

# HW4 report

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

- Part 1

```
29     def choose_action(self, state):
30         """
31         Choose the best action with given state and epsilon.
32
33         Parameters:
34             state: A representation of the current state of the environment.
35             epsilon: Determines the explore/exploit rate of the agent.
36
37         Returns:
38             action: The action to be evaluated.
39         """
40         # Begin your code
41         # TODO
42         #raise NotImplementedError("Not implemented yet.")
43
44         action = np.argmax(self.qtable[state]) # Choose best action by Q-table
45         if np.random.rand() >= self.epsilon: # Decide whether to explore
46             action = self.env.action_space.sample()
47         return action
48
49         # End your code
```

```
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 environment before taking the action.
        action: The executed action.
        reward: Obtained from the environment after taking the action.
        next_state: The state of the environment after taking the action.
        done: A boolean indicates whether the episode is done.

    Returns:
        None (Don't need to return anything)
    """
    # Begin your code
    # TODO
    #raise NotImplementedError("Not implemented yet.")

    self.qtable[state, action] = (1 - self.learning_rate) * self.qtable[state, action] + \
    self.learning_rate * (reward + self.gamma * self.check_max_Q(next_state)) # Update Q-table

    # End your code
```

```
76     def check_max_Q(self, state):
77         """
78         - Implement the function calculating the max Q value of given state.
79         - Check the max Q value of initial state
80
81         Parameter:
82             state: the state to be check.
83         Return:
84             max_q: the max Q value of given state
85         """
86         # Begin your code
87         # TODO
88         #raise NotImplementedError("Not implemented yet.")
89
90         return np.max(self.qtable[state]) # Calculate max Q-value of given state
91
92         # End your code
```

- Part 2

```
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     #raise NotImplementedError("Not implemented yet.")
58
59     bins = np.linspace(lower_bound, upper_bound, num_bins, endpoint = False)
60     return bins[1:] # Use linspace to split into 7 parts and do not take first part
61
62     # End your code
```

```
64 def discretize_value(self, value, bins):
65     """
66     Discretize the value with given bins.
67     Parameters:
68         value: The value to be discretized.
69         bins: A numpy array of quantiles
70     returns:
71         The discretized value.
72     Example:
73         With given bins [2. 4. 6. 8.] and "5" being the value we're going to discretize.
74         The return value of discretize_value(5, [2. 4. 6. 8.]) should be 3.
75     Hints:
76         1. This can be done with a numpy function.
77     """
78     # Begin your code
79     # TODO
80     #raise NotImplementedError("Not implemented yet.")
81
82     return np.digitize(value, bins) # Use digitize to discretize value
83
84     # End your code
```

```

86 def discretize_observation(self, observation):
87     """
88     Discretize the observation which we observed from a continuous state space.
89     Parameters:
90         observation: The observation to be discretized, which is a list of 4 features:
91             1. cart position.
92             2. cart velocity.
93             3. pole angle.
94             4. tip velocity.
95     Returns:
96         state: A list of 4 discretized features which represents the state.
97     Hints:
98         1. All 4 features are in continuous space.
99         2. You need to implement discretize_value() and init_bins() first
100        3. You might find something useful in Agent.__init__()
101     """
102     # Begin your code
103     # TODO
104     #raise NotImplementedError("Not implemented yet.")
105
106     dis_value = () # Define a empty tuple
107     for obs, bin in zip(observation, self.bins):
108         dis_value += (self.discretize_value(obs, bin),) # Discretized the 4 features of observation
109     return dis_value
110
111     # End your code

```

```

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

```

```

3 def learn(self, state, action, reward, next_state, done):
4     """
5     Calculate the new q-value base on the reward and state transformation observed after taking the action
6     Parameters:
7         state: The state of the environment before taking the action.
8         action: The executed action.
9         reward: Obtained from the environment after taking the action.
10        next_state: The state of the environment after taking the action.
11        done: A boolean indicates whether the episode is done.
12    Returns:
13        None (Don't need to return anything)
14    """
15    # Begin your code
16    # TODO
17    #raise NotImplementedError("Not implemented yet.")
18
19    state_action_pair = state + (action,) # Use tuple to be Q-table's index
20    self.qtable[state_action_pair] = (1 - self.learning_rate) * self.qtable[state_action_pair] + \
21    self.learning_rate * (reward + self.gamma * np.max(self.qtable[next_state])) # Update Q-table
22
23    # End your code

```

```

157 def check_max_Q(self):
158     """
159     - Implement the function calculating the max Q value of initial state(self.env.reset()).
160     - Check the max Q value of initial state
161     Parameter:
162         self: the agent itself.
163         (Don't pass additional parameters to the function.)
164         (All you need have been initialized in the constructor.)
165     Return:
166         max_q: the max Q value of initial state(self.env.reset())
167     """
168     # Begin your code
169     # TODO
170     #raise NotImplementedError("Not implemented yet.")
171
172     initial_state = self.discretize_observation(self.env.reset()) # Discretize the initial state
173     return np.max(self.qtable[initial_state]) # Calculate max Q-value of initial state
174
175     # End your code

```

### • Part 3

```

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     #raise NotImplementedError("Not implemented yet.")

```

```

136
137     # 2. Sample trajectories of batch size from the replay buffer and make some modifications.
138     state_batch, action_batch, reward_batch, next_state_batch, mask_batch = self.buffer.sample(self.batch_size)
139     state_batch = torch.Tensor(state_batch)
140     action_batch = torch.LongTensor(action_batch).reshape(self.batch_size, 1)
141     reward_batch = torch.Tensor(reward_batch).reshape(self.batch_size, 1)
142     next_state_batch = torch.Tensor(next_state_batch)
143     mask_batch = torch.Tensor(mask_batch).reshape(self.batch_size, 1)
144
145     # 3. Forward the data to the evaluate net and the target net and make some modifications.
146     q_eval = self.evaluate_net(state_batch).gather(1, action_batch)
147     next_q = self.target_net(next_state_batch)
148     next_q = torch.Tensor([max(q) for q in next_q]).reshape(self.batch_size, 1)
149     q_target = reward_batch + self.gamma * next_q * (1 - mask_batch)
150     # I use 1 - mask to make sure if episode is done, it don't take additional Q-value
151
152     # 4. Compute the loss with MSE.
153     loss = F.mse_loss(q_eval, q_target)

```

```

154
155         # 5. Zero-out the gradients.
156         self.optimizer.zero_grad()
157
158         # 6. Backpropagation.
159         loss.backward()
160
161         # 7. Optimize the loss function and update count.
162         self.optimizer.step()
163         self.count += 1
164
165         # End your code

```

```

168 def choose_action(self, state):
169     """
170     - Implement the action-choosing function.
171     - Choose the best action with given state and epsilon
172     Parameters:
173         self: the agent itself.
174         state: the current state of the environment.
175         (Don't pass additional parameters to the function.)
176         (All you need have been initialized in the constructor.)
177     Returns:
178         action: the chosen action.
179     """
180     with torch.no_grad():
181         # Begin your code
182         # TODO
183         #raise NotImplementedError("Not implemented yet.")
184
185         state = torch.Tensor(state)
186         action = torch.argmax(self.evaluate_net(state)).item() # Choose best action by Q-table
187         if np.random.rand() >= self.epsilon: # Decide whether to explore
188             action = self.env.action_space.sample()
189
190         # End your code

```

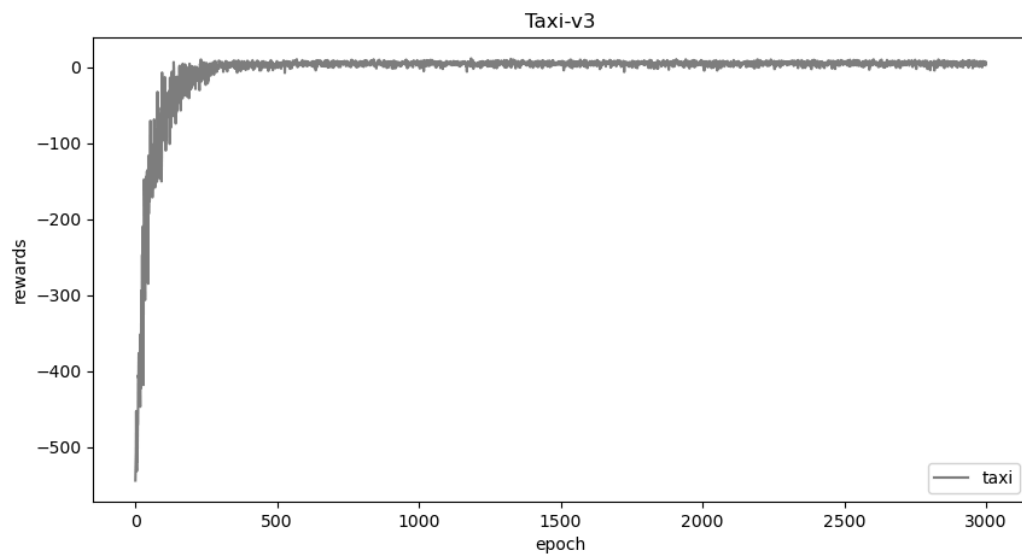
```

193 def check_max_Q(self):
194     """
195     - Implement the function calculating the max Q value of initial state(self.env.reset()).
196     - Check the max Q value of initial state
197     Parameter:
198         self: the agent itself.
199         (Don't pass additional parameters to the function.)
200         (All you need have been initialized in the constructor.)
201     Return:
202         max_q: the max Q value of initial state(self.env.reset())
203     """
204     # Begin your code
205     # TODO
206     #raise NotImplementedError("Not implemented yet.")
207
208     initial_state = torch.Tensor(self.env.reset())
209     q_value = self.target_net(initial_state)
210     return float(max(q_value)) # Calculate max Q-value of initial state
211
212     # End your code

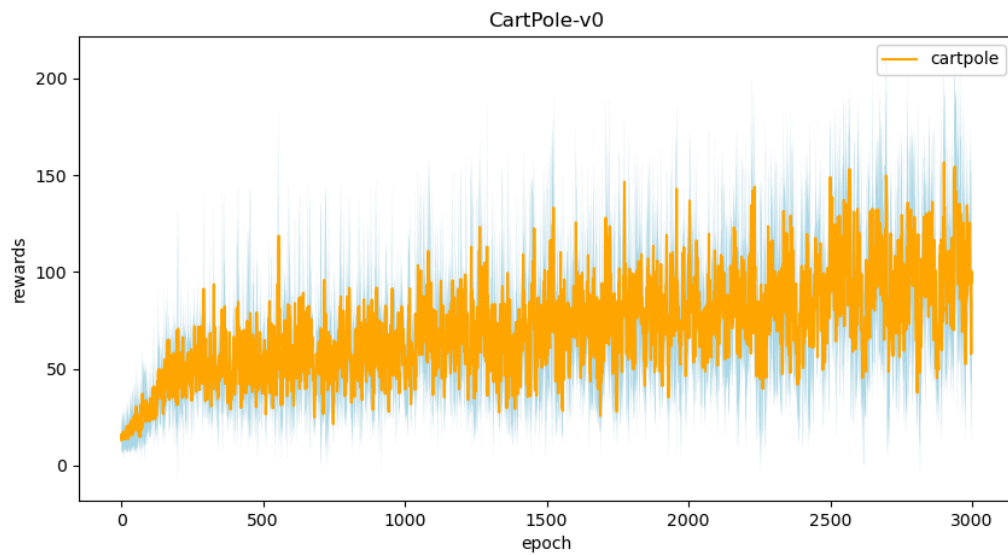
```

## Part II. Experiment Results:

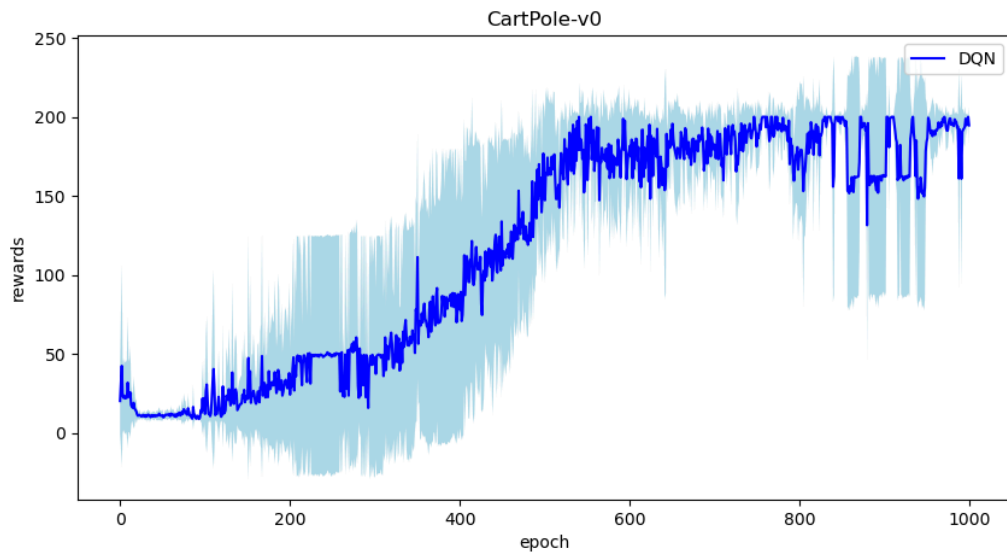
### 1. taxi.png



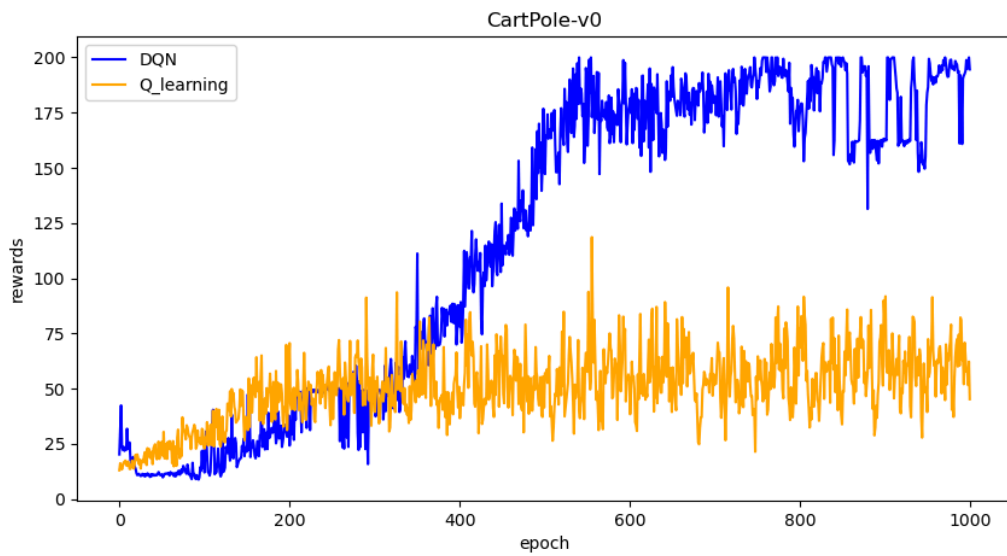
### 2. cartpole.png



### 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%)

Initial state:  
taxi at (2, 2), passenger at Y, destination at R  
max Q:1.6226146699999995

reward {  
-1, per step unless other reward is triggered  
+20, delivering passenger  
-10, executing "pick up" and "drop off" illegally

Need 5 steps to pick up passenger  
and 5 steps to drop off passenger

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

$$= -\left(\frac{1 - 0.9^9}{0.1}\right) + 20 \times 0.9^9 \approx 1.6226$$

2. Calculate the optimal Q-value of the initial state in CartPole-v0, and compare with the Q-value you learned (both cartpole.py and DQN.py). (Please screenshot the result of the "check\_max\_Q" function to show the Q-value you learned) (10%)

max Q:28.95079927460375

reward = +1, per step

Episode end when episode length reaches 200

$$\Rightarrow Q_{opt} = 1 + 0.97 + 0.97^2 + \dots + 0.97^{199} = \frac{1 - (0.97)^{200}}{0.03} \approx 33.58$$



3.

- a. Why do we need to discretize the observation in Part 2? (3%)

Because we use Q-table to store Q-value, if observation is continuous, it will be difficult to store.

- b. How do you expect the performance will be if we increase “num\_bins”? (3%)

I think the performance will be better because it has a more accurate index for Q-table. But it will cost more memory for Q-table to store more accurate information.

- c. Is there any concern if we increase “num\_bins”? (3%)

As I said in (b), it will cost more memory. Moreover, it will spend more time on updating the Q-table due to more states.

4. Which model (DQN, discretized Q learning) performs better in Cartpole-v0, and what are the reasons? (5%)

DQN performs better than discretized Q learning in Cartpole-v0. Because discretized Q learning discretized the state space, which will lose accuracy contrast to original continuous state space.

5.

- a. What is the purpose of using the epsilon greedy algorithm while choosing an action? (3%)

To let the agent explore, not just follow the greedy choice.

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

If we don't use epsilon greedy, but just use a common greedy algorithm, the performance will depend on whether you are lucky or not. If you are lucky, you may get all important information without exploration, otherwise you may not get good performance due to the loss of exploration.

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

For sure you can achieve the same performance without epsilon greedy. As I said in (b), if you are lucky enough, or you can fine-tune the hyperparameter to achieve the similar performance.

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

Because we just want to test our model whether it is good or not, we don't need the epsilon greedy to make exploration.

6. Why does “with torch.no\_grad():” do inside the “choose\_action” function in DQN? (4%)

It means that it will not compute the gradients and do backpropagation for this part because there is no need for this. And it can speed up the computation.