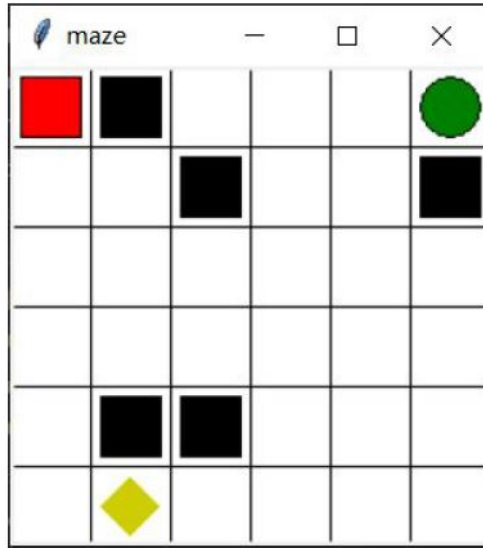


I. INTRODUCTION

A. Purpose

The goal of this homework is to implement a Dyna-Q learning agent to search for the treasure and exit in a grid-shaped maze. The agent will learn by trail and error from interactions with the environment and finally acquire a policy to get as high as possible scores in the game.

Suppose a 6×6 grid-shaped maze in Figure 1. The red rectangle represents the start point and the green circle represents the exit point. You can move upward, downward, leftward and rightward and you should avoid falling into the traps, which are represented by the black rectangles. Finding the exit will give a reward +1 and falling into traps will cause a reward -1, and both of the two cases will terminate current iteration. You will get a bonus reward +3 if you find the treasure, which shown as golden diamond.



(a) The board of Game

B. Environment

There is a minimal amount of equipment to be used in this homework. A few requirements are listed below:

- Python 3.7
- Library: matplotlib,numpy,pandas
- Compiler: Vscode

C. Procedure

The key theories of the algorithm are consist of three parts as following:

1. Dyna-Q learning
2. Epsilon Greedy

II. IMPLEMENTATION

A. Dyna-Q learning

Compared with normal Q learning, Dyna-Q learning in each exploring step that it uses an experience model to generate sub-optimal policies, which speeds up the learning process.

Dyna-Q Learning is composed of model-free Q-Learning and model-based Dyna-Q Learning. Model-free part updates Q-value simply based on its expensive experience, while model-based part updates Q-value when backpropagating the experienced fragment of episode.

And below is Dyna-Q algorithm's pseudo-code.

Tabular Dyna-Q
Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$
Loop forever:
 (a) $S \leftarrow$ current (nonterminal) state
 (b) $A \leftarrow \epsilon$ -greedy(S, Q)
 (c) Take action A ; observe resultant reward, R , and state, S'
 (d) $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$ Direct RL
 (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment) model learning
 (f) Loop repeat n times: planning
 $S \leftarrow$ random previously observed state
 $A \leftarrow$ random action previously taken in S
 $R, S' \leftarrow Model(S, A)$
 $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

知乎 @张文

(b) pseudo-code of Dyna-Q

In this algorithm, I use the normal Dyna-Q algorithm framework, store the data in the tuple, and constantly refresh the state to update the Q value to achieve better score.

B. Epsilon Greedy

The epsilon-greedy algorithm, as a strategy for decision-making, is widely used in many fields of machine learning. For example, when $\epsilon = 0.6$, it means that 60% of the time, the behavior will be selected according to the optimal value of the Q table, and the behavior will be randomly selected 40% of the time.

Using epsilon-greedy algorithm can prevent Dyna-Q algorithm from falling into local optimum to a certain extent, and can help the program to make early decisions.

The most important thing for the epsilon-greedy algorithm is the epsilon attenuation strategy. After trying a variety of different methods and strategies, I finally chose the e-exponential decay strategy, which decreases rapidly at the beginning and slowly decreases afterwards to obtain better results.

C. Parameter Adjustment

In the field of Artificial Intelligence, the selection of suitable parameters can play a great role in a certain sense. In the algorithm used in this experiment, the main adjustable parameters are learning rate alpha, discount factor gamma, and random coefficient epsilon.

1. learning rate alpha

The learning rate affects how much the code adjusts the weights of the network and adjusts the loss gradient. The lower the learning rate, the slower the downward slope. A lower learning rate can ensure that no local minimum is missed, but it also means that it will take a long time to converge.

2. discount factor gamma

The value of the discount factor determines how we view future rewards. A gamma close to 0 means that future rewards are not valued, and a gamma close to 1 means that future rewards are valued and a long-term vision.

3. epsilon

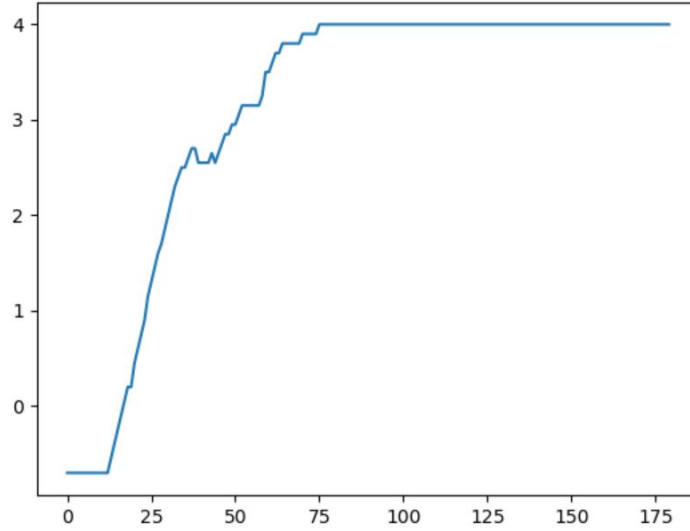
epsilon and its attenuation strategy determine how long the code can find the optimal solution and the ability to jump out of the local optimal solution.

Finally I choose the e exponential decay strategy:

$$\varepsilon = \frac{\varepsilon_0}{e^{\frac{n}{35}}}$$

D. Result

The final experimental results are shown in the figure below, the ordinate is the average of ten consecutive scores.



(c) result score

In the figure, we can see that the score basically continued to rise before 75 times, and converged to full marks after 75 times.

III. CONCLUSION

I designed a agent which's named as My_Agent to participate this tournament. From the test in our PC, I found my agent could complete the mission around 75 times.

Now is the ending time, first I want to thank for the help from teaching assistants and our professor Yue Gao. Through this homework, we have attained lots of skills and experiences. We have almost mastered the Dyna-Q learning. We gain a lot.

I have added an adjunct of source codes with this test report, Thanks again!

IV. EXPERIMENT CODE

This section contains my Agent code block of this homework.

My_Agent:

```
1 import numpy as np
import pandas as pd
3 import random
import math
5
class My_Agent:
7     def __init__(self, actions):
            self.Qvalue_table = {}
            self.actions = actions
9
11    def epsilon_greedy(self, state, epsilon, episode):
            for action in self.actions:
13                if (tuple(state), action) not in self.Qvalue_table:
                    self.Qvalue_table[(tuple(state), action)] = 0

15            epsilon /= math.exp(episode/35)    #change with the episode

17            if random.random() < epsilon:    #case random occur
                action = np.random.choice(self.actions)
19            else:
                # get the action corresponding to max q
21                maxval_list = []
                max_val = -float('inf')
23                for action in self.actions:
                    val = self.Qvalue_table[(tuple(state), action)]
25                    if (max_val < val):
                        max_val = val
27                for action in self.actions:
                    if (self.Qvalue_table[(tuple(state), action)] == max_val):
29                    maxval_list.append(action)
                    action = np.random.choice(maxval_list)
31            return action

33    def Qvalue_refresh(self, new_action, state, new_state, reward, gamma, alpha):
            max_tmp = -float('inf')
35            for action in self.actions:
                step_Qvalue = (tuple(new_state), action)
37                if step_Qvalue not in self.Qvalue_table:
                    self.Qvalue_table[step_Qvalue] = 0
39                if max_tmp < self.Qvalue_table[step_Qvalue]:
                    max_tmp = self.Qvalue_table[step_Qvalue]
41
            self.Qvalue_table[(tuple(state), new_action)] = (1 - alpha) * self.Qvalue_table
                [(tuple(state), new_action)] + alpha * (reward + gamma * max_tmp)
```

Main:

```
1 from maze_env import Maze
from RL_brain import My_Agent
3 import random
from time import sleep
5 import numpy as np
import matplotlib.pyplot as plt
```

```

7
9 if __name__ == "__main__":
    ### START CODE HERE ###
11     env = Maze()
    agent = My_Agent(list(range(env.n_actions)))
13
    epsilon=0.6
15     gamma=0.85
    alpha=0.3
17     training_count=200
19
    training_reward=np.arange(0,training_count)
    reward_plot=np.zeros(training_count)
21     for episode in range(training_count):
        state = env.reset()
23         episode_reward = 0
        while True:
25             #env.render()
            action = agent.epsilon_greedy(state, epsilon, episode)
27             new_state, reward, done = env.step(action)
            episode_reward += reward
29             agent.Qvalue_refresh(action, state, new_state, reward, gamma, alpha)
            state = new_state
31             if done:
                #env.render()
33                 sleep(0.001)
                break
35
            print('episode:', episode, 'episode_reward:', episode_reward)
37             training_reward[episode] = episode_reward
39
        for i in range(training_count):
            for j in range(0,20):
41                 reward_plot[i]+=training_reward[i-j]
43
45     plt.plot(reward_plot[20:]/20)
    plt.show()
    ### END CODE HERE ###
47     print('\ntraining_over\n')

```

maze_env:

```

1 import numpy as np
  np.random.seed(1)
3 import tkinter as tk
  import time
5
7 UNIT = 40
  MAZE_H = 6
9 MAZE_W = 6
11
13 class Maze(tk.Tk, object):
    def __init__(self):
        super(Maze, self).__init__()

```

```

15     self.action_space = ['u', 'd', 'l', 'r']
16     self.n_actions = len(self.action_space)
17     self.title('maze')
18     self.geometry('{0}x{1}'.format(MAZE_H * UNIT, MAZE_H * UNIT))
19     self._build_maze()
20     self.bonusFlag = False
21     self.eating = True

22
23     def _build_maze(self):
24         self.canvas = tk.Canvas(self, bg='white',
25                                height=MAZE_H * UNIT,
26                                width=MAZE_W * UNIT)

27
28         for c in range(0, MAZE_W * UNIT, UNIT):
29             x0, y0, x1, y1 = c, 0, c, MAZE_H * UNIT
30             self.canvas.create_line(x0, y0, x1, y1)
31         for r in range(0, MAZE_H * UNIT, UNIT):
32             x0, y0, x1, y1 = 0, r, MAZE_H * UNIT, r
33             self.canvas.create_line(x0, y0, x1, y1)

34
35         origin = np.array([20, 20])

36
37         # hell1
38         hell1_center = origin + np.array([UNIT * 1, UNIT * 0])
39         self.hell1 = self.canvas.create_rectangle(
40             hell1_center[0] - 15, hell1_center[1] - 15,
41             hell1_center[0] + 15, hell1_center[1] + 15,
42             fill='black')
43
44         # hell2
45         hell2_center = origin + np.array([UNIT * 2, UNIT * 1])
46         self.hell2 = self.canvas.create_rectangle(
47             hell2_center[0] - 15, hell2_center[1] - 15,
48             hell2_center[0] + 15, hell2_center[1] + 15,
49             fill='black')
50
51         # hell3
52         hell3_center = origin + np.array([UNIT * 5, UNIT * 1])
53         self.hell3 = self.canvas.create_rectangle(
54             hell3_center[0] - 15, hell3_center[1] - 15,
55             hell3_center[0] + 15, hell3_center[1] + 15,
56             fill='black')
57
58         # hell4
59         hell4_center = origin + np.array([UNIT * 1, UNIT * 4])
60         self.hell4 = self.canvas.create_rectangle(
61             hell4_center[0] - 15, hell4_center[1] - 15,
62             hell4_center[0] + 15, hell4_center[1] + 15,
63             fill='black')
64
65         # hell5
66         hell5_center = origin + np.array([UNIT * 2, UNIT * 4])
67         self.hell5 = self.canvas.create_rectangle(
68             hell5_center[0] - 15, hell5_center[1] - 15,
69             hell5_center[0] + 15, hell5_center[1] + 15,
70             fill='black')
71
72         oval_center = origin + np.array([UNIT * 5, UNIT * 0])
73         self.oval = self.canvas.create_oval(
74             oval_center[0] - 15, oval_center[1] - 15,
75             oval_center[0] + 15, oval_center[1] + 15,
76             fill='green')

```

```

75         self.rect = self.canvas.create_rectangle(
            origin[0] - 15, origin[1] - 15,
            origin[0] + 15, origin[1] + 15,
77             fill='red')

79         # create bonus
        bonus_center = origin + np.array([UNIT * 1, UNIT * 5])
81         self.bonus = self.canvas.create_polygon(
            [bonus_center[0]+15, bonus_center[1],
83             bonus_center[0], bonus_center[1]-15,
            bonus_center[0]-15, bonus_center[1],
85             bonus_center[0], bonus_center[1]+15],
            fill='#CD0000')
87         self.bonus_location = self.canvas.create_polygon(
            bonus_center[0] - 15, bonus_center[1] - 15,
89             bonus_center[0] + 15, bonus_center[1] + 15,
            fill='white')
91
92         self.canvas.pack()
93
94         def reset(self):
95             self.canvas.delete(self.rect)
96             self.canvas.delete(self.bonus)
97             origin = np.array([20, 20])
98             self.rect = self.canvas.create_rectangle(
99                 origin[0] - 15, origin[1] - 15,
100                 origin[0] + 15, origin[1] + 15,
101                 fill='red')
102             bonus_center = origin + np.array([UNIT * 1, UNIT * 5])
103             self.bonus = self.canvas.create_polygon(
104                 [bonus_center[0]+15, bonus_center[1],
105                 bonus_center[0], bonus_center[1]-15,
106                 bonus_center[0]-15, bonus_center[1],
107                 bonus_center[0], bonus_center[1]+15],
108                 fill='#CD0000')
109             self.bonusFlag = False
110             self.eating = True
111             return self.canvas.coords(self.rect) + [self.bonusFlag]

112         def step(self, action):
113             s = self.canvas.coords(self.rect)
114             base_action = np.array([0, 0])
115             if action == 0: # up
116                 if s[1] > UNIT:
117                     base_action[1] -= UNIT
118             elif action == 1: # down
119                 if s[1] < (MAZE_H - 1) * UNIT:
120                     base_action[1] += UNIT
121             elif action == 2: # right
122                 if s[0] < (MAZE_W - 1) * UNIT:
123                     base_action[0] += UNIT
124             elif action == 3: # left
125                 if s[0] > UNIT:
126                     base_action[0] -= UNIT
127
128             self.canvas.move(self.rect, base_action[0], base_action[1]) # move agent
129
130             s_ = self.canvas.coords(self.rect) # next state

```



```

133     if s_ == self.canvas.coords(self.oval):
134         reward = 1
135         done = True
136     elif (s_ == self.canvas.coords(self.bonus_location)) and (self.bonusFlag ==
137         False):
138         self.bonusFlag = True
139         reward = 3
140         done = False
141     elif s_ in [self.canvas.coords(self.hell1), self.canvas.coords(self.hell2),
142         self.canvas.coords(self.hell3), self.canvas.coords(self.hell4),
143         self.canvas.coords(self.hell5)]:#, self.canvas.coords(self.hell6), self.canvas
144         .coords(self.hell7)]:
145         reward = -1
146         done = True
147     else:
148         reward = 0
149         done = False
150     s_.append(self.bonusFlag)
151
152     return s_, reward, done
153
154 def render(self):
155     self.update()
156     if self.bonusFlag and self.eating:
157         self.eating = False
158         time.sleep(0.5)
159         self.canvas.delete(self.bonus)

```