

Planning, Learning and Decision Making:

Homework 2. Markov decision problems

Mykhaylo Marfeychuk, Roy De Prins

16 October 2018

1 Exercise 1) a

State Space $\mathcal{X} =$

{TopLeft, TopRight, BottomLeft, BottomRight,
Red, TopLeftRed, TopRightRed, BottomLeftRed, BottomRightRed,
Blue, TopLeftBlue, TopRightBlue, BottomLeftBlue, BottomRightBlue,
Goal }.

Action Space $\mathcal{A} =$

{Up, Down, Left, Right }.

2 Exercise 1) b

Transition Probability Matrix for action 'right':

	TL	TR	BL	BR	Red	TLR	TRR	BLR	BRR	Blue	TLB	TRB	BLB	BRB	Goal
TL	0.2	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
TR	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
BL	0.0	0.0	0.2	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
BR	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Red	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
TLR	0.0	0.0	0.0	0.0	0.0	0.2	0.8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
TRR	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
BLR	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.8	0.0	0.0	0.0	0.0	0.0	0.0
BRR	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0	0.0	0.0	0.0
Blue	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.0	0.8	0.0	0.0	0.0
TLB	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.8	0.0	0.0	0.0
TRB	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0	0.0	0.0	0.0
BLB	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.8	0.0
BRB	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.2	0.8
Goal	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.0

Cost Function for action 'right':

	cost_right
TL	0.5
TR	1.0
BL	0.5
BR	1.0
Red	1.0
TLR	0.5
TRR	1.0
BLR	0.5
BRR	1.0
Blue	0.0
TLB	0.0
TRB	1.0
BLB	0.0
BRB	0.0
Goal	0.0

3 Exercise 1) c

We compute the cost-to-go function with the 'Value Iteration, 2.0' algorithm. So by computing the following formula:

$$\mathbf{J}^t = \mathbf{c}_{right} + \gamma \mathbf{P}_{right} \mathbf{J}^{t-1} \quad (1)$$

The cost-to-go function (calculated in Python) associated with the policy 'right' is given in the next figure.

```

1  import numpy as np
2
3  cost_right = np.array([[0.5], [1], [0.5], [1], [1], [0.5], [1], [0.5], [1], [0], [0], [1], [0], [0], [0]])
4  prob_right = np.array([[0.2, 0.8, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
5                          [0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
6                          [0.0, 0.0, 0.2, 0.8, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
7                          [0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
8                          [0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
9                          [0.0, 0.0, 0.0, 0.0, 0.0, 0.2, 0.8, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
10                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
11                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.2, 0.8, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
12                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0],
13                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.2, 0.8, 0.0, 0.0, 0.0, 0.0],
14                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.2, 0.8, 0.0, 0.0, 0.0],
15                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0, 0.0, 0.0, 0.0],
16                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.2, 0.8, 0.0],
17                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.2, 0.8],
18                         [0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 0.0, 1.0]])
19
20  gamma = 0.9
21
22  J = np.zeros((15,1))
23  err = 1
24  i = 0
25
26  while err > 1e-15:
27      Jnew = cost_right + gamma * prob_right.dot(J)
28      err = np.linalg.norm(Jnew - J)
29      i += 1
30      J = Jnew
31
32  print("The cost-to-go function is \n", J)

```

Run: HW2 x

```

The cost-to-go function is
[[ 9.3902439]
 [10.      ]
 [ 9.3902439]
 [10.      ]
 [10.      ]
 [ 9.3902439]
 [10.      ]
 [ 9.3902439]
 [10.      ]
 [ 8.7804878]
 [ 8.7804878]
 [10.      ]
 [ 0.      ]
 [ 0.      ]
 [ 0.      ]]

```

Figure 1: The Python program for calculating to cost-to-go function.