

# Homework 4:

## Reinforcement Learning

### Report Template

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

Part 1:

Choose Action

```
40         # Begin your code
41
42         r = np.random.uniform(0,1) # choose a random r from 0~1
43         if(r > self.epsilon): # choose best action from the current state
44             action = np.argmax(self.qtable[state])
45         else: # choose a random action
46             action = env.action_space.sample()
47
48         return action
49
50         # End your code
```

Learn

```
68         # Begin your code
69
70         if done: max_next_Qopt = 0
71         else: max_next_Qopt = np.max(self.qtable[next_state])
72
73         Qopt_old = self.qtable[state, action]
74         # Qopt(s, a) = (1-eta) * Qopt(s, a) + eta * (reward + gamma * max_Qopt(s', a'))
75         Qopt_new = (1 - self.learning_rate) * Qopt_old + self.learning_rate * (reward + self.gamma * max_next_Qopt)
76         self.qtable[state, action] = Qopt_new # update qtable to new Qopt
77
78         # End your code
```

Check Max Q

```
92         # Begin your code
93
94         return np.max(self.qtable[state]) # return the max Q value of given state
95
96         # End your code
```

## Part 2:

### Init Bins

```
55     # Begin your code
56
57     # Returns num evenly spaced samples, calculated over the interval [lower, upper] without the upper point.
58     bins = np.linspace(lower_bound, upper_bound, num_bins, endpoint=False)
59     bins = np.delete(bins, 0) # Delete the lower point --> num-1 evenly spaced samples.
60
61     return bins
62
63     # End your code
```

### Discretized Values

```
79     # Begin your code
80
81     # Return the indices of the bins to which each value in input array belongs.
82     return np.digitize(value, bins, right=False) # bins[i-1] ≤ value < bins[i], return i
83
84     # End your code
```

### Discretize Observation

```
102     # Begin your code
103
104     # state = [0~6, 0~6, 0~6, 0~6] because there are 7 bins for each observation value
105     state = [0, 0, 0, 0]
106
107     for i in range(len(observation)):
108         # discretize each observation value accordingly to put each value in one of the 7 bins
109         state[i] = self.discretize_value(observation[i], self.bins[i])
110
111     return state
112
113     # End your code
```

### Choose Action

```
124     # Begin your code
125
126     r = np.random.uniform(0,1) # choose a random r from 0~1
127     if(r > self.epsilon): # choose best action from the current state
128         action = np.argmax(self.qtable[state[0], state[1], state[2], state[3]])
129     else: # choose a random action
130         action = env.action_space.sample()
131
132     return action
133
134     # End your code
```

## Learn

```
148     # Begin your code
149
150     # max_Qopt(s', a')
151     if done: max_next_Qopt = 0
152     else: max_next_Qopt = np.max(self.qtable[next_state[0], next_state[1], next_state[2], next_state[3]])
153
154     Qopt_old = self.qtable[state[0], state[1], state[2], state[3], action]
155     # Qopt(s, a) = (1-eta) * Qopt(s, a) + eta * (reward + gamma * max_Qopt(s', a'))
156     Qopt_new = (1 - self.learning_rate) * Qopt_old + self.learning_rate * (reward + self.gamma * max_next_Qopt)
157     self.qtable[state[0], state[1], state[2], state[3], action] = Qopt_new # update qtable to new Qopt
158
159     # End your code
```

## Check Max Q

```
173     # Begin your code
174
175     state = self.env.reset() # reinitialize state
176     state = self.discretize_observation(state) # discretize the initial state values
177
178     return np.max(self.qtable[state[0], state[1], state[2], state[3]]) # return the max Q value of the initial state
179
180     # End your code
```

## Part 3:

## Learn:

```
133     # Begin your code
134
135     # sample trajectories of batch size from the replay buffer
136     sample = self.buffer.sample(self.batch_size)
137     state = torch.tensor(np.array(sample[0]), dtype=torch.float)
138     action = torch.tensor(sample[1], dtype=torch.long)
139     reward = torch.tensor(sample[2], dtype=torch.float)
140     next_state = torch.tensor(np.array(sample[3]), dtype=torch.float)
141     done = sample[4]
142
143     # forward the data to the evaluate net and the target net
144     m = torch.LongTensor([0]*32)
145     eval = self.evaluate_net(state).gather(1, action.unsqueeze(1)) + m
146     # m = torch.LongTensor([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0])
147     # 0, 0, 0, 0, 0, 0, 0, 0])
148     # eval = eval + m
149     next = self.target_net(next_state).detach()
150     for i in range (self.batch_size) :
151         if done[i]: next[i] = 0
152     target = reward + self.gamma * next.max(1).values.unsqueeze(-1)
153
154     # compute loss with MSE
155     loss_func = nn.MSELoss()
156     loss = loss_func(eval.float(), target.float())
157
158     # zero-out the gradients
159     self.optimizer.zero_grad()
160
161     # backpropagation
162     loss.backward()
163     self.optimizer.step()
164
165     # End your code
```

## Choose Action

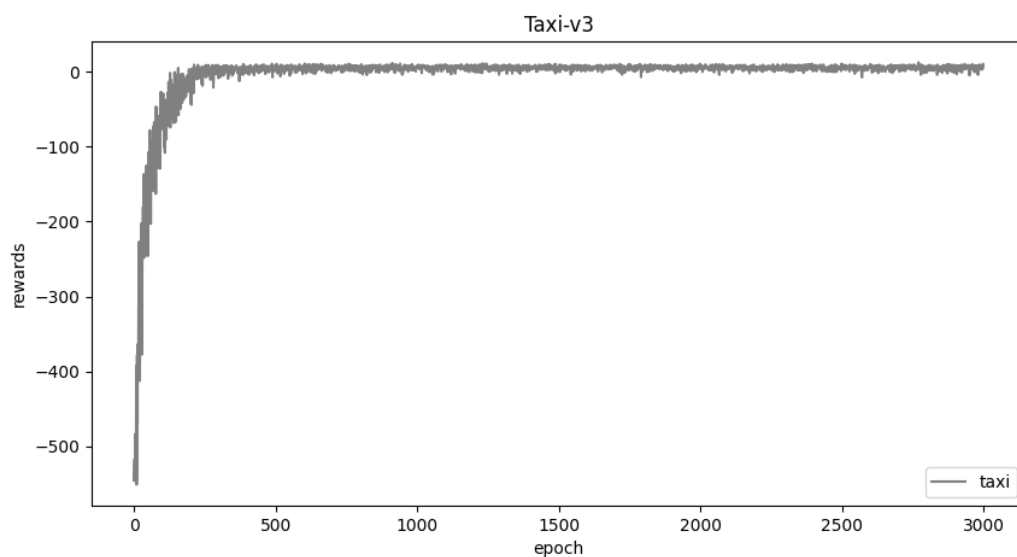
```
181         # Begin your code
182
183         x = torch.unsqueeze(torch.tensor(state, dtype=torch.float), 0)
184         r = np.random.uniform(0,1) # choose a random r from 0~1
185         if(r > self.epsilon): # choose best action from the current state
186             action_values = self.evaluate_net(x)
187             action = torch.argmax(action_values).item()
188         else: # choose a random action
189             action = env.action_space.sample()
190
191         # End your code
```

## Check Max Q

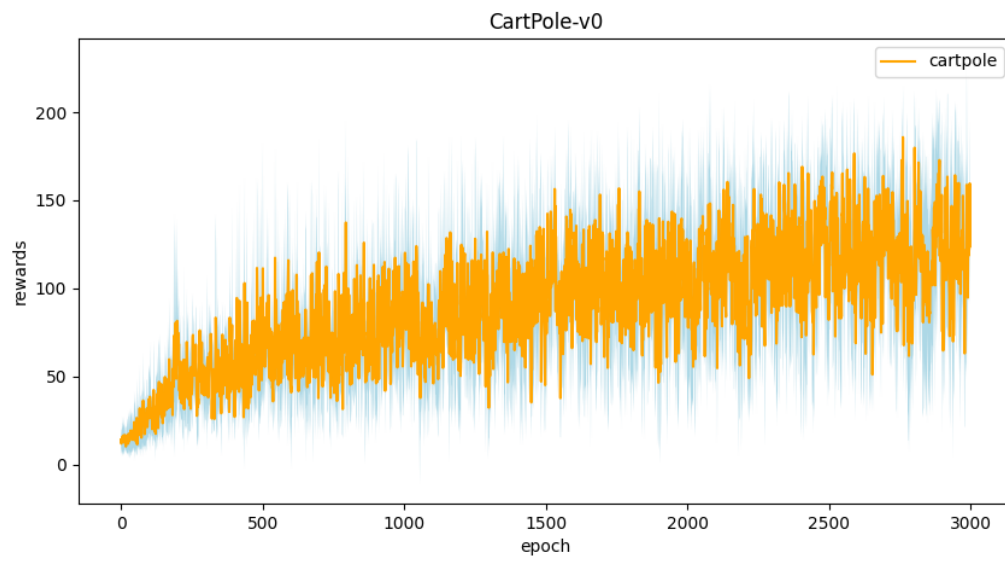
```
205         # Begin your code
206
207         state = self.env.reset() # reinitialize state
208         x = torch.unsqueeze(torch.tensor(state, dtype=torch.float), 0)
209         action_values = self.evaluate_net(x) # evaluate state to get Q values
210
211         return torch.max(action_values).item() # return max Q value
212
213         # End your code
```

## Part II. Experiment Results:

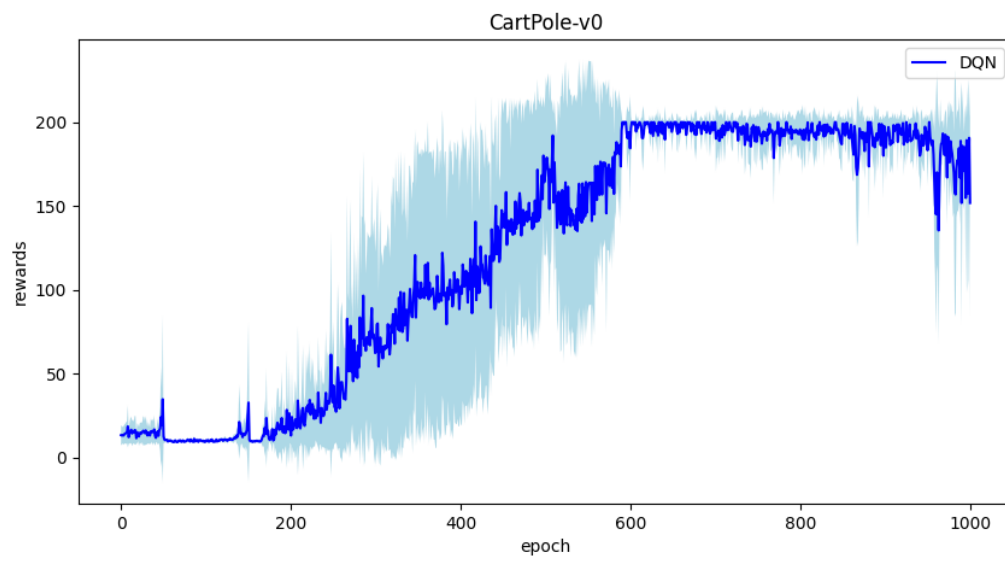
### 1. taxi.png:



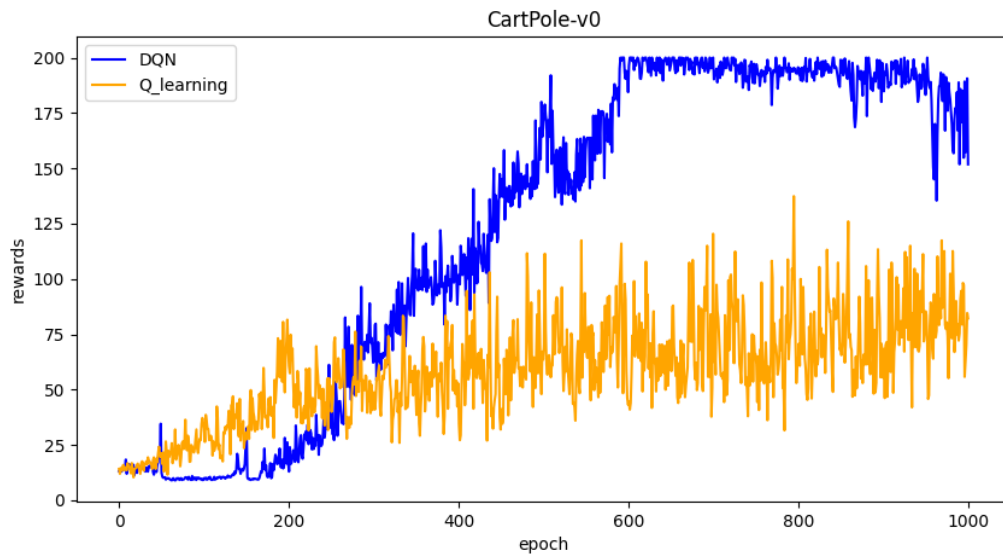
## 2. cartpole.png



## 3. DQN.png:



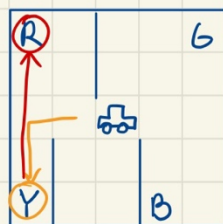
4. compare.png:



### 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



Taxi at (2,2)  
Passenger at Y  
Destination at R

5 steps to pick up (move  $\times 4$  + pickup  $\times 1$ )  
5 steps to drop off (move  $\times 4$  + dropoff  $\times 1$ )  
9 rewards of (-1) and 1 reward of (20)

$$Q_{\text{opt}} = 0.9^0(-1) + 0.9^1(-1) + 0.9^2(-1) + \dots + 0.9^8(-1) + 0.9^9(20) \quad \text{successful dropoff}$$

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

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%)

**max Q: 30.86389868745296**

Handwritten calculation on a grid background:

$$Q_{\text{opt}} = 0.99^0 + 0.99^1 + 0.99^2 + \dots + 0.99^{100}$$

Annotations: "gamma = 0.99" and "max steps taken" with an arrow pointing to the exponent 100.

$$= \left( \frac{1 - 0.99^{100}}{1 - 0.99} \right) = 33.2579$$

3.
  - a. Why do we need to discretize the observation in Part 2? (3%)  
Because the observation is continuous, we have discretize them to put them into separate bins to be able to create a proper qtable.
  - b. How do you expect the performance will be if we increase “num\_bins” ? (3%)  
The performance would become better since there will be more qtable states that fit the data more properly.
  - c. Is there any concern if we increase “num\_bins” ? (3%)  
Increasing the num\_bins would also increase the space and memory to save the increasing states in the qtable. Calculation time would also become longer as it takes more effort to go thorough all the 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. The reason is stated in its name, discretized. Q learning discretizes its data to fit them onto a qtable. This increases efficiency but lowers performance since the data is not so accurate anymore after discretizing. DQN, on the other hand, uses raw, continuous data. It uses more training time to obtain higher performance.
5.
  - a. What is the purpose of using the epsilon greedy algorithm while choosing an action? (3%)  
To maintain a balance in exploitation and exploration. When exploiting, choose the max(Q) action of current state from the constructed qtable. When exploring, choose a random action from the action space. This way, we can achieve both exploitation and exploration with balance.
  - b. What will happen, if we don't use the epsilon greedy algorithm in the CartPole-v0 environment? (3%)

All explore: the result would be random and there will be no learning curve.  
All exploit: can only work on known information and may miss unknown high-rewarded results.

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

If we can find another method to achieve balance in exploration and exploitation, then yes. Else, no, as reasons stated in (b).

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

The purpose of the epsilon greedy algorithm is to balance exploration and exploitation when constructing the qtable. When testing, we will only be using the qtable to find the best result without overwriting it. So we don't need the algorithm during testing.

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

There's no need to calculate the gradient and back propagation in `choose_action`, so we disable gradient calculation to increase performance with the function `torch.no_grad()`.