

Exercise 4.1

In Example 4.1, if π is the equiprobable random policy, what is $q_\pi(11, \text{down})$? What is $q_\pi(7, \text{down})$?

Answer

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

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \cdot q_\pi(s, a)$$

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

$$Q_{\pi}(11, \text{down}) = r + 1 \cdot 0 = -1.$$

$$Q_{\pi}(7, \text{down}) = r + 1 \cdot v_{\pi}(11) = -1 - 1 \cdot 14 = -15.$$

$$\begin{aligned} q_\pi(7, \text{down}) &= \mathbb{E}_\pi[G_t | S_t = s, A_t = a] \\ &= R_t + \underbrace{\gamma}_{1} \underbrace{\mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']}_{v(11)} \\ &= -1 - 14 \\ &= -15. \end{aligned}$$

Exercise 4.2

In Example 4.1, suppose a new state 15 is added to the gridworld just below state 13, and its actions, left, up, right, and down, take the agent to states 12, 13, 14, and 15, respectively. Assume that the transitions from the original states are unchanged. What, then, is $v_\pi(15)$ for the equiprobable random policy? Now suppose the dynamics of state 13 are also changed, such that action down from state 13 takes the agent to the new state 15. What is $v_\pi(15)$ for the equiprobable random policy in this case?

Answer

unchanged policy

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

$$v_\pi(15) = -20.$$

Changed policy

Suppose the dynamics of 13 also changed. We do the iterative policy evaluation. The initialization is natural: we let $V_0(s) = v_\pi(s)$, where v_π is the value when the dynamics are unchanged, as shown in the $k = \infty$ case in Figure 4.1.

$V_0(15) = -20$, where -20 is what we just derived for $v_\pi(15)$.

Particularly, we do the “immediate overwriting”. First we update the value of 13.

$$v(13) = 0.25 * (-1 - 20 - 1 - 22 - 1 - 14 - 1 - 20) = 20.$$

We then proceed to update $V(15)$ with $V(13)$:

$$V(13) = 0.25 * (-1 - 22 - 1 - 20 - 1 - 14 - 1 - 20) = 20.$$

As we proceed, we observe no change – the initialization is good enough. Hence the iterative policy evaluation ends with 1 iteration. $V_1(15) = -20$ for this case.

Exercise 4.3

What are the equations analogous to (4.3), (4.4), and (4.5) for the action- value function q_π and its successive approximation by a sequence of functions q_0, q_1, q_2, \dots ?

$$q_\pi(s, a) = r(s, a) + \sum_{a'} \mathbb{E}_\pi [\gamma q_\pi(S_{t+1}, a') | S_t = s, A_t = a]$$

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

$$q_{k+1}(s, a) = \sum_{s', a} p(s', r | s, a) \left[r + \sum_{a'} \pi_k(a' | s') \gamma q_k(s', a') \right]$$