

# Homework 4:

## Reinforcement Learning

### Report Template

Please keep the title of each section and delete examples. Note that please keep the questions listed in Part III.

#### Part I. Implementation (-5 if not explain in detail):

- Please screenshot your code snippets of **Part 1 ~ Part 3**, and explain your implementation.

##### 1.taxi.py

we first use random to determine the action will choosed

if the random number less than self.epsilon, we choose the action by random in the possible action

ifthe random number greater than self.epsilon, we choose the action with maximun value in qtable

```
8 def choose_action(self, state):
9     """
10     Choose the best action with given state and epsilon.
11
12     Parameters:
13         state: A representation of the current state of the enviornment.
14         epsilon: Determines the explore/exploit rate of the agent.
15
16     Returns:
17         action: The action to be evaluated.
18     """
19     # Begin your code
20     # TODO
21     epsilon = np.random.rand()
22     if epsilon < self.epsilon:
23         action = self.env.action_space.sample()
24     else:
25         action = np.argmax(self.qtable[state])
26     return action
27     raise NotImplementedError("Not implemented yet.")
28     # End your code
```

from the formula of Q-learning we update the qtable by the equation.

$$newQ_{S,A} = Q_{S,A} + \alpha \left( R_{S,A} + \gamma * \max Q'(s', a') - Q_{S,A} \right)$$

基于状态和  
行动的新Q值
c当前Q值
基于状态和  
行动的奖励
在给定新的状态和  
行动下未来最大的奖励

学习效率
折扣因子

```
def learn(self, state, action, reward, next_state, done):
    """
    Calculate the new q-value base on the reward and state transformation observed after taking the action.

    Parameters:
        state: The state of the enviornment before taking the action.
        action: The exacuted action.
        reward: Obtained from the enviornment after taking the action.
        next_state: The state of the enviornment after taking the action.
        done: A boolean indicates whether the episode is done.

    Returns:
        None (Don't need to return anything)
    """
    # Begin your code
    max_val = np.max(self.qtable[next_state])
    self.qtable[state,action] = self.qtable[state,action] + self.learning_rate*(reward+self.gamma*max_val-self.qtable[state,action])
    if done:
        np.save("../Tables/taxi_table.npy", self.qtable)
    # End your code
```

just return the maxvalue in the specific state in q-table.

```
def check_max_Q(self, state):
    """
    - Implement the function calculating the max Q value of given state.
    - Check the max Q value of initial state

    Parameter:
        state: the state to be check.

    Return:
        max_q: the max Q value of given state
    """
    # Begin your code
    # TODO
    return np.max(self.qtable[state])
    raise NotImplementedError("Not implemented yet.")
    # End your code
```

## 2.cartpole.py

we use function linspace to separate the upper bound and lower bound into n interval, but we don't the first and last element, so we eliminate the first and last element with

```
ans[1:-1]
```

```

39 def init_bins(self, lower_bound, upper_bound, num_bins):
40     """
41     Slice the interval into #num_bins parts.
42     Parameters:
43         lower_bound: The lower bound of the interval.
44         upper_bound: The upper bound of the interval.
45         num_bins: Number of parts to be sliced.
46     Returns:
47         a numpy array of #num_bins - 1 quantiles.
48     Example:
49         Let's say that we want to slice [0, 10] into five parts,
50         that means we need 4 quantiles that divide [0, 10].
51         Thus the return of init_bins(0, 10, 5) should be [2. 4. 6. 8.].
52     Hints:
53         1. This can be done with a numpy function.
54     """
55     # Begin your code
56     # TODO
57     ans = np.linspace(lower_bound, upper_bound, num_bins+1)
58     ans = ans[1:-1];
59     return ans
60     #raise NotImplementedError("Not implemented yet.")
61     # End your code

```

we can easily use function digitize in numpy to calculate the target value in which interval.

```

63 def discretize_value(self, value, bins):
64     """
65     Discretize the value with given bins.
66     Parameters:
67         value: The value to be discretized.
68         bins: A numpy array of quantiles
69     returns:
70         The discretized value.
71     Example:
72         With given bins [2. 4. 6. 8.] and "5" being the value we're going to discretize.
73         The return value of discretize_value(5, [2. 4. 6. 8.]) should be 2, since 4 <= 5 < 6 where [4, 6) is the 3rd bin.
74     Hints:
75         1. This can be done with a numpy function.
76     """
77     # Begin your code
78     # TODO
79     ans = np.digitize(value, bins)
80     return ans
81     raise NotImplementedError("Not implemented yet.")
82     # End your code

```

when we get the observation value, we need to turn it into discrete value. So for every value in observation, we use the discretize function we write before to turn it into discrete value. 1

```

84 def discretize_observation(self, observation):
85     """
86     Discretize the observation which we observed from a continuous state space.
87     Parameters:
88         observation: The observation to be discretized, which is a list of 4 features:
89             1. cart position.
90             2. cart velocity.
91             3. pole angle.
92             4. tip velocity.
93     Returns:
94         state: A list of 4 discretized features which represents the state.
95     Hints:
96         1. All 4 features are in continuous space.
97         2. You need to implement discretize_value() and init_bins() first
98         3. You might find something useful in Agent.__init__()
99     """
100    # Begin your code
101    # TODO
102    state = []
103    n = len(observation)
104    #print(n)
105    for i in range(0, n):
106        state.append(self.discretize_value(observation[i], self.bins[i]))
107    return state
108    #raise NotImplementedError("Not implemented yet.")
109    # End your code

```

we first use random to determine the action will be chosen

if the random number less than self.epsilon, we choose the action by random in the possible action  
 if the random number greater than self.epsilon, we choose the action with maximum value in qtable

```

111     def choose_action(self, state):
112         """
113         Choose the best action with given state and epsilon.
114         Parameters:
115             state: A representation of the current state of the environment.
116             epsilon: Determines the explore/exploit rate of the agent.
117         Returns:
118             action: The action to be evaluated.
119         """
120         # Begin your code
121         # TODO
122         # print(state)
123         epsilon = np.random.rand()
124         if epsilon < self.epsilon:
125             action = self.env.action_space.sample()
126         else:
127             action = np.argmax(self.qtable[state[0], state[1], state[2], state[3]])
128         return action
129         raise NotImplementedError("Not implemented yet.")
130         # End your code

```

from the formula of Q-learning we update the qtable by the equation.

$$\underbrace{new Q_{S,A}}_{\text{基于状态和动作的新Q值}} = \underbrace{Q_{S,A}}_{\text{c当前Q值}} + \underbrace{\alpha}_{\text{学习效率}} \left( \underbrace{R_{S,A}}_{\text{基于状态和动作的奖励}} + \underbrace{\gamma}_{\text{折扣因子}} * \underbrace{\max Q'(s', a')}_{\text{在给定新的状态和动作下未来最大的奖励}} \right) - Q_{S,A}$$

```

132     def learn(self, state, action, reward, next_state, done):
133         """
134         Calculate the new q-value base on the reward and state transformation observed after taking the action.
135         Parameters:
136             state: The state of the environment before taking the action.
137             action: The executed action.
138             reward: Obtained from the environment after taking the action.
139             next_state: The state of the environment after taking the action.
140             done: A boolean indicates whether the episode is done.
141         Returns:
142             None (Don't need to return anything)
143         """
144         # Begin your code
145         # TODO
146         max_val = np.max(self.qtable[next_state[0], next_state[1], next_state[2], next_state[3]])
147         self.qtable[state[0], state[1], state[2], state[3], action] = self.qtable[state[0], state[1], state[2], state[3], action] + self.learning_rate * (reward + sel
148         if done:
149             np.save("../Tables/cartpole_table.npy", self.qtable)
150             # raise NotImplementedError("Not implemented yet.")
151         # End your code

```

just return the maxvalue in the specific state in q-table.

```

155     def check_max_Q(self):
156         """
157         - Implement the function calculating the max Q value of initial state(self.env.reset()).
158         - Check the max Q value of initial state
159         Parameter:
160             self: the agent itself.
161             (Don't pass additional parameters to the function.)
162             (All you need have been initialized in the constructor.)
163         Return:
164             max_q: the max Q value of initial state(self.env.reset())
165         """
166         # Begin your code
167         # TODO
168         state = self.discretize_observation(self.env.reset())
169         return np.max(self.qtable[state[0],state[1],state[2],state[3]])
170         #raise NotImplementedError("Not implemented yet.")
171         # End your code

```

### 3.DQN.py

from the steps in spec:

- 1.every 100 time we load(update) the target\_net (write by TA)
- 2.we first get the state,action reqard next\_state and done from self.buffer and we turn this values into nparray and turn into tensor matrix and ensure they have the same batch size
- 4.we put the state and next state into the network to get the Q value and next Q value note that we only consider the next Q value with haven' t done. Last we use the formula to calculate the target Q value
- 5.we call the loss function from network and put Qvalue and target Q value into it.
- 6.we use loss.backward() to reach backproagation
- 7.we use self.optimizer.step() to optimize the loss function

```

109  ✓ def learn(self):
110  ✓     """
111         - Implement the learning function.
112         - Here are the hints to implement.
113         Steps:
114         -----
115         1. Update target net by current net every 100 times. (we have done this for you)
116         2. Sample trajectories of batch size from the replay buffer.
117         3. Forward the data to the evaluate net and the target net.
118         4. Compute the loss with MSE.
119         5. Zero-out the gradients.
120         6. Backpropagation.
121         7. Optimize the loss function.
122         -----
123         Parameters:
124             self: the agent itself.
125             (Don't pass additional parameters to the function.)
126             (All you need have been initialized in the constructor.)
127         Returns:
128             None (Don't need to return anything)
129         """
130  ✓     if self.count % 100 == 0:
131         self.target_net.load_state_dict(self.evaluate_net.state_dict())
132

```

```

133     # Begin your code
134     # TODO
135     state, action, reward, next_state, done = self.buffer.sample(self.batch_size)
136     #print(state)
137     #memory = self.buffer.sample(self.batch_size)
138     state = torch.FloatTensor(np.asarray(state))
139     action = torch.LongTensor(np.asarray(action)).view(self.batch_size, 1)
140     reward = torch.IntTensor(np.asarray(reward)).view(self.batch_size, 1)
141     next_state = torch.FloatTensor(np.asarray(next_state))
142     done = torch.IntTensor(np.asarray(done)).view(self.batch_size, 1)
143
144     q_value = self.evaluate_net(state).gather(1, action)
145
146     q_next_value = self.target_net.forward(next_state).detach()
147
148     for i in range(self.batch_size):
149         if done[i]:
150             q_next_value[i] = 0
151     q_next_value = q_next_value.max(1)[0].view(self.batch_size, 1)
152     target_q_value = reward + self.gamma * q_next_value
153
154     loss_function = nn.MSELoss()
155     loss = loss_function(q_value, target_q_value)
156
157     self.optimizer.zero_grad()
158     loss.backward()
159     self.optimizer.step()
160     torch.save(self.target_net.state_dict(), "./Tables/DQN.pt")
161     # End your code

```

we first use random to determine the action will choosed  
if the random number less than self.epsilon, we choose the action by  
random in the possible action  
if the random number greater than self.epsilon, we first turn the state into  
tensor matrix and use unsqueeze to make it to higher dimension. Then we  
can put it into the nueral network to see what may the value get after  
action, we choose the action with max value.

```

163     def choose_action(self, state):
164         """
165         - Implement the action-choosing function.
166         - Choose the best action with given state and epsilon
167         Parameters:
168             self: the agent itself.
169             state: the current state of the enviornment.
170             (Don't pass additional parameters to the function.)
171             (All you need have been initialized in the constructor.)
172         Returns:
173             action: the chosen action.
174         """
175
176         with torch.no_grad():
177             # Begin your code
178             # TODO
179             epsilon = np.random.rand()
180             if epsilon < self.epsilon:
181                 action = self.env.action_space.sample()
182             else:
183                 print(state)
184                 tensor = torch.Tensor(state)
185                 temp = torch.unsqueeze(tensor, 0)
186                 value = self.evaluate_net.forward(temp)
187                 action = torch.argmax(value).item()
188         return action

```

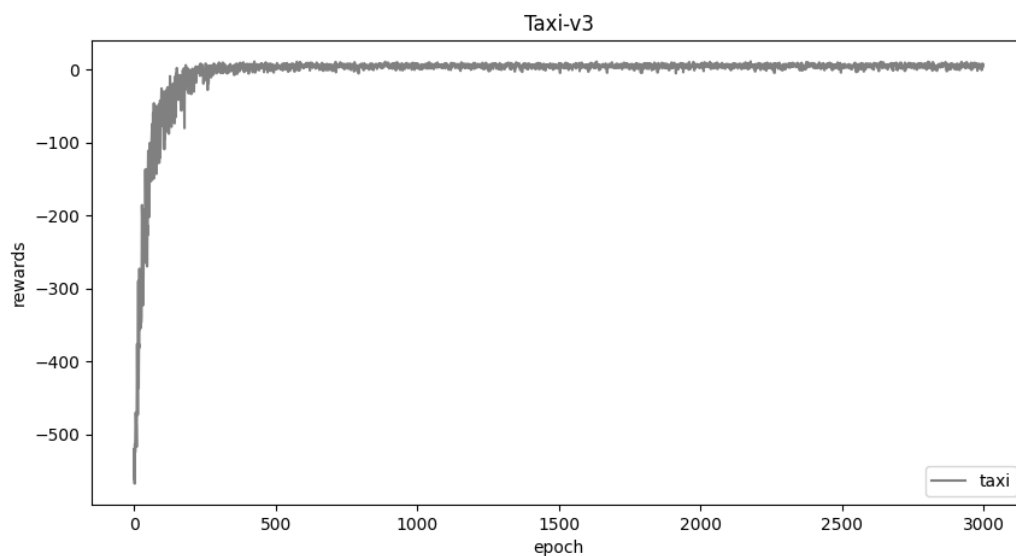
we first turn the initial state(`self.env.reset()`) into tensor matrix and use `unsqueeze` to make it to higher dimension. and use the neural network `target_net` to calculate the values we get after the calculation and we return the max Q value.

```
190 def check_max_Q(self):
191     """
192     - Implement the function calculating the max Q value of initial state(self.env.reset()).
193     - Check the max Q value of initial state
194     Parameter:
195         self: the agent itself.
196         (Don't pass additional parameters to the function.)
197         (All you need have been initialized in the constructor.)
198     Return:
199         max_q: the max Q value of initial state(self.env.reset())
200     """
201     # Begin your code
202     # TODO
203
204     tensor = torch.FloatTensor(self.env.reset())
205     temp = torch.unsqueeze(tensor,0)
206     action_values = self.target_net(temp)
207     max_Q = torch.max(action_values).item()
208     return max_Q
```

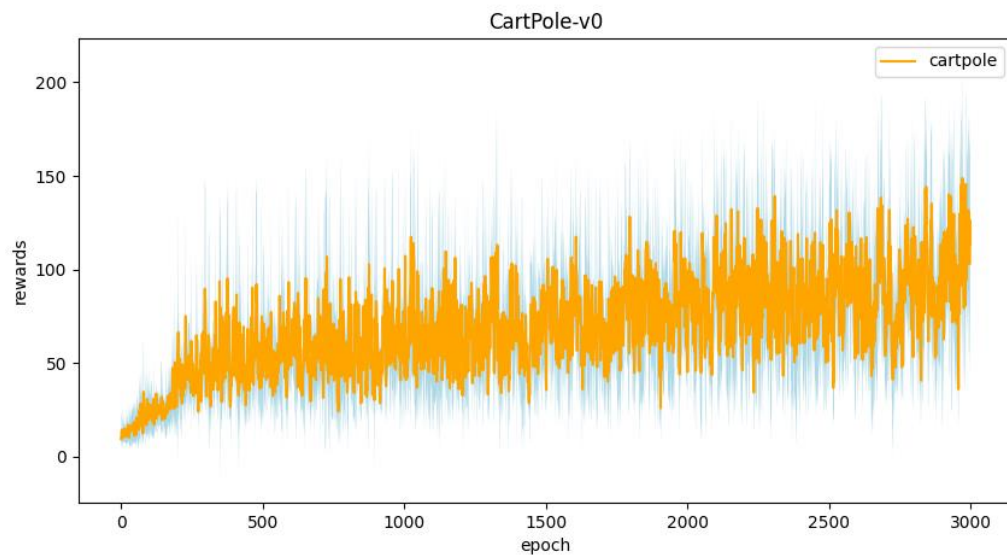
## Part II. Experiment Results:

Please paste [taxi.png](#), [cartpole.png](#), [DQN.png](#) and [compare.png](#) here.

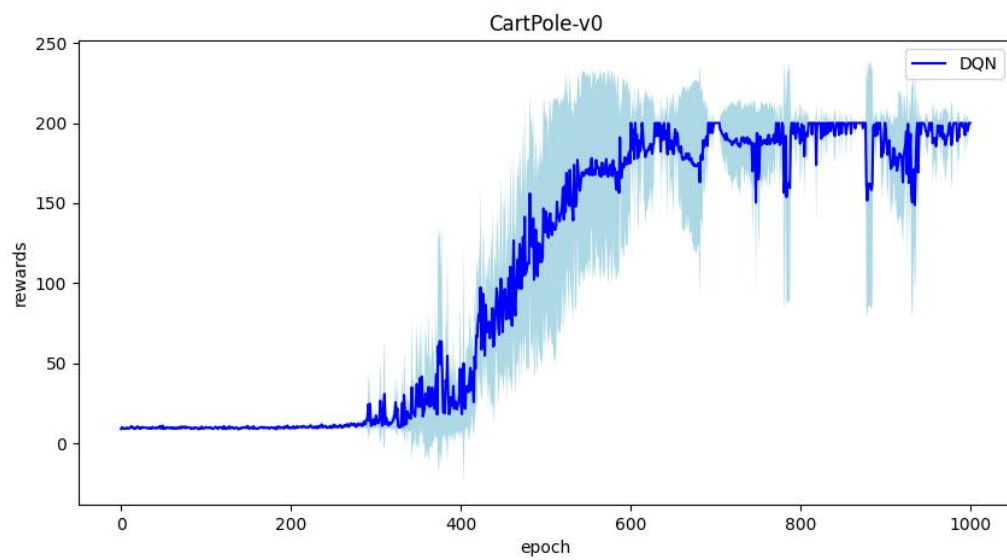
### 1. taxi.png:



### 2. cartpole.png

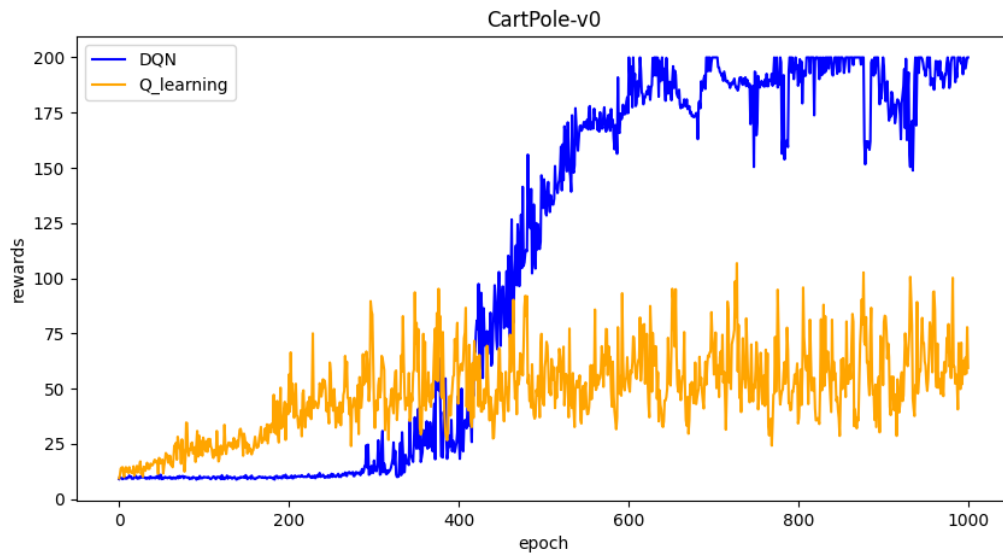


### 3. DQN.png



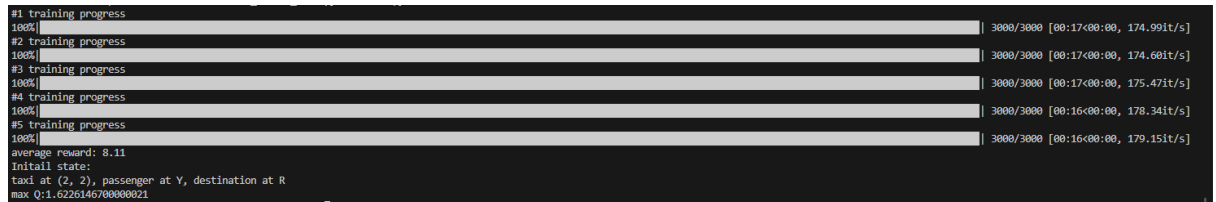
### 4. compare.png





### Part III. Question Answering (50%):

1. Calculate the optimal Q-value of a given state in Taxi-v3, and compare with the Q-value you learned (Please screenshot the result of the `check_max_Q` function to show the Q-value you learned). **(10%)**



$\gamma = 0.9$   
 taxi at (2,4) passenger at Y destination at R

$\Rightarrow L \rightarrow L \rightarrow D \rightarrow D \rightarrow \text{pick} \rightarrow U \rightarrow U \rightarrow U \rightarrow U \rightarrow \text{drop}$   
 (L: left, D: down, U: up) (9 step -1, 1 step 20)

$$Q_{\text{opt}} = -1(1 + 0.9 + 0.9^2 + \dots + 0.9^8) + 20 \times (0.9)^9$$

$$= -1 \times \frac{1 - 0.9^9}{1 - 0.9} + 20 \times 0.9^9$$

$$= 1.62261467000000017$$

2. Calculate the max Q-value of the initial state in CartPole-v0, and compare with the Q-value you learned. (Please screenshot the result of the "check max Q" function to show the Q-value you learned) **(10%)**

#1 training progress	100%	3000/3000	[00:19:00:00, 153.551t/s]
#2 training progress	100%	3000/3000	[00:20:00:00, 149.491t/s]
#3 training progress	100%	3000/3000	[00:18:00:00, 158.161t/s]
#4 training progress	100%	3000/3000	[00:19:00:00, 154.981t/s]
#5 training progress	100%	3000/3000	[00:22:00:00, 135.401t/s]
average reward: 177.2			
max Q: 31.72239218308627			

$$\gamma = 0.99, \text{reward} = 177.2$$

$$Q_{opt} = 1 + 0.99 + 0.99^2 + \dots + 0.99^{196}$$

$$= \frac{1 - (0.99)^{197}}{1 - 0.99}$$

make it more precise

replace 197 with 177.2

$$\approx \frac{1 - (0.99)^{177.2}}{1 - 0.99} \approx 33.18238437674456$$

3.
  - a. Why do we need to discretize the observation in Part 2? (3%)  
**Because the data we observe is continuous, but we need to make it into discrete states, so we have to discretize it to get the data in which interval.**
  - b. How do you expect the performance will be if we increase “num\_bins”? (3%)  
**The performance will become greater, because we will have more states, and the data after discretize will be more closer to the true data.**
  - c. Is there any concern if we increase “num\_bins”? (3%)  
**If we increase “num\_bins” the process speed may become slower, because we increase the number of states and also make the size of Q-table bigger and need more space to calculate.**
4. Which model (DQN, discretized Q learning) performs better in Cartpole-v0, and what are the reasons? (5%)  
**DQN perform better in cartpole-v0, in Qlearning we should divided the observation to limited discrete data and DQN can directly use the continuous data, which is more familiar to real data. so DQN can have higher performance.**
5.
  - a. What is the purpose of using the epsilon greedy algorithm while choosing an action? (3%)  
**the purpose is to balance exploration and exploitation, it can prevent the extreme cases, always choose the random action to**

**get more information, or choose the best action in with max Q value,(there may be other better action we havn' t explore).**

- b. What will happen, if we don't use the epsilon greedy algorithm in the CartPole-v0 environment? **(3%)**

**if we always use explore, choose action randomly, the episode will be easily ended and since we don' t choose the best action in qtable , we may not construct a good Q table.**

**if we always use exploit, we only choose the action that is known, this may cause us can' t not explore the better performance action that we haven' t explore and can' t get a good Q table .**

- c. Is it possible to achieve the same performance without the epsilon greedy algorithm in the CartPole-v0 environment? Why or Why not? **(3%)**

**Yes it is possible. Since the CartPole-v0 environment is easily, if we can find another method to balance the explore and exploit rate, we can get the same performance in the CartPole-v0 environment.**

- d. Why don't we need the epsilon greedy algorithm during the testing section? **(3%)**

**Because in the training section we need to find new action that we haven' t explore. But in the testing section we only want to test the q table we got after training, so we not need the epsilon greedy algorithm during the testing section.**

6. Why does "`with torch.no_grad() :`" do inside the "`choose_action`" function in DQN? **(4%)**

**we don' t need to calculate the gradient and back paropagation in choose action, so we use with torch.no\_grad() to disable gradient tracking.**