

Chapter 3

EXERCISE 3.1

Devise three example tasks of your own that fit into the MDP framework, identifying for each its states, actions, and rewards. Make the three examples as different from each other as possible. The framework is abstract and flexible and can be applied in many different ways. Stretch its limits in some way in at least one of your examples.

Controlling a rocket. A reinforcement learning agent is to control a rocket that should fly into orbit.

The agent will want to know all about the forces acting on the rocket, and the current position, velocity, acceleration, heading, and rotation forces of the rocket. It might need to know the fuel level, and the temperature of the engine in order to complete its task.

Finally, most rockets have discrete stages; the first transition might perhaps be to discard the level 1 boosters.

Thus the state might be composed of

- Position
- Velocity
- Acceleration
- Fuel status
- Engine temperature
- Heading
- Rotation
- The phase of the launch

The agent should control the engine and controls of the rocket. Maybe the agent also controls when the rocket releases its boosters. The agent might also want to warn the humans that something is wrong, so it might want to send a distress signal.

Thus the actions could be

- Accelerate engine
- Turn (or something to control the heading of the engine)
- Transition into next phase of launch
- Send distress signal

The agent should fly the rocket so that it heads forward, with minimum turbulence, with fuel efficiency, and a nice ascent into orbit. Its reward might be some weighted combination of these factors.

Programming. We could face a reinforcement learner create programs for tasks, given a description in some form. In fact, this has already been done.

The agent might receive as a state the programming language it should use, if it should not decide that itself. Of course it should want a description of the program.

An important thing to know when programming is the current state of the code, so the agent should receive that as well.

Now, the agent might want to run the program to see if it produces the desired output. So that might be a possible part of the state.

For the actions, the agent might choose to input a character to the file, delete text, or any other action that modifies the text. The agent might also decide to run the program to see the results. Finally, when the agent is happy, it should be able to submit its code.

Cook food. We might produce a chef agent. It might control some kind of robot that has access to all necessary appliances in a kitchen, and the ingredients, and the agent's goal is to produce the requested dish.

The states might consist of visual input, temperature of stoves and ovens, and the sensorimotor data from the robot's actuators.

The agent needs to control the robot, so the actions should include all the possible controls of the robot. Also, for better automation, stove and oven controls might be actions separate from the robot's movement.

The rewards should be negative for breaking stuff, or otherwise reeking havoc in the kitchen, and there should be positive rewards for creating the requested dish, where the size of the reward is a function of the dish's quality.

EXERCISE 3.2

Is the MDP framework adequate to usefully represent all goal-directed learning tasks? Can you think of any clear exceptions?

I cannot find any counterexamples.

EXERCISE 3.3

The actions that the agent can perform instantaneously at any time step should be the defined actions for the agent. The choice of where to drive should not be an output for a driving agent, as that is not how you drive a car, in the immediate sense.

The choice of where to drive should be something that the agent decides makes plans on how to execute, and then uses its actuators/available actions to perform.

A driving agent that would output a coordinate of where to drive would not be a driving agent, rather a guidance agent.

Of course, a driving agent could have some other agent that controls where it should drive, and the driving agent implements that plan, but that would be a separate agent, distinct from the driving agent.

Now, as stated, the actions should be something that can be outputted, and executed immediately. If the agent outputs a rotational acceleration of the steering wheel, or simply its goal angle should be something that could be immediately implemented by some microcontroller or other machinery.

Therefore that is an acceptable action.

Now, if there is some automatic system that can better automate a task than the agent can, say, a good brake system in a car that prevents skidding, the agent should not directly control the force of the brake pads, but should simply gain access to the brake pedal, and let the brake system do the rest.

EXERCISE 3.4

TABLE 1. The Table

s	a	s'	r	$p(s', r \mid s, a)$
high	search	high	1	0.4
high	search	high	0	0.6
high	search	low	1	0.4
high	search	low	0	0.6
high	wait	high	1	0.1
high	wait	high	0	0.9
high	wait	low	1	0.1
high	wait	low	0	0.9
low	search	low	1	0.2
low	search	low	0	0.8
low	wait	low	1	0
low	wait	low	0	1
low	recharge	low	1	0
low	recharge	low	0	1

EXERCISE 3.5

The formula simply becomes

$$p(s', r \mid s, a) = 1, \text{ for all } s \in \mathcal{S} \setminus \mathcal{S}^+$$

EXERCISE 3.6

The return would not just be related to $-\gamma^K$, it would equal $-\gamma^K$, where K is the number of steps until failure.

In the continuing task formulation, the agent might also take into account the next failure, in the next reset, in a given state.

EXERCISE 3.7

The expected return is always one, as the agent always escapes the maze eventually, and when it does, there is no discrimination between states. All states get an expected return of one.

To communicate to the agent that it should navigate the maze efficiently, we should penalize it for every move it makes, so that it wants to escape the maze quickly. We only instructed it to *eventually* escape the maze, no matter how long it takes.

EXERCISE 3.8

$$G_5 = 0 \quad (1)$$

$$G_4 = R_5 + \gamma G_5 = 2 + 0.5 \cdot 0 = 2 \quad (2)$$

$$G_3 = R_4 + \gamma G_4 = 3 + 0.5 \cdot 2 = 4 \quad (3)$$

$$G_2 = R_3 + \gamma G_3 = 6 + 0.5 \cdot 4 = 8 \quad (4)$$

$$G_1 = R_2 + \gamma G_2 = 2 + 0.5 \cdot 8 = 6 \quad (5)$$

$$G_0 = R_1 + \gamma G_1 = -1 + 0.5 \cdot 6 = 2 \quad (6)$$

EXERCISE 3.9

We know that

$$G_1 = \sum_{i=0}^{\infty} 7\gamma^i = 7 \sum_{i=0}^{\infty} \gamma^i \quad (7)$$

$$= 7 \frac{1}{1-\gamma} = 7 \frac{1}{0.1} \quad (8)$$

$$= 7 \cdot 10 = 70 \quad (9)$$

So

$$G_0 = 2 + \gamma G_1 \quad (10)$$

$$= 2 + 0.9 \cdot 70 \quad (11)$$

$$= 2 + 7 = 65 \quad (12)$$

EXERCISE 3.10

Let

$$\Gamma = \sum_{k=0}^{\infty} \gamma^k$$

Then

$$\Gamma - \gamma\Gamma = \sum_{k=0}^{\infty} \gamma^k - \sum_{k=1}^{\infty} \gamma^k \quad (13)$$

$$= \gamma^0 = 1 \quad (14)$$

And hence

$$(1-\gamma)\Gamma = 1 \quad (15)$$

$$\implies \Gamma = \frac{1}{1-\gamma} \quad (16)$$

EXERCISE 3.11

If A is the set of actions, and R the set of possible rewards, it is

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

EXERCISE 3.12

We want to give an equation for v_π in terms of q_π and π . This can be done by taking the action-value function, and averaging out the action, that is

$$v_\pi(s) = \mathbb{E}[q_\pi(s, a)] = \sum_{a \in A} \pi(a | s) q_\pi(s, a)$$

EXERCISE 3.13

We want to give an equation for q_π in terms of v_π and $p(s', r | s, a)$. We define $\mathbb{E}_x[f(x)]$ as the expectation of f with respect to x . With that definition, we see that

$$q_\pi = \mathbb{E}_{s', r} [r + \gamma v(s')] \quad (17)$$

$$= \sum_{s' \in S} \pi(a | s) \sum_{r \in R} p(s', r | s, a) (r + \gamma v(s')) \quad (18)$$

EXERCISE 3.14

We represent the state s as its coordinate $s = (x, y)$.

Thus,

$$v_\pi(2, 2) = 0.7 \quad (19)$$

$$v_\pi(2, 1) = 2.5 \quad (20)$$

$$v_\pi(1, 2) = 0.7 \quad (21)$$

$$v_\pi(3, 2) = 0.4 \quad (22)$$

$$v_\pi(2, 3) = -0.4 \quad (23)$$

Now, the Bellman equation for the value function v_π is

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

We now show that the Bellman equation for the value function holds for $v_\pi(2, 2)$ with respect to its neighbors. The set of actions is $A = (UP, RIGHT, DOWN, LEFT)$. However, all the actions from the state $s = (2, 2)$ result in a reward of 0.

Also, the environment is deterministic. We can now perform the calculation:

$$v_\pi(2, 2) = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')] \quad (24)$$

$$= \frac{1}{4} (\gamma v_\pi(2, 1) + \gamma v_\pi(1, 2) + \gamma v_\pi(1, 2) + \gamma v_\pi(3, 2) + \gamma v_\pi(2, 3)) \quad (25)$$

$$= \frac{0.9}{4} (2.5 + 0.7 + 0.4 + (-0.4)) \quad (26)$$

$$= 0.72 \approx 0.7 \quad (27)$$

EXERCISE 3.15

If we add a constant c to all the rewards, each reward becomes $R_t + c$, so the return becomes

$$G'_t = (R_{t+1} + c) + (R_{t+2} + c) + \dots \quad (28)$$

$$= \sum_{k=0}^{\infty} \gamma^k (R_{t+k+1} + c) \quad (29)$$

$$= \sum_{k=0}^{\infty} \gamma^k c + \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (30)$$

$$= \frac{c}{1-\gamma} + \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \quad (31)$$

$$= G_t + \frac{c}{1-\gamma} \quad (32)$$

So the value v'_π becomes

$$v'_\pi(s) = \mathbb{E}[G'_t \mid S_t = s] \quad (33)$$

$$= \mathbb{E}\left[G_t + \frac{c}{1-\gamma} \mid S_t = s\right] \quad (34)$$

$$= \mathbb{E}[G_t \mid S_t = s] + \mathbb{E}\left[\frac{c}{1-\gamma} \mid S_t = s\right] \quad (35)$$

$$= \mathbb{E}[G_t \mid S_t = s] + \frac{c}{1-\gamma} \quad (36)$$

$$= v_\pi(s) + \frac{c}{1-\gamma} \quad (37)$$

$$= v_\pi(s) + v_c \quad (38)$$

where the constant added is $v_c = \frac{c}{1-\gamma}$

EXERCISE 3.16

Lets say the task is episodic, and some episode lasts k steps. Then, at time t , the return is

$$G'_t = \sum_{i=0}^{k-t-1} \gamma^i (R_{t+i+1} + c) \quad (39)$$

$$= \sum_{i=0}^{k-t-1} \gamma^i c + \sum_{i=0}^{k-t-1} R_{t+i+1} \quad (40)$$

$$= \frac{c(\gamma^{k-t-1} - 1)}{\gamma - 1} + \sum_{i=0}^{k-t-1} R_{t+i+1} \quad (41)$$

$$(42)$$

So here we see that the return depends on the number of steps remaining until the end of the episode.

The feature of the continuing tasks that caused there to be a constant added to the value (or return), is that the sum was infinite.

EXERCISE 3.17

We derive the Bellman equation for the action-value function:

$$q(s, a) = \mathbb{E}_\pi [G_t \mid S_t = s, A_t = a] \quad (43)$$

$$= E_\pi [R_{t+1} + \gamma G_{t+1} \mid S_t = s, A_t = a] \quad (44)$$

$$= \sum_{s'} \sum_r p(s', r \mid s, a) (r + \gamma \mathbb{E}_\pi [G_{t+1} \mid S_{t+1} = s']) \quad (45)$$

$$= \sum_{s', r} p(s', r \mid s, a) (r + \gamma \sum_{a'} \pi(a' \mid s') \mathbb{E} [G_{t+1} \mid s', a']) \quad (46)$$

$$= \sum_{s', r} p(s', r \mid s, a) \left(r + \gamma \sum_{a'} \pi(a' \mid s) q(a', s') \right) \quad (47)$$

$$(48)$$

EXERCISE 3.18

This is simply

$$v_\pi(s) = \mathbb{E}_\pi [q_\pi(a, s) \mid S_t = s] \quad (49)$$

$$= \sum_a \pi(a \mid s) q_\pi(a, s) \quad (50)$$

$$(51)$$

EXERCISE 3.19

This is

$$q_\pi(s, a) = \mathbb{E} [R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s, A_t = a] \quad (52)$$

$$= \sum_{s', r} p(s', r \mid s, a) (r + \gamma v_\pi(s')) \quad (53)$$

$$(54)$$

EXERCISE 3.20

It is the same as the figure of $q_*(s, driver)$, but the whole green area has value -1 .

EXERCISE 3.21

The same as v_{putt} .

EXERCISE 3.22

We consider the top state to be the only state, as otherwise the definition of the optimal policy would not work. This is because the optimal policy must satisfy $\pi_*(s) \geq \pi(s)$ for all other policies π , which would not be the case if this was the given MDP.

f $\gamma = 0$, then π_{left} is the optimal policy, as the return in the top state is then $+1$. If $\gamma = 0.9$, however, then the values of the top state for the two policies are

$$v_{\pi_{right}}(s) = \sum_{i=0}^{\infty} 0.9^{2i+1} \cdot 2 \quad (55)$$

$$= \sum_{i=0}^{\infty} 0.81^i \cdot 0.9 \cdot 2 \quad (56)$$

$$= 1.8 \sum_{i=1}^{\infty} 0.81^i \quad (57)$$

$$= 1.8 \left(\frac{1}{1 - 0.81} - 0.81 \right) = 8.19 \quad (58)$$

$$v_{\pi_{left}}(s) = \sum_{i=0}^{\infty} 0.9^{2i} \cdot 1 \quad (59)$$

$$= \sum_{i=0}^{\infty} 0.81^i \quad (60)$$

$$= \frac{1}{1 - 0.81} = 5.26 \quad (61)$$

$$(62)$$

So in that case π_{right} is better. For $\gamma = 0.5$, we have

$$v_{\pi_{right}}(s) = \sum_{i=0}^{\infty} 0.5^{2i+1} \cdot 2 \quad (63)$$

$$= \sum_{i=0}^{\infty} 0.25^i \cdot 0.5 \cdot 2 \quad (64)$$

$$= \sum_{i=1}^{\infty} 0.25^i \quad (65)$$

$$= \frac{1}{1 - 0.25} - 0.25 = 1.08 \quad (66)$$

$$v_{\pi_{left}}(s) = \sum_{i=0}^{\infty} 0.5^{2i} \quad (67)$$

$$= \sum_{i=0}^{\infty} 0.25^i \quad (68)$$

$$= \frac{1}{1 - 0.25} = 1.333 \quad (69)$$

$$(70)$$

so π_{left} is better.

EXERCISE 3.23

TODO

EXERCISE 3.24

We first need to find the value of γ used. Notice that since

$$G_t = R_{t+1} + \gamma G_{t+1} \quad (71)$$

$$\implies \gamma = \frac{G_t - R_{t+1}}{G_{t+1}} \quad (72)$$

we can use two given values to find γ . We use the state A . Then $G_t = 24.4$, $R_{t+1} = 10$, and $G_{t+1} = 16$. Thus

$$\gamma = \frac{24.4 - 10}{16} = 0.9 \quad (73)$$

The following python program will find v_* to any given accuracy:

```

1  import numpy as np
2
3  UP=1
4  DOWN=2
5  LEFT=3
6  RIGHT=4
7
8  GAMMA = 0.9
9
10 pistar = [
11     [RIGHT, DOWN, LEFT, UP, LEFT],
12     [UP, UP, UP, LEFT, LEFT],
13     [UP for _ in range(5)],
14     [UP for _ in range(5)],
15     [UP for _ in range(5)]
16 ]
17
18 def action(x, y, dir):
19     if x == 0 and y == 1:
20         return 4, 1, 10
21
22     if x == 0 and y == 3:
23         return 2, 3, 5
24
25     if dir == UP:
26         x -= 1
27     elif dir == DOWN:
28         x += 1
29     elif dir == LEFT:
30         y -= 1
31     elif dir == RIGHT:
32         y += 1
33     return x, y, 0
34
35 def get_valuestar(x, y):
36     cur_gamma = 1
37     cur_value = 0
38     for i in range(100):
39         x, y, r = action(x, y, pistar[x][y])
40         cur_value += cur_gamma * r
41         cur_gamma *= GAMMA
42     return cur_value
43
44 values = np.zeros((5, 5))
45

```

```

46  for i in range(5):
47      for j in range(5):
48          values[i][j] = get_valuestar(i, j)
49
50  print(values)

```

And we obtained

TABLE 2. Gridworld values

21.977	24.419	21.977	19.419	17.477
19.780	21.977	19.780	17.802	16.022
17.802	19.780	17.802	16.022	14.419
16.022	17.802	16.022	14.419	12.977
14.419	16.022	14.419	12.977	11.680

EXERCISE 3.25

We of course have

$$v_*(s) = \max_{a \in \mathcal{A}(s)} q_*(s, a) \quad (74)$$

EXERCISE 3.26

Since we have (74), we have

$$q_*(s, a) = \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')] \quad (75)$$

EXERCISE 3.27

This is simply

$$\pi_*(s) = \operatorname{argmax}_{a \in \mathcal{A}(s)} q_*(s, a) \quad (76)$$

EXERCISE 3.28

This is

$$\pi_*(s) = \operatorname{argmax}_{a \in \mathcal{A}(s)} \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_*(s')] \quad (77)$$

EXERCISE 3.29

In terms of $p(s' \mid s, a)$ and $r(s, a)$, we have

$$v_\pi(s) = \sum_a \pi(a \mid s) \left(r(s, a) + \gamma \sum_{s'} p(s' \mid s, a) v_\pi(s') \right) \quad (78)$$

$$v_*(s) = \max_a \left(r(s, a) + \gamma \sum_{s'} p(s' \mid s, a) v_*(s') \right) \quad (79)$$

$$q_\pi(s, a) = r(s, a) + \gamma \left(\sum_{s'} p(s' \mid s, a) \sum_{a'} \pi(a' \mid s') q_*(s', a') \right) \quad (80)$$

$$q_*(s, a) = r(s, a) + \gamma \left(\max_{a'} \sum_{s'} p(s' \mid s, a') q_*(s', a') \right) \quad (81)$$