

Homework 4: Reinforcement Learning

學號: 111550129 姓名: 林彥亨

Part I. Implementation

Q-learning in taxi

- `choose_action()`

```
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 enviornment.
35         epsilon: Determines the explore/explot 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     # use random number from 0 to 1
45     # if the number is bigger than the epsilon, then return the maximum Q in the Q table.
46     # if the number is smaller than the epsilon, then return the random action.
47     if np.random.uniform(0, 1) > self.epsilon:
48         return np.argmax(self.qtable[state])
49     else:
50         return env.action_space.sample()
51
52     # End your code
```

- `learn()`

```
54 def learn(self, state, action, reward, next_state, done):
55     """
56     Calculate the new q-value base on the reward and state transformation observed after taking the action.
57
58     Parameters:
59         state: The state of the enviornment before taking the action.
60         action: The exacuted action.
61         reward: Obtained from the enviornment after taking the action.
62         next_state: The state of the enviornment after taking the action.
63         done: A boolean indicates whether the episode is done.
64
65     Returns:
66         None (Don't need to return anything)
67     """
68     # Begin your code
69     # TODO
70     # raise NotImplementedError("Not implemented yet.")
71
72     # calculate the new q-value with th (variable) learning_rate: float
73     self.qtable[state, action] += self.learning_rate * (
74         reward + self.gamma * np.max(self.qtable[next_state]) - self.qtable[state, action]
75     )
76
77     # End your code
78     if done:
79         np.save("./Tables/taxi_table.npy", self.qtable)
```

- choose_max_Q()

```

81  def check_max_Q(self, state):
82      """
83      - Implement the function calculating the max Q value of given state.
84      - Check the max Q value of initial state
85
86      Parameter:
87      state: the state to be check.
88      Return:
89      max_q: the max Q value of given state
90      """
91      # Begin your code
92      # TODO
93      # raise NotImplementedError("Not implemented yet.")
94      # to return the maximum Q
95      return np.max(self.qtable[state])
96
97      # End your code

```

Q-learning in Cartpole

- init_bins()

```

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      # use the linspace to slice the interval with given lower bound, upper bound and number of bins.
60      # return the list from index 1.
61      return np.linspace(lower_bound, upper_bound, num_bins, endpoint=False)[1:]
62
63      # End your code

```

- discretize_value()

```

65     def discretize_value(self, value, bins):
66         """
67         Discretize the value with given bins.
68         Parameters:
69             value: The value to be discretized.
70             bins: A numpy array of quantiles
71         returns:
72             The discretized value.
73         Example:
74             With given bins [2. 4. 6. 8.] and "5" being the value we're going to discretize.
75             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.
76         Hints:
77             1. This can be done with a numpy function.
78         """
79         # Begin your code
80         # TODO
81         # raise NotImplementedError("Not implemented yet.")
82
83         # Discretize the value with given bins.
84         # Using np.digitize() to determine the value in which interval of bins.
85         return np.digitize(value, bins, right=False)
86
87
88         # End your code

```

- discretize_observation()

```

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

```

- choose_action()

```

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 enviornment.
118             epsilon: Determines the explore/explot 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         if np.random.uniform(0, 1) > self.epsilon:
127             return np.argmax(self.qtable[tuple(state)])
128         else:
129             return env.action_space.sample()
130
131         # End your code

```

- learn()

```

133 ✓ def learn(self, state, action, reward, next_state, done):
134 ✓     """
135         Calculate the new q-value base on the reward and state transformation observed after taking the action.
136         Parameters:
137             state: The state of the enviornment before taking the action.
138             action: The exacuted action.
139             reward: Obtained from the enviornment after taking the action.
140             next_state: The state of the enviornment after taking the action.
141             done: A boolean indicates whether the episode is done.
142         Returns:
143             None (Don't need to return anything)
144         """
145         # Begin your code
146         # TODO
147         # raise NotImplementedError("Not implemented yet.")
148
149         self.qtable[tuple(state)][action] += self.learning_rate * (
150             reward
151             + self.gamma * np.max(self.qtable[tuple(next_state)])
152             - self.qtable[tuple(state)][action]
153         )
154
155         # End your code
156         if done:
157             np.save("./Tables/cartpole_table.npy", self.qtable)

```

- check_max_Q()

```

159 def check_max_Q(self):
160     """
161     - Implement the function calculating the max Q value of initial state(self.env.reset()).
162     - Check the max Q value of initial state
163     Parameter:
164         self: the agent itself.
165         (Don't pass additional parameters to the function.)
166         (All you need have been initialized in the constructor.)
167     Return:
168         max_q: the max Q value of initial state(self.env.reset())
169     """
170     # Begin your code
171     # TODO
172     # raise NotImplementedError("Not implemented yet.")
173
174     # check the maximum value of Q
175     return np.max(self.qtable[tuple(self.discretize_observation(self.env.reset()))])
176
177     # End your code

```

DQN in Cartpole

- learn()

```

133     # Begin your code
134     # TODO
135     # raise NotImplementedError("Not implemented yet.")
136
137     # Use function "sample" defined in class "replay_buffer" to get the sampled data.
138     sample = self.buffer.sample(self.batch_size)
139     # convert these sampled data into tensor.
140     states = torch.tensor(np.array(sample[0]), dtype=torch.float)
141     actions = torch.tensor(sample[1], dtype=torch.long).unsqueeze(1)
142     rewards = torch.tensor(sample[2], dtype=torch.float)
143     next_states = torch.tensor(np.array(sample[3]), dtype=torch.float)
144     done = torch.tensor(sample[4], dtype=torch.bool)
145
146     # "q_eval" is predicted values from evaluate network which is extracted based on "action".
147     q_eval = self.evaluate_net(states).gather(1, actions)
148     # "q_next" is actual values from target network.
149     q_next = self.target_net(next_states).detach() * (~done).unsqueeze(-1)
150     # "q_target" is the expected Q-values obtained from the formula "reward + gamma * max(q_next)".
151     q_target = rewards.unsqueeze(-1) + self.gamma * q_next.max(1)[0].view(self.batch_size, 1)
152
153     # use nn.MSELoss() func. to evaluate the loss of q_eval and q_target
154     loss_func = nn.MSELoss()
155     loss = loss_func(q_eval, q_target)
156
157     # zero-out the gradients before doing backpropagation.
158     self.optimizer.zero_grad()
159     loss.backward()
160     # update the parameters.
161     self.optimizer.step()
162
163     # End your code
164     torch.save(self.target_net.state_dict(), "./Tables/DQN.pt")

```

- choose_action()

```

166     def choose_action(self, state):
167         """
168         - Implement the action-choosing function.
169         - Choose the best action with given state and epsilon
170         Parameters:
171             self: the agent itself.
172             state: the current state of the environment.
173             (Don't pass additional parameters to the function.)
174             (All you need have been initialized in the constructor.)
175         Returns:
176             action: the chosen action.
177         """
178         with torch.no_grad():
179             # Begin your code
180             # TODO
181             # raise NotImplementedError("Not implemented yet.")
182
183             if np.random.uniform(0, 1) > self.epsilon:
184                 action = torch.argmax(
185                     self.evaluate_net(torch.tensor(state, dtype=torch.float))
186                 ).item()
187             else:
188                 action = env.action_space.sample()
189
190             # End your code
191             return action

```

- check_max_Q().

```

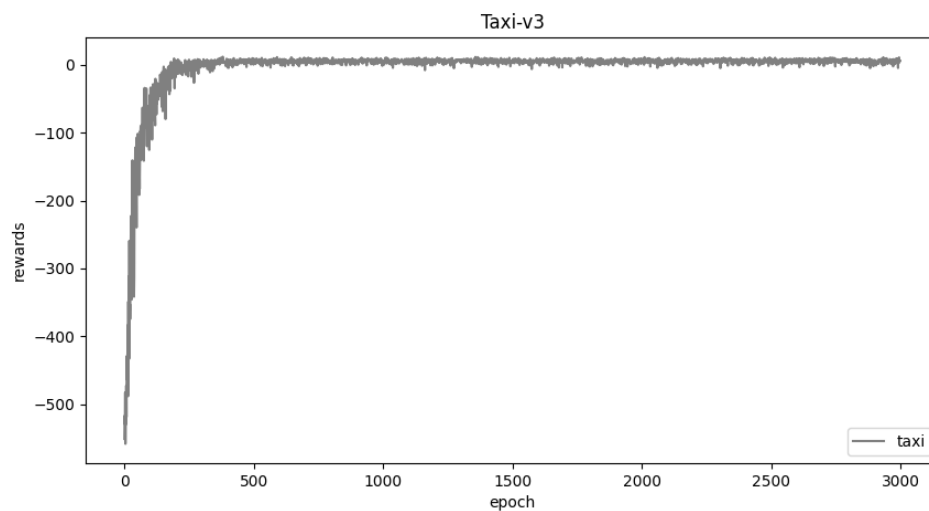
192 def check_max_Q(self):
193     """
194     - Implement the function calculating the max Q value of initial state(self.env.reset()).
195     - Check the max Q value of initial state
196     Parameter:
197         self: the agent itself.
198         (Don't pass additional parameters to the function.)
199         (All you need have been initialized in the constructor.)
200     Return:
201         max_q: the max Q value of initial state(self.env.reset())
202     """
203     # Begin your code
204     # TODO
205     # raise NotImplementedError("Not implemented yet.")
206
207     # change the initial state to tensor, forward to the NN.
208     x = torch.unsqueeze(torch.tensor(self.env.reset(), dtype=torch.float), 0)
209     # find out the max Q and return it.
210     return torch.max(self.target_net(x)).item()
211
212     # End your code

```

Part II. Experiment Results

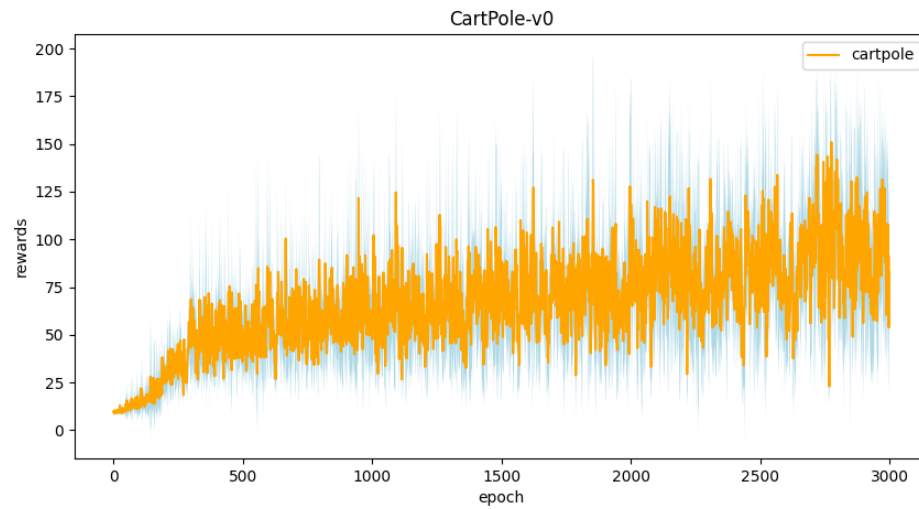
1. Q-learning in Taxi-v3

- taxi.png



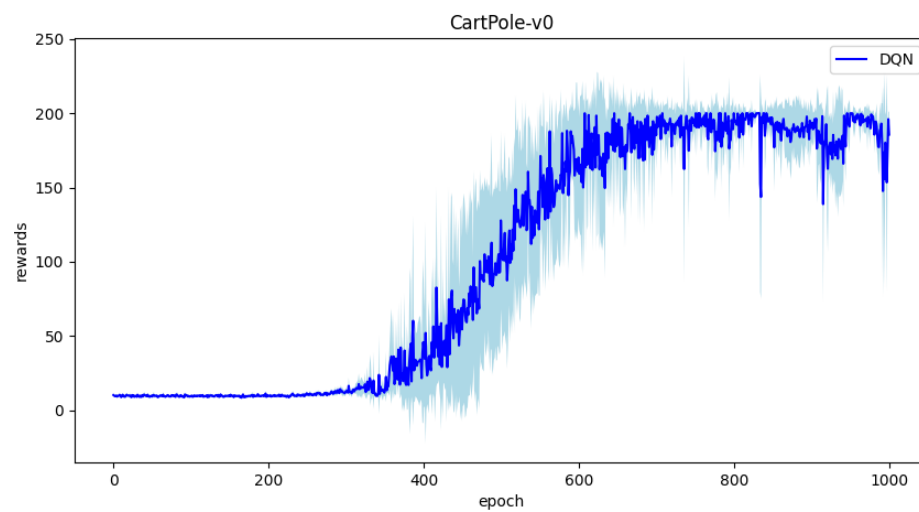
2. Q-learning in Cartpole-v0

- cartpole.png



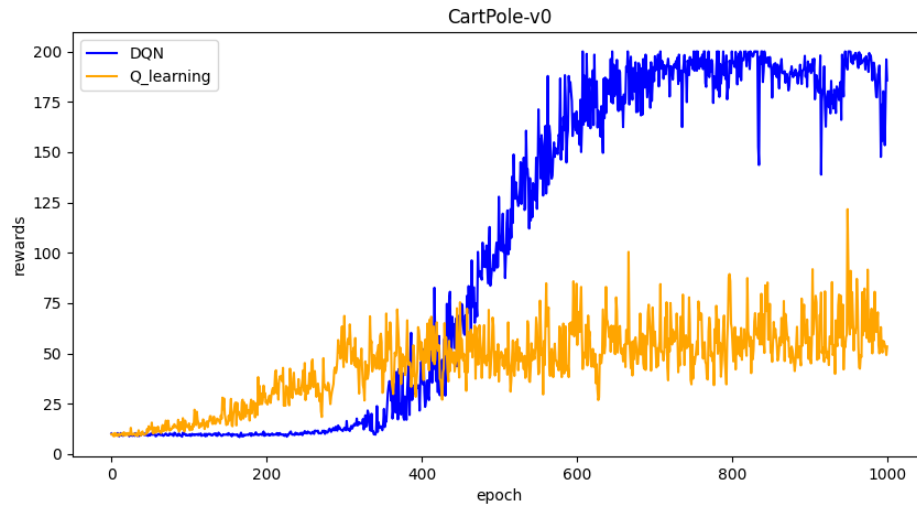
3. DQN in Cartpole-v0

- DQN.png



4. Compare Q-learning with DQN

- compare.png



Part III. Question Answering

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%)

A: Firstly, the taxi is at (2, 2), the passenger at Y, and the destination is at R, the reward are 9 steps of -1 and 1 step of 20.

optimal Q-value: $-1 * (1 - r^9) / 1 - r + 20 * (r^9) = -1 * (1 - 0.9^9) / 1 - 0.9 + 20(0.9^9) = 1.6226...$

```
PS C:\Users\user\Desktop\AI_HW4> python taxi.py
#1 training progress
100%|██████████| 3000/3000 [00:01<00:00, 1786.02it/s]
#2 training progress
100%|██████████| 3000/3000 [00:01<00:00, 1807.06it/s]
#3 training progress
100%|██████████| 3000/3000 [00:01<00:00, 1704.44it/s]
#4 training progress
100%|██████████| 3000/3000 [00:01<00:00, 1686.54it/s]
#5 training progress
100%|██████████| 3000/3000 [00:01<00:00, 1847.35it/s]
average reward: 7.81
Initail state:
taxi at (2, 2), passenger at Y, destination at R
max Q:1.6226146700000021
```

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%)

optimal q-value: $1-r^{(\text{avg} - \text{reward})/1-r}$

⇒ the result is around 33.3333...

The result of the maxQ of the DQN is closer to the optimal value than the Qmax of Q-learning.

```
PS C:\Users\user\Desktop\AI_HW4> python cartpole.py

#1 training progress
100%|██████████| 3000/3000 [00:07<00:00, 395.25it/s]
#2 training progress
100%|██████████| 3000/3000 [00:07<00:00, 381.98it/s]
#3 training progress
100%|██████████| 3000/3000 [00:07<00:00, 392.46it/s]
#4 training progress
100%|██████████| 3000/3000 [00:08<00:00, 369.58it/s]
#5 training progress
100%|██████████| 3000/3000 [00:08<00:00, 336.03it/s]
average reward: 120.48
max Q:30.630427737043938
```

```
PS C:\Users\user\Desktop\AI_HW4> python DQN.py

#1 training progress
100%|██████████| 1000/1000 [02:55<00:00, 5.69it/s]
#2 training progress
100%|██████████| 1000/1000 [03:09<00:00, 5.28it/s]
#3 training progress
100%|██████████| 1000/1000 [03:00<00:00, 5.54it/s]
#4 training progress
100%|██████████| 1000/1000 [03:54<00:00, 4.27it/s]
#5 training progress
100%|██████████| 1000/1000 [03:07<00:00, 5.34it/s]
reward: 200.0
max Q:34.19133377075195
```

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

A: Because the states are continuous, it makes Q-learning complex. We have to discretize the observation first.

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

A: In my opinion, the performance may become better, since the discretized data is more similar for observation.

c. Is there any concern if we increase “num_bins”? (3%)

A: “Oversampling” may cause the complexity of the computation.

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

A:

The DQN perform better.

Here is the reasons:

- a. DQN can deal with the continuous data; Q-learning deal with discretized data.
- b. Since the neural network, DQN can relieve the problem .
- c. DQNs, being more complex, might be prone to overfitting the training data in such a small problem setting, while discretized Q-learning remains straightforward and resilient to overfitting.

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

A: the epsilon-greedy approach allows the agent to explore new actions while still primarily relying on the best-known actions, thus helping it adaptively improve its policy over time.

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

A: Without exploration, the agent will take much longer to discover the optimal policy because it only tries a limited set of actions. It may cause poor performance.

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

A: It is possible. Since it is some alternative algorithm can be used.

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

A: During testing, the focus is on exploiting the learned policy to its fullest extent, which means relying on the highest-value actions identified during training, rather than exploring potentially less optimal ones.

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

In this step, we have to choose the next action which is estimated by neural network, and we don't need to calculate the gradient.

The torch.no_grad() can optimize the memory usage, speed up computation and ensure the correctness.