

Homework 4:

Reinforcement Learning

110550080何曉嫻

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

Part 1:

```
40 # Begin your code
41 if np.random.uniform(0,1) > self.epsilon:
42     action= np.argmax(self.qtable[state]) #exploitation
43 else:
44     action= env.action_space.sample()      # exploration
45 return action
46 raise NotImplementedError("Not implemented yet.")
47 # End your code
```

```
63 # Begin your code
64 # calculate q value and update the qtable
65 self.qtable[state, action]= self.qtable[state, action]+ self.learning_rate*\
66     (reward+ (self.gamma* np.max(self.qtable[next_state]))- self.qtable[state, action])
67
68 # End your code
```

```
81 # Begin your code
82 max_q= max(self.qtable[state])
83 return max_q
84 raise NotImplementedError("Not implemented yet.")
85 # End your code
```

Part 2:

```
55 # Begin your code
56 # create an list of arithmetic sequence(include lower_bound and upper_bound)
57 bins= np.linspace(lower_bound, upper_bound, num_bins+1)
58 # delete the head and tail
59 bins= np.delete(bins, num_bins)
60 bins= np.delete(bins, 0)
61 return bins
62 raise NotImplementedError("Not implemented yet.")
63 # End your code
```

```
79 # Begin your code
80 # discretize the value, find the value is in which interval
81 x= np.searchsorted(bins, value, side='right')
82 return x
83 raise NotImplementedError("Not implemented yet.")
84 # End your code
```

```

102     # Begin your code
103     # Discretize the observation which we observed from a continuous state space.
104     # Observation is a list of 4 features.
105     state=[]
106     for i in range(len(observation)):
107         state.append(self.discretize_value(observation[i], self.bins[i]))
108     return tuple(state)
109     raise NotImplementedError("Not implemented yet.")
110     # End your code

```

```

121     # Begin your code
122     if np.random.uniform(0,1) > self.epsilon:
123         action= np.argmax(self.qtable[state]) # exploitation
124     else:
125         action= env.action_space.sample()      # exploration
126     return action
127     raise NotImplementedError("Not implemented yet.")
128     # End your code

```

```

142     # Begin your code
143     # calculate q value and update the qtable (using tuple)
144     self.qtable[state+ (action,)] = self.qtable[state+ (action,)] + self.learning_rate*\
145         (reward+ (self.gamma* np.max(self.qtable[next_state])) - self.qtable[state+ (action,)] )
146
147     # raise NotImplementedError("Not implemented yet.")
148     # End your code

```

```

162     # Begin your code
163     # Discretize and find the max q
164     max_q= max(self.qtable[self.discretize_observation(self.env.reset())])
165     return max_q
166     raise NotImplementedError("Not implemented yet.")
167     # End your code

```

Part 3:

```

170     # Begin your code
171     if random.uniform(0,1) > self.epsilon:
172         # exploitation
173         action= torch.argmax(self.evaluate_net.forward(torch.FloatTensor(state))).item()
174     else:
175         # exploration
176         action = self.env.action_space.sample()
177     return action
178     raise NotImplementedError("Not implemented yet.")
179     # End your code

```

```

133     # Begin your code
134     # Sample trajectories of batch size from the replay buffer.
135     states, actions, rewards, next_states, done= self.buffer.sample(self.batch_size)
136     states= torch.FloatTensor(np.array(states))
137     actions= torch.LongTensor(actions)
138     rewards= torch.FloatTensor(rewards)
139     next_states= torch.FloatTensor(np.array(next_states))
140     # Forward the data to the evaluate net and the target net.
141     q_eval= self.evaluate_net(states).gather(1, actions.reshape(self.batch_size, 1))
142     q_next= self.target_net(next_states).detach()
143     for i in range(len(done)):
144         if (done[i]): q_next[i]= 0
145     q_target= rewards.reshape(self.batch_size, 1)+ self.gamma* q_next.max(1)[0].view(self.batch_size, 1)
146     # Compute the loss with MSE
147     func= nn.MSELoss()
148     loss= func(q_eval, q_target)
149     # Zero-out the gradients
150     self.optimizer.zero_grad()
151     # Backpropagation
152     loss.backward()
153     # Optimize the loss function
154     self.optimizer.step()
155     # End your code

```

```

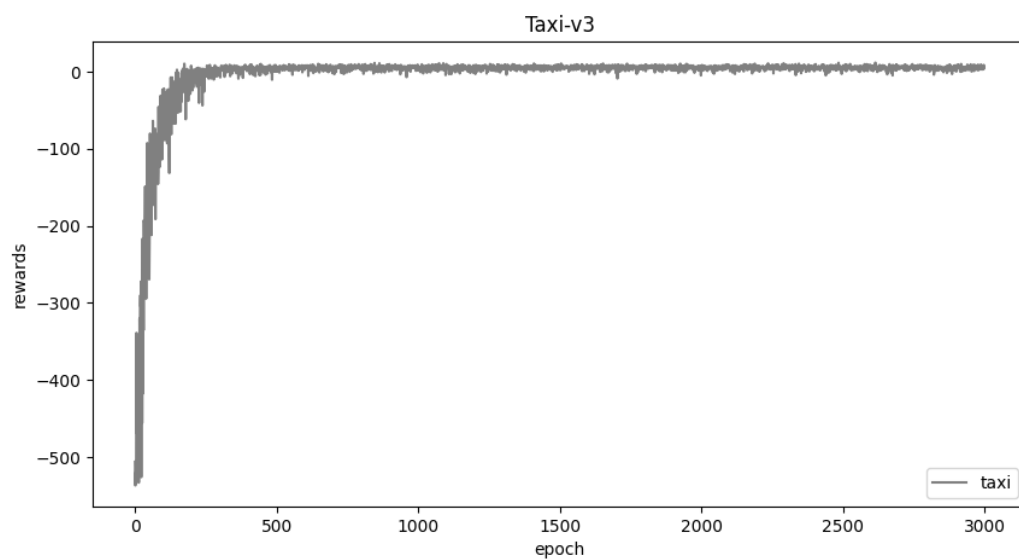
194     # Begin your code
195     # Use target network to estimate the maximum Q-value of the actions in the next state.
196     return torch.max(self.target_net(torch.FloatTensor(self.env.reset()))))
197     raise NotImplementedError("Not implemented yet.")
198     # End your code

```

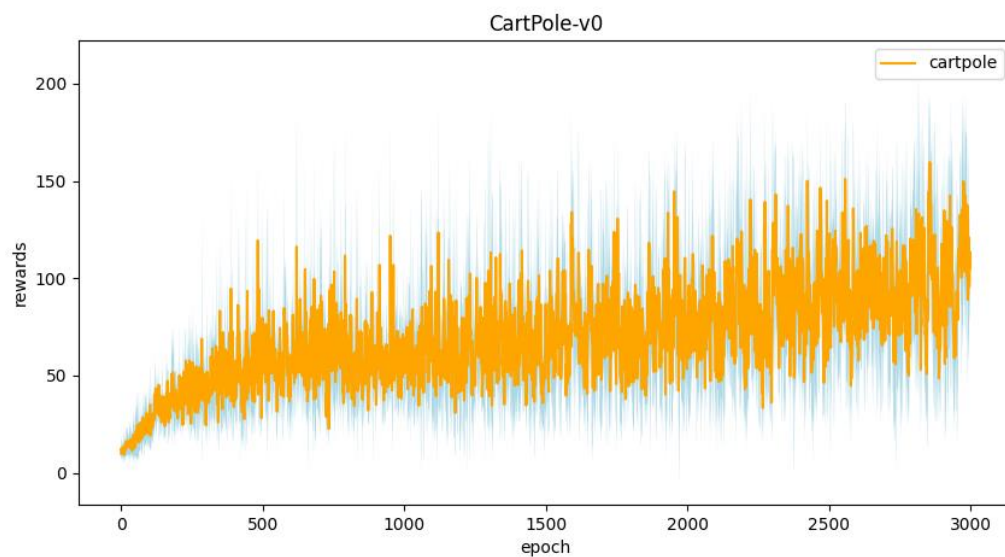
Part II. Experiment Results:

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

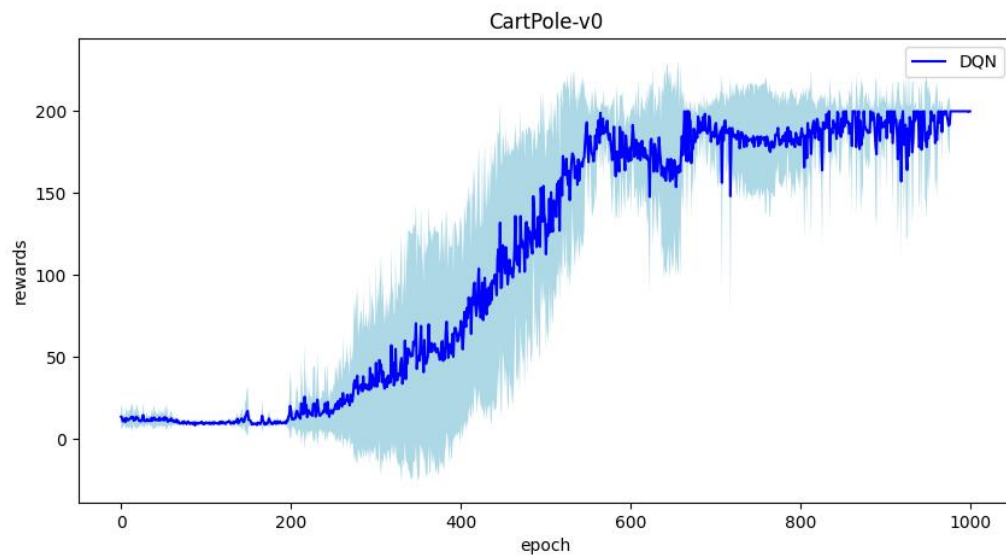
1. taxi.png



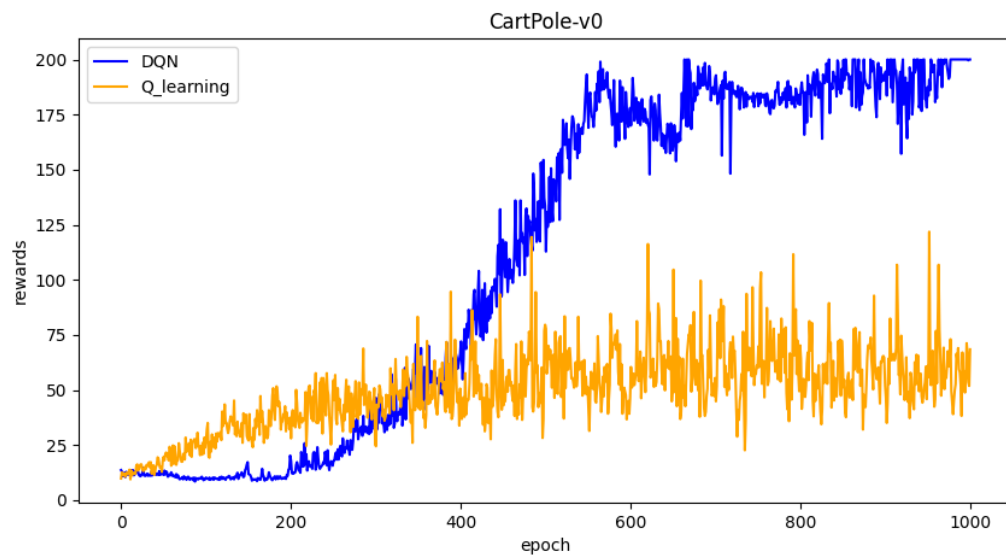
2. cartpole.png



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%)

```
PS C:\Users\lydia\OneDrive\桌面\AI\AI_HW4> python taxi.py
#1 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [00:24<00:00, 123.10it/s]
#2 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [00:24<00:00, 123.02it/s]
#3 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [00:24<00:00, 122.13it/s]
#4 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [00:25<00:00, 117.87it/s]
#5 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [00:25<00:00, 119.85it/s]
average reward: 7.69
Initail state:
taxi at (2, 2), passenger at Y, destination at R
max Q:1.6226146700000021
```

action : $\leftarrow \leftarrow \downarrow \downarrow$ pick $\uparrow \uparrow \uparrow \uparrow$ drop $\gamma=0.9$

optimal Q-value = $(-1) * (1 - 0.9^9) / (1 - 0.9) + 20 * 0.9^9 = 1.62$

The value is close !

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%)

```
PS C:\Users\lydia\OneDrive\桌面\AI\AI_HW4> python cartpole.py
#1 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [01:54<00:00, 26.18it/s]
#2 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [01:53<00:00, 26.52it/s]
#3 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [01:46<00:00, 28.04it/s]
#4 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [02:04<00:00, 24.18it/s]
#5 training progress
100%|██████████████████████████████████████████████████████████████████████████████| 3000/3000 [02:10<00:00, 22.96it/s]
average reward: 181.62
max Q:31.144320018134984
```

$\gamma=0.97$ reward=181.62

$$\text{optimal Q-value} = (1 - 0.97^{181.62}) / (1 - 0.97) = 33.2013$$

The value is close !

3.

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

The original observations are continuous. However, Q-learning works in a finite environment, so we need to discretize them.

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

The more bins would help the model account for more specific state space, and the performance will be better.

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

More bins would require training with more time, and the computation needs more memory.

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

DQN.

(1) DQN is capable of learning continuous state-action spaces without discretizing them.

(2) DQN uses experience replay and a target network to improve the stability of the learning process and prevent overfitting. Experience replay allows the agent to learn from a diverse set of experiences, while the target network is periodically updated to reduce the correlation between the target and estimated Q-values.

5.

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

The epsilon greedy algorithm selects the action with the highest estimated reward most of the time. The aim is to have a balance between exploration and exploitation.

With a small probability of epsilon, we choose to explore (not to exploit what we have learned so far). In this case, the action is selected randomly, independent of the action-value estimates.

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

The exploration and exploitation will not be balanced. It may cause extreme scenarios, like not exploiting what it has learned or not exploring the new things.

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, because we just need a method to balance exploitation and exploration, and there are many action selection methods. For example, the softmax action selection strategy controls the relative levels of exploration and exploitation by mapping values into action probabilities.

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

The agent has already got the information of the environment.

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

Operations performed within `with torch.no_grad():` will not be tracked for backpropagation, which can save memory and computational resources. This disables gradient calculation during the forward pass through the neural network.