

COURSEWORK 1

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

Reinforcement Learning

Author:

Arnold Cheung (CID: 01184493)

Date: October 29, 2019

1 Understanding of MDPS

1.1

The personalised trace of CID:01184493 is:

$$\tau = s_0 0 s_0 0 s_1 0 s_0 0 s_2 1 s_2 1 \quad (1)$$

1.2

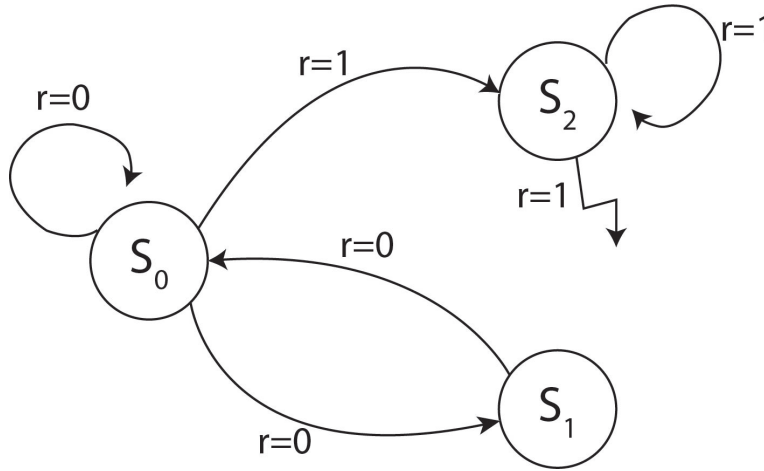


Figure 1: MDP graph of trace τ

Figure one is a minimal MDP graph drawn based on trace τ , three observed states: s_0 , s_1 and s_2 are first drawn. The transitions and their respective rewards are then added to the graph, as it is known that they are possible transitions based on the trace. It has been shown in the end of the trace that there is a reward corresponding with a transition out from s_2 , it is however not shown which state is the transition leading to, there is therefore a path with reward $r = 1$ out from s_2 but not going into any observed states.

- (a) 1. The transition matrix is not deterministic, states can transit to themselves, s_0 can go to s_0, s_1, s_2 , while s_1 has only been shown to be able to go to s_0 , and s_2 has only been shown to go to itself and an unknown state. None of s_0 , s_1 , nor s_2 are terminal states.
2. All of the reward shown has been related to s_2 , either by going into s_2 , out from s_2 or staying in s_2 .
- (b) Most of the reward shown is related to transitions in or out from s_2 , it is therefore expected to be the highest value state from what can be observed. s_0 however is also valuable (positive) as it is shown to be able to transit to s_2 , despite there is no reward corresponding to transitions to itself or with s_1 . Since the trace is not a complete episode, nor is it in a fully known environment with known transition probabilities or rewards, the value of each state can only be estimated

using the temporal difference technique (TD). With a single iteration through the trace, assuming the learning rate $\alpha = 0.5$ the value of s_0 is 0 as the trace has not visited s_0 after the first transition with positive reward. However if the trace is assumed to be traversed twice, the TD algorithm has shown that the value of s_0 has been increased to 1. The code implementation of the algorithm can be found in the appendix.

2 Understanding of Grid Worlds

2.1

The personalised p and γ is:

$$\begin{aligned} p &= 0.45 \\ \gamma &= 0.35 \end{aligned}$$

2.2

Both the optimal value function and the optimal policy is computed by value iteration with dynamic programming, as the transition probability and reward of every state is known (given a perfect model of the environment as a Markov Decision Process). The algorithm is as follows:

1. Initialise all state values as 0
2. Set the tolerance, θ to a small value (0.01)
3. For each state, s in the environment:
 1. Record the current state value, v
 2. Compute the new state value, $V(s)$:

$$V(s) = \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')] \quad (2)$$

3. Repeat 3. until the difference magnitude $|v - V(s)|$ is smaller than θ

The optimal state values and policy computed using a Python implementation is shown in Figure 2.

2.3

The optimal policy execute in state s_9 is West (towards s_8) and $p(a, s_9)$ is:

$$p(a, s_9) = \{N : 0.183, W : 0.45, S : 0.183, E : 0.183\} \quad (3)$$

West is a sensible choice of action at s_9 as s_8 is the state with highest value state directly accessible from s_9 , since it is closer to the only reward state s_2 than other directly accessible state. As long as γ is greater than 0 (values future reward) and p is greater than 0.25 (probability of successful desired action) both the optimal action and $p(a, s_9)$ will remain the same. Even if the agent values future reward very slightly, the s_8 will always be the state with highest value accessible form s_9 , and the probability of executing a desired action will always be the highest out of the four possible actions when $p \geq 0.25$.

2.4

One observation from the optimal value and policy is that none of the optimal policy executes are South, this is due to the fact that the penalty state s_{11} is located South relative to every other state, the state optimal policy execute will therefore avoid South to reach a state of higher value. It should also be pointed out that states closer the the reward state s_2 have higher values, with s_1 being the highest at 4.53, slightly higher than the other states that are also directly adjacent to s_2 , this is due to the non-deterministic execution of optimal policy, s_1 is surrounded by 2 walls, which means s_1 will arrive back at itself if the actual execution of the policy results in a movement towards the wall, allowing another chance to move into the desired state. The difference in values between adjacent states to s_2 is expected to decrease as p increases (increasingly deterministic execution of policy)

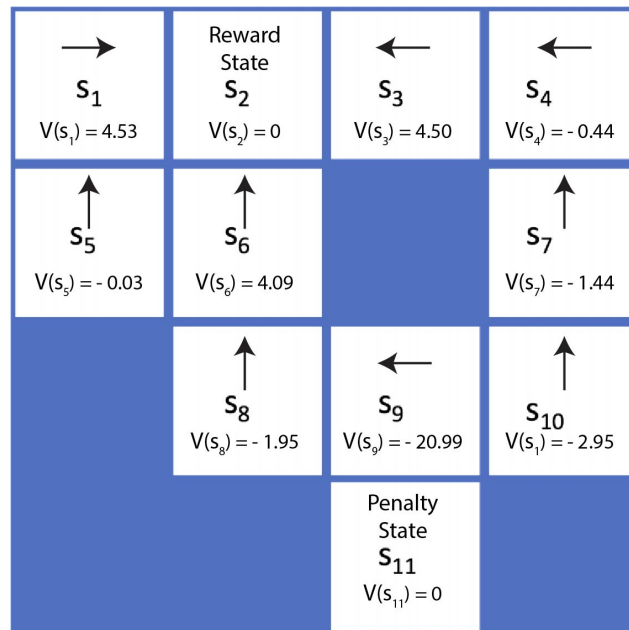


Figure 2: The personalised grid world for CID:01184493. Where s_2 is the reward state and s_{11} is the penalty state. The optimal value of each state is below each state name with the optimal policy at each state above the state name

3 Appendix: Code implementation

3.1 Understanding of MDPs

```
1 states = {0, 1, 2} # observed states
2
3 trace = [0, 0, 1, 0, 0, 2, 2] # personalised trace
4 reward = [0, 0, 0, 0, 0, 1, 1] # personalised reward
5
6
7 # temporal difference algorithm
8 def td_estimation(t, r, gamma, alpha):
9
10     V = {0: 0, 1: 0, 2: 2}
11
12     for _ in range(1):
13         for step in range(len(t) - 1):
14             delta = r[step] + gamma * V[t[step + 1]] - V[t[step]]
15             V[t[step]] += alpha * delta
16
17     return V
18
19
20 print(td_estimation(trace, reward, 1, 0.5))
```

3.2 Understanding of Grid Worlds

```

1 import numpy as np
2
3 num_rows = 4
4 num_cols = 4
5 reward_states = [(0, 1)] # define reward state position
6 penalty_states = [(3, 2)] # define penalty state poistion
7
8 walls = [(1, 2), (2, 0), (3, 0), (3, 1), (3, 3)] # define walls position
9
10 # The environment
11 class Game:
12     def __init__(self, rows, cols, walls, reward_states, penalty_states,
13                 start_pos, action_success_prob):
14         self.rows = rows
15         self.cols = cols
16
17         self.board = np.ndarray((rows, cols), dtype=object)
18         for wall in walls:
19             self.board[wall] = -1 # mark walls (-1)
20
21         for reward_state in reward_states:
22             self.board[reward_state] = 1 # mark reward state (1)
23
24         for penalty_state in penalty_states:
25             self.board[penalty_state] = 2 # mark penalty state (2)
26
27         state_num = 1
28
29         # replace markings with state state objects
30         for i in range(self.rows):
31             for j in range(self.cols):
32                 if not self.board[i, j]: # normal states are unmarked
33                     self.board[i, j] = Position(state_num, -1, is_state=
34                                             True)
35                     state_num += 1
36                 elif self.board[i, j] == -1:
37                     self.board[i, j] = Position(None, None, is_wall=True)
38                 elif self.board[i, j] == 1:
39                     self.board[i, j] = Position(state_num, 10, is_state=
40                                             True, is_terminal=True)
41                     state_num += 1
42                 elif self.board[i, j] == 2:
43                     self.board[i, j] = Position(state_num, -100, is_state
44                                             =True, is_terminal=True)
45                     state_num += 1
46
47         self.actions = ["N", "S", "W", "E"] # define possible actions
48         self.action_success_prob = action_success_prob # define p (
49             personalised p from CID)
50
51         # give a starting position if needed (not used in value
52             iteration)
53         if start_pos:

```

```

48         self.current_pos = start_pos
49
50         # score of the current episode (not used)
51         self.score = 0
52
53         # returns the next position for a given current position and
successfully executed action)
54         def get_next_position(self, from_pos, action):
55
56             if action == "N":
57                 next_pos = (from_pos[0] - 1, from_pos[1])
58             elif action == "S":
59                 next_pos = (from_pos[0] + 1, from_pos[1])
60             elif action == "W":
61                 next_pos = (from_pos[0], from_pos[1] - 1)
62             else:
63                 next_pos = (from_pos[0], from_pos[1] + 1)
64
65             # check if next position is legal
66             if (next_pos[0] >= 0) and (next_pos[0] <= 3):
67                 if (next_pos[1] >= 0) and (next_pos[1] <= 3):
68                     if self.board[next_pos].is_state:
69                         return next_pos
70
71             return from_pos # returns current position if move is not legal
(i.e move into a wall)
72
73         # returns the actual executed action (non-deterministic policy
execution)
74         def get_next_action(self, intended_action):
75             success = np.random.choice([1, 0], p=[self.action_success_prob, 1
76                 - self.action_success_prob])
77             if success:
78                 return intended_action
79             else:
80                 remaining_actions = self.actions.remove(intended_action)
81                 return np.random.choice(remaining_actions)
82
83         # returns the transition reward for moving into a certain next state
84         def get_transition_reward(self, next_pos):
85             return self.board[next_pos].reward
86
87         # returns a transition probability for value iteration
88         def get_transition_prob(self, intended_action, actual_action):
89             if intended_action == actual_action:
90                 return self.action_success_prob
91             else:
92                 return (1 - self.action_success_prob) / (len(self.actions) -
93                     1)
94
95         # show the grid world
96         def show_board(self):
97             for i in range(self.rows):
98                 print('-----')
99                 out = '| '

```

```

98         for j in range(self.cols):
99             box = str(self.board[i, j])
100             box = box.ljust(3, ' ')
101             out += box + ' | '
102         print(out)
103     print('-----')
104
105
106 # State, Wall or Terminal states
107 class Position:
108     def __init__(self, state_num, reward, is_state=False, is_terminal=
109         False, is_wall=False):
110
111         self.is_state = is_state
112         self.is_terminal = is_terminal
113         self.is_wall = is_wall
114
115         if self.is_terminal:
116             self.name = 't' + str(state_num)
117
118         elif self.is_state:
119             self.name = 's' + str(state_num)
120
121         self.reward = reward
122
123     def __str__(self):
124         if self.is_state:
125             return self.name
126         elif self.is_wall:
127             return 'X'
128
129 # the agent for value iteration or other algorithms (not implemented)
130 class Agent:
131     def __init__(self, gamma, alpha, environment):
132         self.gamma = gamma # personalised gamma
133         self.alpha = alpha # learning rate (not used)
134
135         self.environment = environment # load into the created grid
136             world
137         self.actions = self.environment.actions # ["N", "S", "W", "E"]
138
139         # initialise state values as 0
140         self.state_values = {}
141         for i in range(self.environment.rows):
142             for j in range(self.environment.cols):
143                 self.state_values[(i, j)] = 0
144
145         # initialise policy as North
146         self.policy = {}
147         for i in range(self.environment.rows):
148             for j in range(self.environment.cols):
149                 self.policy[(i, j)] = 'N'
150
151         # value iteration algorithm

```



```

151     def value_iteration(self):
152         tolerance = 0.01 # tolerance
153         delta = 1000 # initialise difference to be a large value
154
155         while delta > tolerance:
156             delta = 0
157
158             # for each state in the grid if it is a state and not a
159             terminal state
160             for i in range(self.environment.rows):
161                 for j in range(self.environment.cols):
162                     if self.environment.board[(i, j)].is_state and not
163                         self.environment.board[(i, j)].is_terminal:
164
165                         old_max_value = self.state_values[(i, j)] # v
166
167                         this_state_values = [] # possible value of the
168                             current state for different action
169
170                         for ia in self.actions: # intended action
171
172                             sum_action_values = 0
173
174                             for aa in self.actions: # actual action (for
175                                 different possible next states)
176
177                                 next_pos = self.environment.
178                                     get_next_position((i, j), aa) # next
179                                     position
180                                 trans_prob = self.environment.
181                                     get_transition_prob(ia, aa) #
182                                     transition probability
183                                 trans_reward = self.environment.
184                                     get_transition_reward(next_pos) #
185                                     transition reward
186
187                                 sum_action_values += trans_prob * (
188                                     trans_reward + self.gamma * self.
189                                     state_values[next_pos])
190
191                                 this_state_values.append(sum_action_values)
192
193             self.state_values[(i, j)] = max(this_state_values)
194             # update state value
195             self.policy[(i, j)] = self.actions[np.argmax(
196                 this_state_values)] # update policy
197
198             delta = max(delta, abs(old_max_value - self.
199                 state_values[(i, j)])) # calculate
200                 difference
201
202             # self.show_values()
203
204     # show state values
205     def show_values(self):

```

```

190         for i in range(self.environment.rows):
191             print('-----')
192             out = '| '
193             for j in range(self.environment.cols):
194                 if self.environment.board[(i, j)].is_terminal:
195                     box = 'T'
196                 elif self.environment.board[(i, j)].is_wall:
197                     box = 'X'
198                 else:
199                     box = "{:.2f}".format(self.state_values[i, j])
200
201                 box = box.ljust(7, ' ')
202                 out += box + '| '
203             print(out)
204         print('-----')
205
206     # show optimal policy
207     def show_policy(self):
208         for i in range(self.environment.rows):
209             print('-----')
210             out = '| '
211             for j in range(self.environment.cols):
212                 if self.environment.board[(i, j)].is_terminal:
213                     box = 'T'
214                 elif self.environment.board[(i, j)].is_wall:
215                     box = 'X'
216                 else:
217                     box = self.policy[(i, j)]
218
219                 box = box.ljust(3, ' ')
220                 out += box + '| '
221             print(out)
222         print('-----')
223
224
225 gridworld = Game(num_rows, num_cols, walls, reward_states, penalty_states
226                 , None, 0.45) # initiate gridworld
227 gridworld.show_board() # show board
228 agent = Agent(0.35, None, gridworld) # initiate the agent
229 agent.value_iteration() # value iteration
230 agent.show_values() # show state values corresponding to board position
231 agent.show_policy() # show optimal policy corresponding to board
232                        position

```