

AU 342 PRINCIPLES OF ARTIFICIAL INTELLIGENCE

By: Shen Zhen (518021910149)
Cheng Xinwen (518021910031)

HW#: 2

November 1, 2020

I. INTRODUCTION AND IMPLEMENTATION

In this assignment, we will implement Reinforcement Learning with Dyna-Q in Maze Environment and with Deep Q Network on Atari Game.

A. Reinforcement Learning in Maze Environment

In this assignment, I will 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. This part is finished individually.

1. Game Description

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.

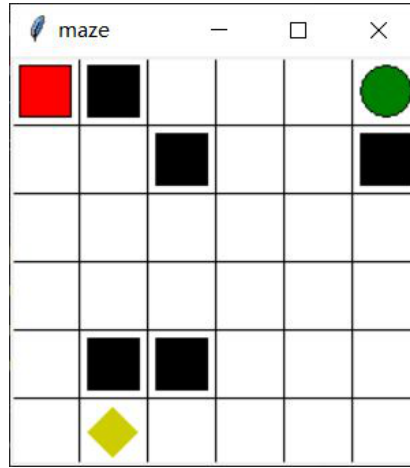


FIG. 1: Maze environment

2. State and Action Space

State(5-dimension): The current position of the agent (4D) and a bool variable (1D) that indicates whether the treasure has been found.

Action(1-dimension): A discrete variable and 0, 1, 2, 3 respectively represent move upward, downward, rightward and leftward.

3. Codes and Explanations

I implemented an agent based on Dyna-Q learning. It is an algorithm different from Q-Learning in that it tries to build the model based on accessible data, which leads to less interaction with the maze environment and needs less data. In general, it is cheaper and saves space. Here is Dyna-Q algorithm shown with pseudocodes.

Tabular Dyna-Q

Initialize $Q(s, a)$ and $Model(s, a)$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}(s)$

Do forever:

- (a) $S \leftarrow$ current (nonterminal) state
- (b) $A \leftarrow \epsilon$ -greedy(S, Q)
- (c) Execute action A ; observe resultant reward, R , and state, S'
- (d) $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$
- (e) $Model(S, A) \leftarrow R, S'$ (assuming deterministic environment)
- (f) Repeat n times:
 - $S \leftarrow$ random previously observed state
 - $A \leftarrow$ random action previously taken in S
 - $R, S' \leftarrow Model(S, A)$
 - $Q(S, A) \leftarrow Q(S, A) + \alpha[R + \gamma \max_a Q(S', a) - Q(S, A)]$

FIG. 2: Dyna-Q

In my code, I choose an N for Dyna-Q, which means that after N actions, the Dyna-Q algorithm will be called and update the Q value table according to existing experience. In my way to calculate the Q value, based on the traditional formula, I added a factor k related to the count of times of which the agent has been to the next state. Also, Temporal-Difference Learning is used. The new formula can be illustrate as following.

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} (Q(s', a') - C(s', a')) - Q(s, a) \right]$$

Here $C(s', a')$ denotes the count of times agent has chosen state s' with action a' . When next state has been rarely been to, the agent tends to choose it. With this formula, the agent can fully explore the whole maze environment. The code is below.

```

1 for j in range(4):
2     self.cnt_of_state[state_][j] = self.cnt_of_state[state_][j] * 0.92
3
4 self.cnt_of_state[state][a] += 1
5 self.update_cnt_of_state(s)
6 self.has_been_to_this_state[float((state[0] + state[2]) / 2), float((state[1] + state[3]) / 2)] = True
7 self.q_dict[state][a] += self.alpha * (
8     r + self.gamma * max(((self.q_dict[state_][i]) - self.cnt_of_state[state_][i] * self.k) for i in range(4))
9     - self.q_dict[state][a])
10 self.model_dict[(state, a)] = [state_, r]
11
12 N = self.Dyna_N
13 for n in range(N): # iteration for N times
14     ms, ma = random.choice(list(self.model_dict))
15     ms_, mr = self.model_dict[(ms, ma)]
16     self.q_dict[ms][ma] += self.alpha * (mr + self.gamma * max(((self.q_dict[ms_][i])

```

Considering the exploration-exploitation dilemma, I completed my Dyna-Q learning agent by implementing ϵ -greedy action selection, meaning it chooses random actions in an epsilon fraction of the time, and follows its current best Q-values otherwise.

Epsilon-Greedy policy selects a random action with probability ϵ or otherwise follows the greedy policy $a = \underset{a}{\operatorname{argmax}} Q^\pi(s, a)$. It can contribute to exploration of the maze environment. However in the real situation, if the probability ϵ is too big, convergence can be a tough problem. Therefore, my plan is to set a major value of ϵ at first. In order to ensure the agent access to every grid of the panel, I estimated whether every grid has been been to. When they are all arrived at, ϵ is set zero and exploitation should take effect. The code is below.

```

2 def epsilon_decay(self, s):
3     '''
4     ---epsilon_is_large_at_first ,then_gradually_decays
5     ---:return:
6     ---'''
7     epsilon = self.epsilon
8     if self.all_arrived == False:
9         epsilon = 0.1
10    if self.all_arrived == True:
11        epsilon = 0
12    return epsilon
13
14 self.epsilon = self.epsilon_decay(state)
15 # use epsilon greedy
16 if random.random() < self.epsilon:
17     action = np.random.choice(self.actions) # a random action
18 else:
19     max_actions=[]
20     for i in range(4):
21         if self.q_dict[state][i] == max(self.q_dict[state]):
22             max_actions.append(i)
23     action = random.choice(max_actions) #choose action of maximum value

```

4. Results and Analysis

In general, my agent is able to finish the game and find the treasure then leave with a total reward of 4. In most cases, my agent can get convergence in 100 steps. However in a few cases, it takes about 200 steps to be converged, I estimated and I think it has something to do with my exploration-exploitation strategy, which can be further optimized.

On the other hand, how to choose the best parameters also matters. The most important parameters in this environment is the ϵ , the discount γ , the learning rate α , and the exploration factor k. ϵ is chosen based on the principal that it should be big at first and then decay to zero. γ is the discount of the contribution of Q values in next steps, which should be a proper value to guarantee convergence. α together with k also should be a proper value. Each parameter is significant in the learning process.

Below is my average learning result, in most cases it can converge in less than 50 episodes.

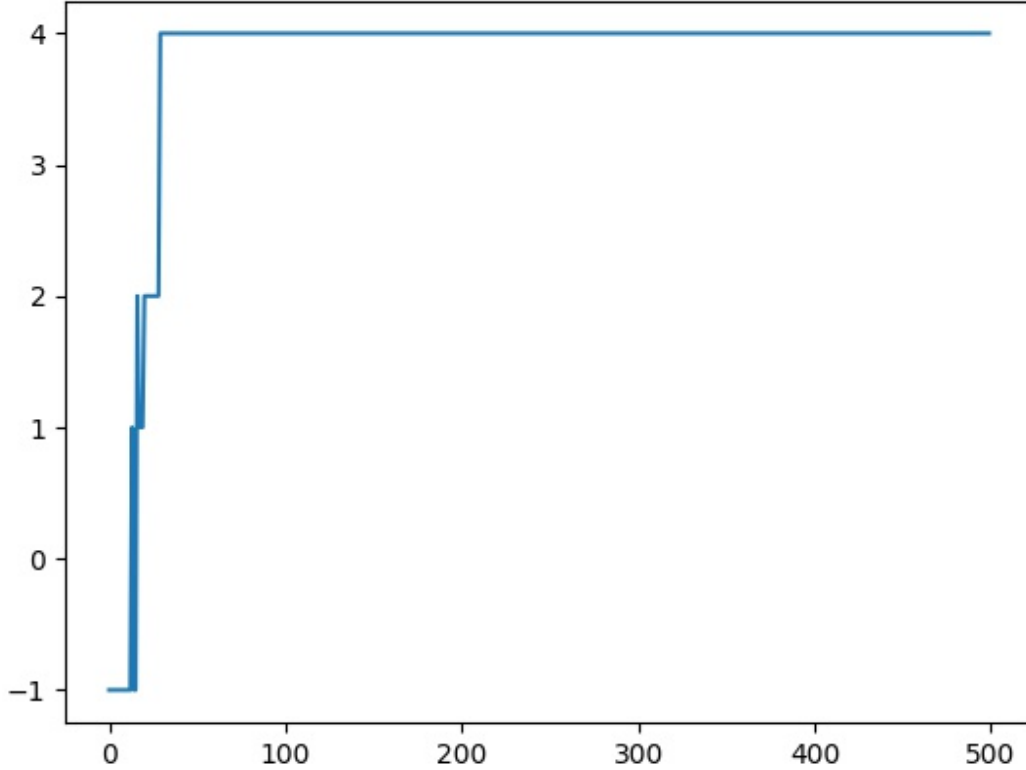


FIG. 3: Q-Learning Result

B. Reinforcement Learning on Atari Game

In this part, we will implement a DQN agent to play the atari game pacman. We first completed the given code and get a simple DQN agent. After tuning the hyper-parameters of the agent using both our experience and suggestions from skillful engineers, we still could not get a satisfying set of parameters. So we decided to improve our DQN algorithm by changing the uniform sample with prioritized experience replay and using Double DQN instead of Nature DQN. Although we still could not reach a satisfying score, in this process, we developed our unique understanding of those hyper-parameters and had a primary knowledge of how their values will influence the convergence and the outcome of learning.

1. Game Description

Pacman is one of the classic and leading games. You need to guide the pac-Man to eat all the dots and avoid the ghosts. In this assignment, you are asked to design a DQN agent to learn control policies directly from the visual information of the game.

As shown in the figure, the ‘MsPacman-ram-v0’ gym environment is utilized as the training environment. This environment provides the ram(128 bytes) of the atari console as model input. Each time, the

agent should choose an action from 9 available actions, corresponding to the 8 buttons on the handle and “do nothing”.

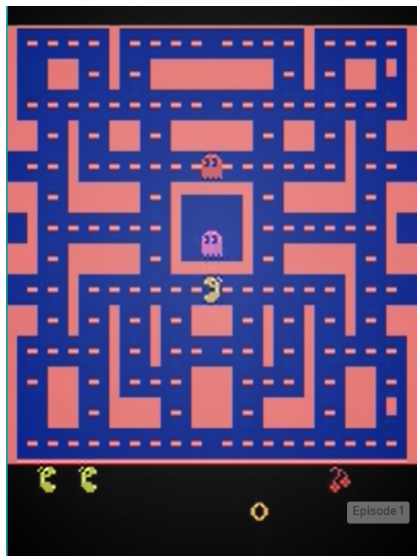


FIG. 4: Atari MsPacman

2. Codes and Explanations

There are already many functions in the given code, so the first step is to understand and use these functions and complete DQN training process. Firstly, we initialized the parameters we would use, and then tuned them according to the training and learning curve. To speed up the learning efficiency and training process, we decided to improve the DQN algorithm first and after that we began to tune the agent. So we will first introduce the primary DQN algorithm frame and next is our improving approach on prioritized experience replay and Double DQN, and finally it comes to our tuning of parameters.

```
# These are hyper parameters for the DQN
2 self.discount_factor = 0.9
  self.learning_rate = 0.01
4 self.epsilon = 1.0
  self.epsilon_min = 0.01
6 self.epsilon_decay = (self.epsilon - self.epsilon_min) / 50000
  self.batch_size = 32
8 self.train_start = 1000
  self.beta = 0.1
10 self.maxlen = 100000
  self.memory = Memory(self.batch_size, self.maxlen, self.beta)
```

Then in each episode, we firstly initialized the score, state and lives. During one game, pacman gets an action using get-action function, which is based on DQN and we will introduce it later. After executing an action, we got next state, reward, done and info as feedback from the environment. Then we judged if the agent was dead, if not the game continues and we save the sample<s, a, r, s'> to replay. To reduce computational overhead, we decided to train the evaluation every 4 steps and updated the target network after 2500 update-times. Update-times here increases only when the agent has taken more than 250 steps under the circumstance of each live. The following is our code of main loop.

```
1 if __name__ == "__main__":
```

```

3  # load the gym env
  env = gym.make('MsPacman-ram-v0')
  # set random seeds to get reproduceable result(recommended)
5  set.random.seed(0)
  # get size of state and action from environment
7  state_size = env.observation_space.shape[0]
  action_size = env.action_space.n
9  # create the agent
  agent = DQNAgent(state_size, action_size)
11 # log the training result
  scores, episodes = [], []
13 graph_episodes = []
  graph_score = []
15 avg_length = 100
  sum_score = 0
17 iteration = 0
  update_times = 0
19
  # train DQN
21 for e in range(EPISODES):
    done = False
    score = 0
    state = env.reset()
    state = np.reshape(state, [1, state_size])
    lives = 3
    while not done:
        dead = False
        step = 0
        while not dead:
            # render the gym env
            if agent.render:
            33         env.render()
            step += 1
            # get action for the current state
            action = agent.get_action(state, step)
            # take the action in the gym env, obtain the next state
            37         next_state, reward, done, info = env.step(action)
            next_state = np.reshape(next_state, [1, state_size])
            # judge if the agent dead
            41         dead = info['ale.lives'] < lives
            lives = info['ale.lives']
            # update score value
            43         score += reward
            # save the sample <s, a, r, s'> to the replay memory
            45         agent.append_sample(state, action, reward, next_state, done)
            # train the evaluation network
            47         if step % 4 == 0:
            49             agent.train_model()
            # go to the next state
            51             state = next_state
            # update the target network after some iterations.
            53             if update_times >= 2500:
            55                 update_times = 0
            57                 agent.eval2target()
            # print info and draw the figure.
            if done:
            59                 scores.append(score)
            sum_score += score
            61                 episodes.append(e)
            # plot the reward each episode
            63                 #pylab.plot(episodes, _scores, _'b')
            print("episode:", e, "score:", score, "memory length:",
            65                 agent.memory.sum_tree.size, "epsilon:", agent.epsilon, "learning rate", agent.learning_rate)
            if e % avg_length == 0:
            67                 graph_episodes.append(e)
            graph_score.append(sum_score / avg_length)
            69                 sum_score = 0
            # plot the reward each avg-length episodes
            71                 pylab.plot(graph_episodes, graph_score, _'r')
            pylab.savefig("./pacman_avg.png")

```

Next we will introduce the prioritized experience replay of the agent. To implement the prioritized experience replay, we first implemented the SumTree data structure. SumTree has four functions: adding a node, updating the priority of the node, getting total priority and sampling according to a given weight. Then we created a class called Memory to implement the prioritized experience replay. In Memory, we can store a transition experienced during the game, get mini batches using prioritized experience replay and update the td-error which is used as the priority of each transition.

In class SumTree, firstly we set the tree capacity. Then we initialized self.tree to a list of 0 and self.data to a list of None. Here we used self.tree to save the priority and self.data to save the transitions. When adding a node, we saved it into self.data and set its priority as the maximum to ensure every node can be sampled at least once. Then we increased the curr-point and self.size. In SumTree, only leaves save the transitions and according priority, other nodes' value is the sum of their children's priorities. When updating the priority of a certain leaf node, we update all the ancestors as well. The following is our code for class SumTree:

```

class SumTree:

```

```

2  def __init__(self, capacity):
4      self.capacity = capacity
      self.tree = [0] * (2 * capacity - 1)
6      self.data = [None] * capacity
      self.size = 0
      self.curr_point = 0
10
12  def add(self, data):
      self.data[self.curr_point] = data
      self.update(self.curr_point, max(self.tree[self.capacity - 1:self.capacity + self.size]) + 1)
14
      self.curr_point += 1
      if self.curr_point >= self.capacity:
          self.curr_point = 0
18
      if self.size < self.capacity:
          self.size += 1
20
22
24  def update(self, point, weight):
      idx = point + self.capacity - 1
      change = weight - self.tree[idx]
      self.tree[idx] = weight
      parent = (idx - 1) // 2
28      while parent >= 0:
          self.tree[parent] += change
          parent = (parent - 1) // 2
30
32  def get_total(self):
      return self.tree[0]
34
36  def get_min(self):
      return min(self.tree[self.capacity - 1:self.capacity + self.size - 1])
38
40  def sample(self, v):
      idx = 0
      while idx < self.capacity - 1:
          l_idx = idx * 2 + 1
          r_idx = l_idx + 1
          if self.tree[l_idx] >= v:
              idx = l_idx
          else:
              idx = r_idx
              v = v - self.tree[l_idx]
48
50  point = idx - (self.capacity - 1)
      return point, self.data[point]

```

In class Memory we used SumTree to implement an experience pool. The main function is get-mini-batches, in which we sampled a mini-batch to train the evaluation model using randomized priority, which combines greedy priority preferential and uniform random sampling. By using this algorithm, we sped up the convergence and usually it will converge in no more than 2000 episodes.

```

1  class Memory(object):
      def __init__(self, batch_size, max_size, beta):
          self.batch_size = batch_size
          self.max_size = 2 ** math.floor(math.log2(max_size))
          self.beta = beta
          self._sum_tree = SumTree(max_size)
7
      def store_transition(self, s, a, r, s_, done):
          self._sum_tree.add((s, a, r, s_, done))
9
11     def get_mini_batches(self):
          n_sample = self.batch_size if self._sum_tree.size >= self.batch_size else self._sum_tree.size
          total = self._sum_tree.get_total()
          points = []
          transitions = []
          step = total // n_sample
          for i in range(n_sample):
              v = np.random.uniform(i * step, (i + 1) * step - 1)
              t = self._sum_tree.sample(v)
              points.append(t[0])
              transitions.append(t[1])
          return points, transitions
23
25     def update(self, points, td_error):
          for i in range(len(points)):
              self._sum_tree.update(points[i], td_error[i])

```

Finally, we will introduce our hyper-parameters setting.

- **The discounting factor:** We tested different discounting factors of 0.9, 0.95 and 0.99 under the same other hyper-parameters. After some basic tests, we got the sense that they contribute to convergence to the same extent. However it showed the tendency that a higher discounting factor may cost more steps to converge. We supposed that a higher value of the discounting factor means the agent shows more interest in the future, which is harder than focusing on the present situation. Thus the training process takes more pain. Finally we chose 0.9 for the discounting factor.
- **The learning rate α :** We tested different α of 0.01, 0.005 and 0.02 under the same other hyper-parameters. The outcome showed that a higher learning rate performs well in the initial term when the agent does not have enough samples, however it is obviously over-fitting in that it could not converge later with a heavier burden of learning. But if α is too low, the learning process takes a long time which is so expensive for us. Finally we chosen the learning rate α of 0.01.
- **The batch size:** We tested different α of 128, 64 and 32 under the same other hyper-parameters. A larger batch size will contribute to better use of past experience, however it will take too much time to wait for it to take effect. So we chosen the batch size of 32 in order to save time while it performs well as the others.
- **The epsilon:** We used Epsilon-Greedy Strategy to ensure exploration of the environment. We tested different epsilon and its related parameters such as `epsilon_decay` and `epsilon_min` under the same other hyper-parameters. We concluded that a higher epsilon is needed in the beginning to explore more possibilities. As time goes on, the epsilon value should decays to force the agent to choose actions from the trained model, otherwise it can't converge. So we chosen epsilon of 1.0, `epsilon_min` of 0.01 and `epsilon_decay` equals to $(\text{epsilon} - \text{epsilon_min}) / 50000$.

3. Results and Analysis

It's a pity that we didn't get a satisfying result and by the time limit, we can only run about 2500 episodes. Because after that, the speed decrease quickly and we need nearly one hour to run 100 episodes. From the FIG 5, we consider it converges to about 500 scores.

- The first problem we met is that we did not know how to implement prioritized experience replay. The first data structure we designed was using binary priority heap, all the nodes denoting both the transitions and the priority. However we did not know how to rank according to priority when we zipped the transitions and priority in a tuple, or how could we save them separately and maintain their subscript's accordance. Then we learned about SumTree, which is a new data structure for us. By trial and error, we implemented the prioritized experience replay successfully.
- When it comes to tuning the agent, we got lost in how these parameters influence the effect of the network. Our cognition about these parameters was all about a trend, such as high learning rate will cause vibration, and if the learning rate is too low it may not converge. But we did not have specific concepts about high and low for this net, so we spent a long time to try different numbers to find the feeling. Although we thought to have found a certain suitable value for each parameter, when we combined these suitable values it still could not turn out to perform well.
- At first, we didn't tune the batch-size and it ran so slow that it took us nearly 6 hours to run 4000 episodes. After we decreased the batch-size to 32, we thought we would get lower scores and higher speed, but it surprised us that it worked better with the same parameter setting as the 128 batch-size. This sped up my training process greatly.



FIG. 5: DQN Result

II. EQUIPMENT

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

- macOS High Sierra (v10.13.4)
- Python 3.6 version
- TensorFlow(v1.13.1), Keras(v2.3.1), Matplotlib, Pandas

III. CONCLUSION

In this assignment, we finally tried to implement Reinforcement Learning with Dyna-Q and DQN, which we have learnt about the theories and the principles behind the concepts in class, into mechanical puzzles. And we attempted to implement Epsilon-Greedy to get rid of the explore-exploit dilemma, aiming to compel the agent to explore more possibilities of state and action.

In the first experiment, I firstly used Dyna-Q Learning to explore a maze environment. To deal with the exploration-exploitation dilemma, I used Epsilon-Greedy Strategy, and set a decay for the probability ϵ . Besides, I connected the Q value of each state and action with the times the agent has been to the next state. In that way, the agent has a tendency to enter a state which it has never been to.

In the second experiment, we completed a DQN network ourselves. To tune the agent, we implemented prioritized experience replay. Thus we can fully exploit the most important experience. Also Epsilon-Greedy Strategy is implemented to ensure exploration. Finally our agent can roughly converge.

In the two experiments, the first tough problem is how to choose the parameters. Sometimes they have close relations. When changing one parameter, the others may be influenced. So I spent a lot of time to choose the proper parameters. The another toughest problem is the time to train DQN network. Because of the limitation of the hardware, it took too much time to wait for the training process to end.

After all, I sincerely appreciate our professor Yue Gao for providing us with such a great chance to implement Reinforcement Learning algorithms into practice. And I sincerely appreciate our teaching assistants for answering our doubts and helping us solve problems. It is such a special experience to do this work.

IV. APPENDIX

A. Codes of Experiment 1

```

from maze_env import Maze
2 import numpy as np
import pandas as pd
4 import random

6 UNIT = 40
MAZE_H = 6
8 MAZE_W = 6

10 class Agent:
    ### START CODE HERE ###
12     ### a random agent ###

14     def __init__(self, actions):
        self.actions = actions
16         self.epsilon = 1

18     def choose_action(self, observation):
        action = np.random.choice(self.actions)
20         return action

22     ### END CODE HERE ###

24 class myAgent:
    ### and agent of mine using Dyna-Q ###

26     def __init__(self, actions):
        self.actions = actions #actions
        self.epsilon = 0.1 #initialize for epsilon greedy
30         self.decay = 0.99 #a decay factor on epsilon
        self.gamma = 0.8 #choose a key for gamma
        self.alpha = 0.6 #choose a key for alpha
        self.q_dict = {} #remember Q(s,a)
        self.model_dict = {} #remember Model(s,a)
34         self.Dyna_N = 10 #repeat N times
        self.epsilon_dict = {} #remember the next step it can take when in this state
        self.has_been_to_this_state = {} #whether has been to this state?
        self.greedy_random_actions = [] #random actions for e-greedy
        self.cnt_of_state = {} #remember how many times has come to this state
        self.cnt_of_state_all_actions = {}
        self.new_state = (20.0, 20.0)
        self.all_arrived = False
        self.k = 0.3
        self.epoch_num = 0 # same with episode

46         for is_treasure_be_found in [False, True]:
            for i in range(MAZE_W):
48                 for j in range(MAZE_H):
                    , ,
50 ----- initialize state and some dictionaries and stuff
                    state (origin_for_x, origin_for_y, destination_for_x, destination_for_y, a_bool)
52 ----- every item in q-state is [a,b,c,d], which corresponds to rewards of four directions
                    , ,
54                 state = (float(i*UNIT+5), float(j*UNIT+5), float((i+1)*UNIT-5), float((j+1)*UNIT-5),
                            is_treasure_be_found)
                            self.q_dict[state] = [0, 0, 0, 0]
56                 self.epsilon_dict[state] = [False, False, False, False]
                            self.has_been_to_this_state[float((state[0]+state[2])/2), float((state[1]+state[3])/2)] = False
58                 self.cnt_of_state[state] = [0, 0, 0, 0]
                            self.cnt_of_state_all_actions[state] = 1
60

62     def update_random_actions_for_greedy(self, s):
        state = tuple(s)
64         self.all_arrived = False
        self.has_been_to_this_state[float((state[0]+state[2])/2), float((state[1]+state[3])/2)] = True
66         expected_states = []
        actions_to_get_close = []
68         present_state = ((state[0]+state[2])/2, (state[1]+state[3])/2)

70         #has not been to this state before
        for s1 in self.has_been_to_this_state:
72             if self.has_been_to_this_state[s1] == False:
                expected_states.append(s1)
74

76         if not expected_states:
            self.all_arrived = True
            return []

78         if present_state == self.new_state:
            new_state = random.sample(expected_states, 1)
80             new_state = new_state[0]
        else:
            new_state = self.new_state
82         if new_state[0] > present_state[0]:
            actions_to_get_close.append(2)
84         if new_state[0] < present_state[0]:
            actions_to_get_close.append(3)
86         if new_state[1] < present_state[1]:
            actions_to_get_close.append(0)
88         if new_state[1] > present_state[1]:
            actions_to_get_close.append(1)
90

```

```

100         actions_to_get_close.append(1)
102     for a2 in actions_to_get_close:
104         if self.q_dict[state][a2] == -1:
106             actions_to_get_close.remove(a2)
108     self.new_state = new_state
110     return actions_to_get_close
112
114     def epsilon_decay(self, s):
116         """
117         epsilon is large at first, then gradually decays
118         """
119         return:
120         """
121         epsilon = self.epsilon
122         if self.all_arrived == False:
123             epsilon = 0.1
124         if self.all_arrived == True:
125             epsilon = 0
126         return epsilon
127
128     def update_cnt_of_state(self, s):
129         state = tuple(s)
130         total = 0
131         for l in range(4):
132             total += self.cnt_of_state[state][l]
133         self.cnt_of_state_all_actions[state] = total
134
135     def choose_action(self, s):
136         """
137         choose an action, use epsilon greedy
138         param:s: present state, a_list
139         return: an action (x0,y0,x1,y1)
140         """
141         actions_to_get_close2 = self.update_random_actions_for_greedy(s)
142         state = tuple(s) # transform list s into a tuple
143         self.has_been_to_this_state[float((state[0] + state[2]) / 2), float((state[1] + state[3]) / 2)] = True
144         self.epsilon = self.epsilon_decay(state)
145         # use epsilon greedy
146         if random.random() < self.epsilon:
147             action = np.random.choice(self.actions) # a random action
148         else:
149             max_actions = []
150             for i in range(4):
151                 if self.q_dict[state][i] == max(self.q_dict[state]):
152                     max_actions.append(i)
153             action = random.choice(max_actions) #choose action of maximum value
154
155         if (action == 0 and state[1] == 5) or (action == 1 and state[3] == 235) or (action == 2 and state[2] ==
156             235) or \
157             (action == 3 and state[0] == 5):
158             self.q_dict[state][action] = -100
159             self.epsilon_dict[state][action] = True
160             self.update_cnt_of_state(s)
161             return self.choose_action(list(s))
162         else:
163             return action
164
165     def update(self, s, a, s_, r):
166         """
167         update Q-value
168         param:s: previous state
169         param:a: action
170         param:s_: new state
171         param:r: reward
172         return: q_dict
173         """
174         state = tuple(s)
175         state_ = tuple(s_)
176         if self.all_arrived:
177             self.k = 0
178
179         for j in range(4):
180             self.cnt_of_state[state_][j] = self.cnt_of_state[state_][j] * 0.92
181
182         self.cnt_of_state[state][a] += 1
183         self.update_cnt_of_state(s)
184         self.has_been_to_this_state[float((state[0] + state[2]) / 2), float((state[1] + state[3]) / 2)] = True
185         self.q_dict[state][a] += self.alpha*(r + self.gamma*max(((self.q_dict[state_][i]) - self.cnt_of_state[state_][i])*self.k) for i in range(4)) - self.q_dict[state][a])
186         self.model_dict[(state, a)] = [state_, r]
187
188     N = self.Dyna_N
189     for n in range(N): # iteration for N times
190         ms, ma = random.choice(list(self.model_dict))
191         ms_, mr = self.model_dict[(ms, ma)]
192         self.q_dict[ms][ma] += self.alpha * (mr + self.gamma * max(((self.q_dict[ms_][i]) - self.cnt_of_state[
193             state_][i])*self.k) for i in range(4)) - self.q_dict[ms][ma])
194
195     return self.q_dict

```

B. Codes of Experiment 2

```

1  #- coding:utf-8 -#
2  # DQN homework.
3  import os
4  import sys
5  import gym
6  import pylab
7  import random
8  import numpy as np
9  from collections import deque
10 from keras.layers import Dense
11 from keras.optimizers import Adam
12 from keras.models import Sequential
13 from gym import wrappers
14 from utils import *
15 import math

17 # hyper-parameter.
18 EPISODES = 5000
19
21 class SumTree:
22     def __init__(self, capacity):
23
24         self.capacity = capacity
25         self.tree = [0] * (2 * capacity - 1)
26         self.data = [None] * capacity
27         self.size = 0
28         self.curr_point = 0
29
31     def add(self, data):
32         self.data[self.curr_point] = data
33         self.update(self.curr_point, max(self.tree[self.capacity - 1:self.capacity + self.size]) + 1)
34
35         self.curr_point += 1
36         if self.curr_point >= self.capacity:
37             self.curr_point = 0
38
39         if self.size < self.capacity:
40             self.size += 1
41
43     def update(self, point, weight):
44         idx = point + self.capacity - 1
45         change = weight - self.tree[idx]
46         self.tree[idx] = weight
47         parent = (idx - 1) // 2
48         while parent >= 0:
49             self.tree[parent] += change
50             parent = (parent - 1) // 2
51
53     def get_total(self):
54         return self.tree[0]
55
57     def get_min(self):
58         return min(self.tree[self.capacity - 1:self.capacity + self.size - 1])
59
61     def sample(self, v):
62         idx = 0
63         while idx < self.capacity - 1:
64             l_idx = idx * 2 + 1
65             r_idx = l_idx + 1
66             if self.tree[l_idx] >= v:
67                 idx = l_idx
68             else:
69                 idx = r_idx
70                 v = v - self.tree[l_idx]
71
72         point = idx - (self.capacity - 1)
73         return point, self.data[point]
74
76 class Memory(object):
77     def __init__(self, batch_size, max_size, beta):
78         self.batch_size = batch_size
79         self.max_size = 2 ** math.floor(math.log2(max_size))
80         self.beta = beta
81         self._sum_tree = SumTree(max_size)
82
84     def store_transition(self, s, a, r, s_, done):
85         self._sum_tree.add((s, a, r, s_, done))
86
88     def get_mini_batches(self):
89         n_sample = self.batch_size if self._sum_tree.size >= self.batch_size else self._sum_tree.size
90         total = self._sum_tree.get_total()
91         points = []
92         transitions = []
93         step = total // n_sample
94         for i in range(n_sample):
95             v = np.random.uniform(i * step, (i + 1) * step - 1)
96             t = self._sum_tree.sample(v)
97             points.append(t[0])
98             transitions.append(t[1])
99         return points, transitions

```

```

97     def update(self, points, td_error):
98         for i in range(len(points)):
99             self._sum_tree.update(points[i], td_error[i])
100
101 class DQNAgent:
102     def __init__(self, state_size, action_size):
103         self.render = False
104         self.state_size = state_size
105         self.action_size = action_size
106
107         # These are hyper parameters for the DQN
108         self.discount_factor = 0.9
109         self.learning_rate = 0.01
110         self.epsilon = 1.0
111         self.epsilon_min = 0.01
112         self.epsilon_decay = (self.epsilon - self.epsilon_min) / 50000
113         self.batch_size = 32
114         self.train_start = 1000
115         self.beta = 0.1
116         self.maxlen = 100000
117         self.memory = Memory(self.batch_size, self.maxlen, self.beta)
118
119         # create main model
120         self.model_target = self.build_model()
121         self.model_eval = self.build_model()
122
123     # approximate Q function using Neural Network
124     def build_model(self):
125         model = Sequential()
126         model.add(Dense(128, input_dim=self.state_size, activation='relu',
127                        kernel_initializer='he_uniform'))
128         model.add(Dense(32, activation='relu',
129                        kernel_initializer='he_uniform'))
130         model.add(Dense(self.action_size, activation='linear',
131                        kernel_initializer='he_uniform'))
132         model.summary()
133         model.compile(loss='mse', optimizer=Adam(lr=self.learning_rate))
134         return model
135
136     # get action from model using epsilon-greedy policy
137     def get_action(self, state, step):
138         if step > 250:
139             if self.epsilon > self.epsilon_min:
140                 self.epsilon -= self.epsilon_decay
141             if self.epsilon < self.epsilon_min:
142                 self.epsilon = self.epsilon_min
143             if np.random.rand() <= self.epsilon:
144                 return random.randrange(self.action_size)
145             else:
146                 q_value = self.model_eval.predict(state)
147                 return np.argmax(q_value[0])
148
149         # save sample <s,a,r,s'> to the replay memory
150     def append_sample(self, state, action, reward, next_state, done):
151         self.memory.store_transition(state, action, reward, next_state, done)
152
153     # pick samples randomly from replay memory (with batch size)
154     def train_model(self):
155         if self.memory._sum_tree.size < self.train_start:
156             return
157         batch_size = min(self.batch_size, self.memory._sum_tree.size)
158         points, mini_batch = self.memory.get_mini_batches()
159         update_input = np.zeros((batch_size, self.state_size))
160         update_target = np.zeros((batch_size, self.action_size))
161         td_error = np.zeros(batch_size)
162         action, reward, done = [], [], []
163
164         for i in range(self.batch_size):
165             update_input[i] = mini_batch[i][0]
166             action.append(mini_batch[i][1])
167             reward.append(mini_batch[i][2])
168             update_target[i] = mini_batch[i][3]
169             done.append(mini_batch[i][4])
170
171         target = self.model_eval.predict(update_input)
172         target_val = self.model_target.predict(update_target)
173
174         for i in range(self.batch_size):
175             # Q-Learning: get maximum Q-value at s' from model
176             if done[i]:
177                 td_error[i] = abs(reward[i] - target[i][action[i]])
178                 target[i][action[i]] = reward[i]
179             else:
180                 td_error[i] = abs(reward[i] + self.discount_factor * (np.amax(target_val[i])) - target[i][action[i]])
181                 target[i][action[i]] = reward[i] + self.discount_factor * (
182                     np.amax(target_val[i]))
183
184         # and do the model fit!
185         self.model_eval.fit(update_input, target, batch_size=self.batch_size,
186                             epochs=1, verbose=0)
187         self.memory.update(points, td_error)
188
189     def eval2target(self):
190         self.model_target.set_weights(self.model_eval.get_weights())
191
192 if __name__ == "__main__":
193     # load the gym env

```

```

197 env = gym.make('MsPacman-ram-v0')
198 # set random seeds to get reproduceable result (recommended)
199 set_random_seed(0)
200 # get size of state and action from environment
201 state_size = env.observation_space.shape[0]
202 action_size = env.action_space.n
203 # create the agent
204 agent = DQNAgent(state_size, action_size)
205 # log the training result
206 scores, episodes = [], []
207 graph_episodes = []
208 graph_score = []
209 avg_length = 100
210 sum_score = 0
211 iteration = 0
212 update_times = 0
213
214 # train DQN
215 for e in range(EPISODES):
216     done = False
217     score = 0
218     state = env.reset()
219     state = np.reshape(state, [1, state_size])
220     lives = 3
221     while not done:
222         dead = False
223         step = 0
224         while not dead:
225             # render the gym env
226             if agent.render:
227                 env.render()
228                 step += 1
229             # get action for the current state
230             action = agent.get_action(state, step)
231             # take the action in the gym env, obtain the next state
232             next_state, reward, done, info = env.step(action)
233             next_state = np.reshape(next_state, [1, state_size])
234             # judge if the agent dead
235             dead = info['ale.lives'] < lives
236             lives = info['ale.lives']
237             # update score value
238             score += reward
239             # save the sample <s, a, r, s'> to the replay memory
240             agent.append_sample(state, action, reward, next_state, done)
241             # train the evaluation network
242             if step % 4 == 0:
243                 agent.train_model()
244             # go to the next state
245             state = next_state
246
247         # update the target network after some iterations.
248         if update_times >= 2500:
249             update_times = 0
250             agent.eval2target()
251         # print info and draw the figure.
252         if done:
253             scores.append(score)
254             sum_score += score
255             episodes.append(e)
256             # plot the reward each episode
257             # pylab.plot(episodes, scores, 'b')
258             print("episode:", e, "score:", score, "memory length:",
259                   agent.memory.sum_tree.size, "epsilon:", agent.epsilon, "learning rate", agent.learning_rate)
260             if e % avg_length == 0:
261                 graph_episodes.append(e)
262                 graph_score.append(sum_score / avg_length)
263                 sum_score = 0
264             # plot the reward each avg length episodes
265             pylab.plot(graph_episodes, graph_score, 'r')
266             pylab.savefig("./pacman_avg.png")

```