calculated as follows,

$$\begin{aligned}
 G_{\text{right}} &= 0 + 0.5 \cdot 2 + 0 + 0.5^3 \cdot 2 + \dots \\
 &= 0.5 \cdot 2 + 0.5^3 \cdot 2 + \dots \\
 &= 2 \cdot \sum_{k=0}^{\infty} 0.5^{2k+1} = 2 \cdot \sum_{k=0}^{\infty} (0.5)(0.25)^k = 2 \cdot \frac{0.5}{1 - 0.25} \\
 &= \frac{1}{0.75} = 1.333
 \end{aligned}$$

In this case, both the **left** and **right** policies are optimal.