

**RBE 595 — Reinforcement Learning**  
**Week #7 Assignment**  
**Temporal Difference Learning**

Arjan Gupta

## Problem 1

Between DP (Dynamic Programming), MC (Monte-Carlo) and TD (Temporal Difference), which one of these algorithms use bootstrapping? Explain.

### Answer

Bootstrapping is the process of updating the value of a state based on the value of a future state.

- **Dynamic Programming** (DP) uses bootstrapping. This is because DP uses the Bellman equation to update the value of a state based on the value of a future state.
- **Monte-Carlo** (MC) does not use bootstrapping. This is because MC does not use the Bellman equation to update the value of a state based on the value of a future state. Instead, MC uses the actual return value to update the value of a state.
- **Temporal Difference** (TD) uses bootstrapping. This is because TD uses the Bellman equation to update the value of a state based on the value of a future state.

## Problem 2

We mentioned that the target value for TD is  $[R_{t+1} + \gamma V(s_{t+1})]$ . What is the target value for Monte-carlo, Q-learning, SARSA and Expected-SARSA?

### Answer

- **Monte-Carlo** (MC) does not use bootstrapping. Therefore, the target value is the actual return value,  $G_t$ .
- **Q-Learning** is an off-policy TD control algorithm. Therefore, the target value is  $R_{t+1} + \gamma \max_a Q(S_{t+1}, a)$ .
- **SARSA** is an on-policy TD control algorithm. Therefore, the target value is  $R_{t+1} + \gamma Q(S_{t+1}, A_{t+1})$ .
- **Expected-SARSA** is an on-policy TD control algorithm. Therefore, the target value is  $R_{t+1} + \gamma \mathbb{E}_{\pi} [Q(S_{t+1}, A_{t+1}) \mid S_{t+1}]$ .

## Problem 3

What are the similarities of TD and MC?

### Answer

The similarities between TD and MC are as follows:

- Both TD and MC are model-free.
- Both TD and MC are used for prediction and control.

## Problem 4

Assume that we have two states  $x$  and  $y$  with the current value of  $V(x) = 10$ ,  $V(y) = 1$ . We run an episode of  $\{x, 3, y, 0, y, 5, T\}$ . What's the new estimate of  $V(x)$ ,  $V(y)$  using TD (assume step size  $\alpha = 0.1$  and discount rate  $\gamma = 0.9$ ).

### Answer

The new estimate of  $V(x)$  is as follows:

$$\begin{aligned}
 V(x) &= V(x) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(x)] \\
 &= 10 + 0.1 [3 + 0.9 \cdot 1 - 10] \\
 &= 10 + 0.1 [3.9 - 10] \\
 &= 10 + 0.1 [-6.1] \\
 &= 10 - 0.61 \\
 &= 9.39
 \end{aligned}$$

However,  $V(y)$  gets updated twice in this episode. The first update is as follows:

$$\begin{aligned}
 V(y) &= V(y) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(y)] \\
 &= 1 + 0.1 [0 + 0.9 \cdot 1 - 1] \\
 &= 1 + 0.1 [0.9 - 1] \\
 &= 1 + 0.1 [-0.1] \\
 &= 1 - 0.01 \\
 &= 0.99
 \end{aligned}$$

The second update is as follows:

$$\begin{aligned}
 V(y) &= V(y) + \alpha [R_{t+1} + \gamma V(S_{t+1}) - V(y)] \\
 &= 0.99 + 0.1 [5 + 0.9 \cdot 0 - 0.99] \\
 &= 0.99 + 0.1 [5 - 0.99] \\
 &= 0.99 + 0.1 [4.01] \\
 &= 0.99 + 0.401 \\
 &= 1.391
 \end{aligned}$$

Therefore, the new estimate of  $V(x)$  is 9.39 and the new estimate of  $V(y)$  is 1.391.

## Problem 5

Can we consider TD an online (real-time) method and MC an offline method? Why?

### Answer

Yes, we can consider TD an online (real-time) method and MC an offline method. This is because TD learns during the episode, whereas MC learns after the episode has ended. TD updates the value of a state based on the value of the next state (during the episode), whereas MC updates the value of a state based on the actual return value (after the entire episode has ended).

## Problem 6

Does Q-learning learn the outcome of exploratory actions? (Refer to the Cliff walking example).

### Answer

Yes, Q-learning learns the outcome of exploratory actions. This is because Q-learning is an off-policy TD control algorithm. Therefore, Q-learning learns the optimal policy,  $\pi_*$ , which is the policy that maximizes the value function,  $q_*$ , i.e.,  $\pi_* = \arg \max_{\pi} q_*(s, a)$ . This means that Q-learning learns the optimal policy,  $\pi_*$ , even if the behavior policy,  $b$ , is exploratory.

## Problem 7

(Exercise 3.17) What is the Bellman equation for action values, that is, for  $q_\pi$ ? 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  $\pi$  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  $q_\pi$ .

### 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,  $\pi$ , 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  $q_\pi$ , 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,  $\pi_{\text{left}}$  and  $\pi_{\text{right}}$ . 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.