

**Written Exercises: Chapter 4**

## EXERCISE 4.1

We are not using any discount, so we have

$$q_\pi(11, \text{down}) = r + v_\pi(T) \quad (1)$$

$$= -1 + 0 = -1 \quad (2)$$

$$q_\pi(7, \text{down}) = r + v_\pi(11) \quad (3)$$

$$= -1 + -14 = -15 \quad (4)$$

## EXERCISE 4.2

Adding the new state, we get

$$v_\pi(15) = \sum_a \pi(a \mid 15)(-1 + v_\pi(s')) \quad (5)$$

$$= 0.25 \cdot (-1 + v_\pi(13)) + 0.25 \cdot (-1 + v_\pi(12)) + 0.25 \cdot (-1 + v_\pi(14)) + 0.25 \cdot (-1 + v_\pi(15)) \quad (6)$$

$$= 0.25(-4 + v_\pi(13) + v_\pi(12) + v_\pi(14) + v_\pi(15)) \quad (7)$$

$$= -1 + 0.25(-20 - 22 - 14 + v_\pi(15)) \quad (8)$$

$$= -1 - 14 + 0.25v_\pi(15) \quad (9)$$

$$\implies 0.75v_\pi(15) = -15 \quad (10)$$

$$\implies v_\pi(15) = \frac{-15}{0.75} = -20 \quad (11)$$

Now, if we change the dynamics such that moving down in state 13 will move us to 15, the value function will not change. This is because  $v_\pi(13) = v_\pi(15)$ , so  $q_\pi(13, \text{down}) = -20$  hold regardless of whether the move will take us to 15 or not.

## EXERCISE 4.3

The equations analogous to (4.3) and (4.4) are given in Exercise 3.17.

The update equation for  $q_\pi(s, a)$  is given by

$$q_{k+1}(s, a) = \mathbb{E}_\pi [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \quad (12)$$

$$= \sum_{s', r} p(s', r \mid s, a) \left( r + \gamma \sum_{a'} \pi(a' \mid s) q_k(a', s) \right) \quad (13)$$

## EXERCISE 4.4

We need to loop through the computed value function to check whether they are equal; two optimal policies will have the same value function.

#### EXERCISE 4.5

The process is very similar - we use the Bellman equation for  $q_\pi$  to evaluate the policy. To improve the policy, we again create the greedy policy  $\pi'$  with respect to  $q_\pi$ . That is,

- (1) Initialization  
 $q(s, a) \in \mathbb{R}$  and  $\pi(s) \in \mathcal{A}(s)$  arbitrarily for all  $s \in \mathcal{S}$ .
- (2) Policy evaluation  
 Loop:  
  - $\Delta \leftarrow 0$
  - Loop for each  $s \in \mathcal{S}$  and each  $a \in \mathcal{A}(s)$  :  
 -  $q \leftarrow q(s, a)$   
 -  $q(s, a) \leftarrow \sum_{s', r} p(s', r \mid s, a) (r + \gamma q(s', \pi(s')))$   
 -  $\Delta \leftarrow \max(\Delta, |q - q(s, a)|)$
  - until  $\Delta < \theta$
- (3) Policy Improvement  
  - *policy-stable*  $\leftarrow true$
  - For each  $s \in \mathcal{S}$ :  
 - *old-action*  $\leftarrow \pi(s)$   
 -  $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s', r} q(s, a)$   
 - If *old-action*  $\neq \pi(s)$  then *policy-stable*  $\leftarrow false$ .
  - If *policy-stable*, then stop and return  $q \approx q_*$  and  $\pi \approx \pi_*$ , else go to 2.

#### EXERCISE 4.6

We assume that there would be a  $1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|}$  probability on the optimal action, and the remaining  $\epsilon - \frac{\epsilon}{|\mathcal{A}(s)|}$  probability would be evenly spread over the remaining actions.

Now, assume that we store this action in  $A(s)$ . Then,

$$\pi_A(a \mid s) = \begin{cases} 1 - \epsilon + \frac{\epsilon}{|\mathcal{A}(s)|} & \text{if } a = A(s) \\ \epsilon - \frac{\epsilon}{|\mathcal{A}(s)|} & \text{otherwise} \end{cases} \quad (14)$$

Then, step 3 would be

- *policy-stable*  $\leftarrow true$
- For each  $s \in \mathcal{S}$  :  
 - *old-action*  $\leftarrow A(s)$   
 -  $A(s) \leftarrow \operatorname{argmax}_{a'} \sum_a \pi'_a(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')]$   
 -  $\pi(s) \leftarrow \pi_A(s)$   
 - If *old-action*  $\neq A(s)$  then *policy-stable*  $\leftarrow false$
- ...

For step two, the value update would simply be

$$V(s) \leftarrow \sum_a \pi_A(s)(a \mid s) \sum_{s', r} p(s', r \mid s, a) [r + \gamma V(s')] \quad (15)$$

where the rest of the step remains the same.

Step one would consist of letting  $V(s) \in \mathbb{R}$  arbitrarily, and setting  $A(s)$  for  $s \in \mathcal{S}$  arbitrarily, with  $\pi(a \mid s) = \pi_{A(s)}(a \mid s)$  as defined above.