

지능시스템

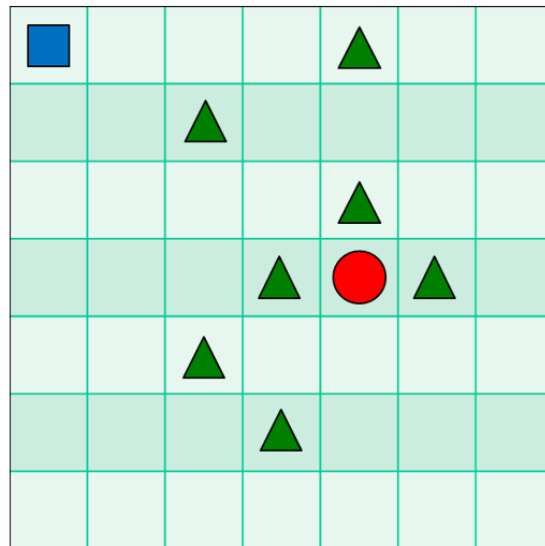
과제5

2019305059

이현수

숙제 5 (First Visit Monte Carlo Method)

- 다음의 7×7 Grid Map에서 출발점, 목표지점, 장애물이 각각 표시된 위치에 있을 때, 목표지점을 찾아가는 행동을 **First Visit Monte Carlo 방법**으로 구하고 학습 진행(epoch)에 따른 결과(학습된 각 상태의 상태함수값 및 정책에 의한 행동 확률을 격자 위에 표시)를 제시하시오. (제출: 5월 29일 6시까지, Google Classroom)



Rewards

- ▲ 상태로 가는 행동: -1
- 상태로 가는 행동: $+1$



그리드 월드에서 좌표값은 위 사진과 같다.

● environment.py 수정코드

```
6 np.random.seed(1)
7 PhotoImage = ImageTk.PhotoImage
8 UNIT = 100 # 픽셀 수
9 HEIGHT = 7 # 그리드 월드 세로
10 WIDTH = 7 # 그리드 월드 가로
```

그리드 월드는 5*5에서 7*7로 수정한다.

```
36 # 캔버스에 이미지 추가
37 self.rectangle = canvas.create_image(50, 50, image=self.shapes[0])
38 self.triangle1 = canvas.create_image(450, 50, image=self.shapes[1])
39 self.triangle2 = canvas.create_image(250, 150, image=self.shapes[1])
40 self.triangle3 = canvas.create_image(450, 250, image=self.shapes[1])
41 self.triangle4 = canvas.create_image(350, 350, image=self.shapes[1])
42 self.triangle5 = canvas.create_image(550, 350, image=self.shapes[1])
43 self.triangle6 = canvas.create_image(250, 450, image=self.shapes[1])
44 self.triangle7 = canvas.create_image(350, 550, image=self.shapes[1])
45 self.circle = canvas.create_image(450, 350, image=self.shapes[2])
```

캔버스에 문제와 같이 추가한다.

```
98 # 보상 함수
99 if next_state == self.canvas.coords(self.circle):
100     reward = 100
101     done = True
102 elif next_state in [self.canvas.coords(self.triangle1),
103                     self.canvas.coords(self.triangle2),
104                     self.canvas.coords(self.triangle3),
105                     self.canvas.coords(self.triangle4),
106                     self.canvas.coords(self.triangle5),
107                     self.canvas.coords(self.triangle6),
108                     self.canvas.coords(self.triangle7)]:
109     reward = -100
110     done = True
111 else:
112     reward = 0
113     done = False
```

102~108줄 코드와 같이 triangle을 추가시킨다.

```

120     def text_value(self, row, col, contents, action, font='Helvetica', size=10,
121                   style='normal', anchor="nw"):
122         if action == 0:
123             origin_x, origin_y = 7, 42
124         elif action == 1:
125             origin_x, origin_y = 85, 42
126         elif action == 2:
127             origin_x, origin_y = 42, 5
128         elif action == 3:
129             origin_x, origin_y = 42, 77
130         else:
131             origin_x, origin_y = 42, 42
132
133         x, y = origin_y + (UNIT * col), origin_x + (UNIT * row)
134         font = (font, str(size), style)
135         text = self.canvas.create_text(x, y, fill="black", text=contents,
136                                       font=font, anchor=anchor)
137         return self.texts.append(text)
138

```

```

139     def print_value_q_all(self, q_table, value_table):
140         for i in self.texts:
141             self.canvas.delete(i)
142         self.texts.clear()
143         for x in range(HEIGHT):
144             for y in range(WIDTH):
145                 state=[x,y]
146                 if str(state) in value_table.keys():
147                     temp=value_table[str(state)]
148                     self.text_value(y,x,round(temp,2),4)
149                 for action in range(0,4):
150                     if str(state) in q_table.keys():
151                         temp=q_table[str(state)][action]
152                         self.text_value(y,x,round(temp,2),action)

```

각 칸에 상태함수값과 행동 확률을 적기위한 함수들이다.

● mc_agent.py 수정코드

```
7 # 몬테카를로 에이전트 (모든 에피소드 각각의 샘플로 부터 학습)
8 class MCAgent:
9     def __init__(self, actions):
10         self.width = 7
11         self.height = 7
12         self.actions = actions
13         self.learning_rate = 0.01
14         self.discount_factor = 0.9
15         self.epsilon = 0.1
16         self.samples = []
17         self.value_table = defaultdict(float)
18         self.percentage_table = defaultdict(lambda: [0.0, 0.0, 0.0, 0.0])
19         self.reward_table = defaultdict(lambda: [0.0, 0.0, 0.0, 0.0])
```

Width=7, height=7로 업데이트 한다.

```
25 # 모든 에피소드에서 에이전트가 방문한 상태의 큐 함수를 업데이트
26 def update(self):
27     visit_state=[]
28     sample_size=len(self.samples)
29     for i, forward in enumerate(self.samples):
30         state, reward, done = forward
31         state = str(state)
32         if state not in visit_state:
33             visit_state.append(state)
34             G_t = 0
35             for j, reverse in enumerate(reversed(self.samples)):
36                 if sample_size-j-1<i:
37                     break
38                 rev_state, rev_reward, rev_done = reverse
39                 rev_state=str(rev_state)
40                 G_t=rev_reward+self.discount_factor*G_t
41                 value = self.value_table[state]
42                 self.value_table[state]=(value+self.learning_rate*(G_t-value))
```

```
70 # 가능한 다음 모든 상태들을 반환
71 def possible_next_state(self, state):
72     col, row = state
73     next_state = [0.0] * 4
74
75     if row!=0:
76         next_state[0]=self.reward_table[str(state)][0]+self.discount_factor+self.value_table[str([col,row-1])]
77     else:
78         next_state[0] = self.reward_table[str(state)][0] + self.discount_factor + self.value_table[str(state)]
79     if row!=self.height-1:
80         next_state[1] = self.reward_table[str(state)][1] + self.discount_factor + self.value_table[str([col, row + 1])]
81     else:
82         next_state[0] = self.reward_table[str(state)][1] + self.discount_factor + self.value_table[str(state)]
83     if col!=0:
84         next_state[2]=self.reward_table[str(state)][2]+self.discount_factor+self.value_table[str([col-1,row])]
85     else:
86         next_state[2]=self.reward_table[str(state)][2]+self.discount_factor+self.value_table[str(state)]
87     if col!=self.width-1:
88         next_state[3]=self.reward_table[str(state)][3]+self.discount_factor+self.value_table[str([col+1,row])]
89     else:
90         next_state[3]=self.reward_table[str(state)][3]+self.discount_factor+self.value_table[str(state)]
91
```

```

94     def reward_state_action(self, state, action, reward):
95         state = str(state)
96         self.reward_table[state][action]=reward
97

```

```

98     def update_percentage_table(self, state, next_state):
99         value = -99999
100         max_index = []
101
102         result = [0.0, 0.0, 0.0, 0.0]
103
104         # 모든 행동에 대해서 [보상 + (감가율 * 다음 상태 가치함수)] 계산
105         for index, action in enumerate([0, 1, 2, 3]):
106
107             reward = self.reward_table[str(state)][action]
108             next_value = self.value_table[str(next_state)]
109             temp = reward + self.discount_factor * next_value
110
111             # 받을 보상이 최대인 행동의 index(최대가 복수라면 모두)를 추출
112             if temp == value:
113                 max_index.append(index)
114             elif temp > value:
115                 value = temp
116                 max_index.clear()
117                 max_index.append(index)
118
119         # 행동의 확률 계산
120         prob = 1 / len(max_index)
121         for index in max_index:
122             result[index] = prob
123
124         self.percentage_table[str(state)][0] = result[0]
125         self.percentage_table[str(state)][1] = result[1]
126         self.percentage_table[str(state)][2] = result[2]
127         self.percentage_table[str(state)][3] = result[3]
128

```

행동에 대한 확률을 나타내기 위한 코드이다.










```

130     # 메인 함수
131     if __name__ == "__main__":
132         env = Env()
133         agent = MCAgent(actions=list(range(env.n_actions)))
134
135         for episode in range(10000):
136             state = env.reset()
137             action = agent.get_action(state)
138
139             while True:
140                 env.render()
141
142                 # 다음 상태로 이동
143                 # 보상은 숫자이고, 완료 여부는 boolean
144                 now_state, next_state, reward, done = env.step(action)
145                 agent.save_sample(next_state, reward, done)
146                 agent.reward_state_action(next_state, action, reward)
147                 # 다음 행동 받아옴
148                 action = agent.get_action(next_state) ✚
149
150                 agent.reward_state_action(now_state, action, reward)
151                 agent.update_percentage_table(now_state, next_state)
152
153                 env.print_value_q_all(agent.percentage_table, agent.value_table)
154                 # 에피소드가 완료됐을 때, 큐 함수 업데이트
155                 if done:
156                     agent.update()
157                     agent.samples.clear()
158                     break

```

메인함수 코드이다.

●실행결과

monte carlo						
0.25	0.25	0.25	0.33		0.0	
0.25 -0.52 0.25	0.25 -0.58 0.25	0.25 -0.81 0.25	0.0 -0.9 0.33			
0.25	0.25	0.25	0.33			
0.25	0.33		0.25	0.0		
0.25 -0.67 0.25	0.33 -0.66 0.0		0.25 -0.53 0.25	0.33 -0.73 0.33		
0.25	0.33		0.25	0.33		
0.25	0.25	0.25	0.25		0.0	
0.25 -0.11 0.25	0.25 -0.18 0.25	0.25 -0.43 0.25	0.25 -0.48 0.25			
0.25	0.25	0.25	0.25			
0.25	0.25	0.25				
0.25 -0.25 0.25	0.25 -0.3 0.25	0.25 -0.39 0.25				
0.25	0.25	0.25				
0.25	0.25		0.0			
0.25 -0.36 0.25	0.25 -0.33 0.25					
0.25	0.25					
0.25	0.25	0.0				
0.25 -0.97 0.25	0.25 -0.83 0.25	0.5 -1.79 0.5		0.0		
0.25	0.25	0.0				
0.25	0.25	0.25				
0.25 -0.43 0.25	0.25 -0.73 0.25	0.25 -0.81 0.25	0.0			
0.25	0.25	0.25				

원하는 결과가 나오지 못함. 완벽하게 못함.