

Reinforcement Learning Lab Exercise 3

Jinay Vora
Student ID: 202201473

Problem Set 2

Description

From [3], study example 4.1 and write a Python code to find $v\pi(s)$.

Solution

	1	2	3
4	5	6	7
8	9	10	11
12	13	14	

Python Code

```
1 import numpy as np
2
3 def in_bounds(row, col):
4     return 0 <= row < 4 and 0 <= col < 4
5
6 def state_to_rowcol(s):
7     s0 = s - 1
8     return divmod(s0, 4)
```

```

1 def rowcol_to_state(row, col):
2     return row * 4 + col + 1
3
4 def next_state(s, action):
5     if s in [1, 16]:
6         return s
7
8     row, col = state_to_rowcol(s)
9
10    if action == 'U':
11        new_row, new_col = row - 1, col
12    elif action == 'D':
13        new_row, new_col = row + 1, col
14    elif action == 'L':
15        new_row, new_col = row, col - 1
16    elif action == 'R':
17        new_row, new_col = row, col + 1
18    else:
19        raise ValueError("Unknown action!")
20
21    if in_bounds(new_row, new_col):
22        return rowcol_to_state(new_row, new_col)
23    else:
24        return s
25
26 def iterative_policy_evaluation(theta=1e-8, max_iterations=100000):
27     V = np.zeros(17)
28     gamma = 1.0
29     actions = ['U', 'D', 'L', 'R']
30
31     for _ in range(max_iterations):
32         delta = 0
33         newV = V.copy()
34
35         for s in range(1, 17):
36             if s in [1, 16]:
37                 continue
38
39             v_new = 0.0
40             for a in actions:
41                 s_next = next_state(s, a)
42                 reward = -1
43                 v_new += 0.25 * (reward + gamma * V[s_next])
44             newV[s] = v_new
45
46         delta = np.max(np.abs(newV - V))

```

```

1      V = newV
2      if delta < theta:
3          break
4
5      return V
6
7  if __name__ == "__main__":
8      V_final = iterative_policy_evaluation()
9      print("Value function with terminal states at (0,0) and (3,3):\n")
10
11     for row in range(4):
12         row_str = []
13         for col in range(4):
14             s = row * 4 + col + 1
15             row_str.append(f"{V_final[s]:6.2f}")
16     print(" ".join(row_str))

```

Solution

Value function with terminal s

0.00	-14.00	-20.00	-22.00
-14.00	-18.00	-20.00	-20.00
-20.00	-20.00	-18.00	-14.00
-22.00	-20.00	-14.00	0.00

Problem Set 3

Description

From [3], solve exercises 4.1 and 4.2. You can choose to solve it using a Python code.

Solution

Exercise 4.1

We know that,

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

But for deterministic actions,

$$q_{\pi}(s, a) = R(s, a, s') + \gamma v_{\pi}(s')$$

Using the above two equations, we get:

$$(a) \quad q_{\pi}(11, \text{down}) = -1 + 1 \cdot (0) = -1$$

$$(b) \quad q_{\pi}(7, \text{down}) = -1 + 1 \cdot (-14) = -15$$

Exercise 4.2

Part 1: Original Dynamics

The value function for state 15 is:

$$v_{\pi}(15) = -1 + \frac{1}{4} (v_{\pi}(12) + v_{\pi}(13) + v_{\pi}(14) + v_{\pi}(15))$$

Part 2: Changed Dynamics

If action down from state 13 leads to 15, the Bellman equations for $v_{\pi}(13)$ and $v_{\pi}(15)$ are:

$$v_{\pi}(13) = -1 + \frac{1}{4} (v_{\pi}(12) + v_{\pi}(13) + v_{\pi}(14) + v_{\pi}(15))$$

$$v_{\pi}(15) = -1 + \frac{1}{4} (v_{\pi}(12) + v_{\pi}(13) + v_{\pi}(14) + v_{\pi}(15))$$

Rearrange both equations:

$$3v_{\pi}(15) = -4 + v_{\pi}(12) + v_{\pi}(13) + v_{\pi}(14)$$

$$3v_{\pi}(13) = -4 + v_{\pi}(12) + v_{\pi}(14) + v_{\pi}(15)$$

Substitute $v_\pi(15)$ into the equation for $v_\pi(13)$:

$$v_\pi(13) = \frac{-4 + v_\pi(12) + v_\pi(14) + \frac{-4 + v_\pi(12) + v_\pi(13) + v_\pi(14)}{3}}{3}$$

Simplifying gives:

$$v_\pi(13) = 0.5v_\pi(12) + 0.5v_\pi(14) - 2$$

By symmetry, $v_\pi(15)$ is:

$$v_\pi(15) = 0.5v_\pi(12) + 0.5v_\pi(14) - 2 = 0.5(-22) + 0.5(-14) - 2 = -20$$

Problem Set 4

Description

From [3], solve exercise 4.3.

Solution

Analogous to Equation (4.3):

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

Analogous to Equation (4.4):

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

Analogous to Equation (4.5):

The iterative update rule for successive approximations of q_π is:

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