

Homework 4:

Reinforcement Learning

Report Template

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

- taxi.py

```
1 def choose_action(self, state):
2     """
3     Choose the best action with given state and epsilon.
4
5     Parameters:
6         state: A representation of the current state of the environment.
7         epsilon: Determines the explore/exploit rate of the agent.
8
9     Returns:
10        action: The action to be evaluated.
11    """
12    # Begin your code
13    if (np.random.uniform(0, 1) > self.epsilon): # exploitation
14        return np.argmax(self.qtable[state])
15    else: # exploration
16        return self.env.action_space.sample()
17    # End your code
```

```
1 def learn(self, state, action, reward, next_state, done):
2     """
3     Calculate the new q-value base on the reward and state transformation observed after taking the action.
4
5     Parameters:
6         state: The state of the environment before taking the action.
7         action: The executed action.
8         reward: Obtained from the environment after taking the action.
9         next_state: The state of the environment after taking the action.
10        done: A boolean indicates whether the episode is done.
11
12    Returns:
13        None (Don't need to return anything)
14    """
15    # Begin your code
16    # Q-learning algorithm
17    self.qtable[state, action] = (1 - self.learning_rate) * self.qtable[state, action] + self.learning_rate * (reward + self.gamma * np.max(self.qtable[next_state]))
18    # End your code
19    if done:
20        np.save("../Tables/taxi_table.npy", self.qtable)
```

```
1 def check_max_Q(self, state):
2     """
3     - Implement the function calculating the max Q value of given state.
4     - Check the max Q value of initial state
5
6     Parameter:
7         state: the state to be check.
8     Return:
9         max_q: the max Q value of given state
10    """
11    # Begin your code
12    # return max Q
13    return np.max(self.qtable[state])
14    # End your code
```

- cartpole.py

```
1 def init_bins(self, lower_bound, upper_bound, num_bins):
2     """
3     Slice the interval into #num_bins parts.
4     Parameters:
5         lower_bound: The lower bound of the interval.
6         upper_bound: The upper bound of the interval.
7         num_bins: Number of parts to be sliced.
8     Returns:
9         a numpy array of #num_bins - 1 quantiles.
10    Example:
11        Let's say that we want to slice [0, 10] into five parts,
12        that means we need 4 quantiles that divide [0, 10].
13        Thus the return of init_bins(0, 10, 5) should be [2. 4. 6. 8.].
14    Hints:
15        1. This can be done with a numpy function.
16    """
17    # Begin your code
18    return np.linspace(lower_bound, upper_bound, num_bins, endpoint=False)[1:]
19    # End your code
```

```
1 def discretize_value(self, value, bins):
2     """
3     Discretize the value with given bins.
4     Parameters:
5         value: The value to be discretized.
6         bins: A numpy array of quantiles
7     returns:
8         The discretized value.
9     Example:
10        With given bins [2. 4. 6. 8.] and "5" being the value we're going to discretize.
11        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.
12    Hints:
13        1. This can be done with a numpy function.
14    """
15    # Begin your code
16    return np.digitize(value, bins)
17    # End your code
```

```
1 def discretize_observation(self, observation):
2     """
3     Discretize the observation which we observed from a continuous state space.
4     Parameters:
5         observation: The observation to be discretized, which is a list of 4 features:
6             1. cart position.
7             2. cart velocity.
8             3. pole angle.
9             4. tip velocity.
10    Returns:
11        state: A list of 4 discretized features which represents the state.
12    Hints:
13        1. All 4 features are in continuous space.
14        2. You need to implement discretize_value() and init_bins() first
15        3. You might find something useful in Agent.__init__()
16    """
17    # Begin your code
18    return tuple([self.discretize_value(obs, bin) for obs, bin in zip(observation, self.bins)])
19    # End your code
```

```

1 def choose_action(self, state):
2     """
3     Choose the best action with given state and epsilon.
4     Parameters:
5         state: A representation of the current state of the environment.
6         epsilon: Determines the explore/exploit rate of the agent.
7     Returns:
8         action: The action to be evaluated.
9     """
10    # Begin your code
11    if (np.random.uniform(0, 1) > self.epsilon): # exploitation
12        return np.argmax(self.qtable[state])
13    else: # exploration
14        return self.env.action_space.sample()
15    # End your code

```

```

1 def learn(self, state, action, reward, next_state, done):
2     """
3     Calculate the new q-value based on the reward and state transformation observed after taking the action.
4     Parameters:
5         state: The state of the environment before taking the action.
6         action: The executed action.
7         reward: Obtained from the environment after taking the action.
8         next_state: The state of the environment after taking the action.
9         done: A boolean indicates whether the episode is done.
10    Returns:
11        None (Don't need to return anything)
12    """
13    # Begin your code
14    # Get the max Q value of the next state
15    max_next = np.max(self.qtable[next_state])
16    if done: max_next = 0
17
18    # Get the Q value of the current state
19    original = self.qtable[state + (action,)]
20
21    # Update Q value with Q-learning
22    self.qtable[state + (action,)] = (1 - self.learning_rate) * original + self.learning_rate * (reward + self.gamma * max_next)
23
24    if done:
25        # End your code
26        np.save("../Tables/cartpole_table.npy", self.qtable)

```

- DQN.py

```
1 def check_max_Q(self):
2     """
3     - Implement the function calculating the max Q value of initial state(self.env.reset()).
4     - Check the max Q value of initial state
5     Parameter:
6         self: the agent itself.
7         (Don't pass additional parameters to the function.)
8         (All you need have been initialized in the constructor.)
9     Return:
10        max_q: the max Q value of initial state(self.env.reset())
11    """
12    # Begin your code
13
14    # Get the Q value of initial state
15    initial_state = torch.tensor(self.env.reset(), dtype=torch.float)
16    q_value = self.target_net(initial_state).detach()
17
18    # Get max Q value
19    max_q = q_value.max(0).values.unsqueeze(-1)
20    max_q = float(max_q[0])
21    return max_q
22    # End your code
```

```
1 def learn(self):
2     """
3     - Implement the learning function.
4     - Here are the hints to implement.
5     Steps:
6     -----
7     1. Update target net by current net every 100 times. (we have done this for you)
8     2. Sample trajectories of batch size from the replay buffer.
9     3. Forward the data to the evaluate net and the target net.
10    4. Compute the loss with MSE.
11    5. Zero-out the gradients.
12    6. Backpropagation.
13    7. Optimize the loss function.
14    -----
15    Parameters:
16        self: the agent itself.
17        (Don't pass additional parameters to the function.)
18        (All you need have been initialized in the constructor.)
19    Returns:
20        None (Don't need to return anything)
21    """
22    if self.count % 100 == 0:
23        self.target_net.load_state_dict(self.evaluate_net.state_dict())
24
25    # Begin your code
26    # Add attribute step to record how many steps have done
27    if not hasattr(self, "step"): self.step = 0
28
29    # 2. Sample trajectories of batch size from the replay buffer.
30    batch_state, batch_action, batch_reward, batch_next_state, batch_done = self.buffer.sample(self.batch_size)
31
32    action=torch.tensor(np.asarray(batch_action).reshape(len(batch_action), 1), dtype=torch.long)
33    reward=torch.tensor(np.asarray(batch_reward).reshape(len(batch_reward), 1), dtype=torch.float)
34
35    # 3. Forward the data to the evaluate net and the target net.
36    q_evaluate = self.evaluate_net(torch.tensor(np.asarray(batch_state), dtype=torch.float)).gather(1, action)
37    q_next = self.target_net(torch.tensor(np.asarray(batch_next_state), dtype=torch.float)).detach()
38    q_target = reward + self.gamma * q_next.max(1).values.unsqueeze(-1) # target Q value
39    for i in range(len(batch_done)):
40        if batch_done[i]:
41            q_target[i][0] = 0
42
43    # 4. Compute the loss with MSE.
44    loss_function = nn.MSELoss()
45    loss = loss_function(q_evaluate, q_target)
46
47    # 5. Zero-out the gradients.
48    self.optimizer.zero_grad()
49
50    # 6. Backpropagation.
51    loss.backward()
52
53    # 7. Optimize the loss function.
54    self.optimizer.step()
55
56    # Update count and step
57    self.count += 1
58    self.step += 1
59
60    # Save the result every 100 steps
61    if self.step % 100 == 0:
62        torch.save(self.target_net.state_dict(), "../Tables/DQN.pt")
63    # End your code
```

```

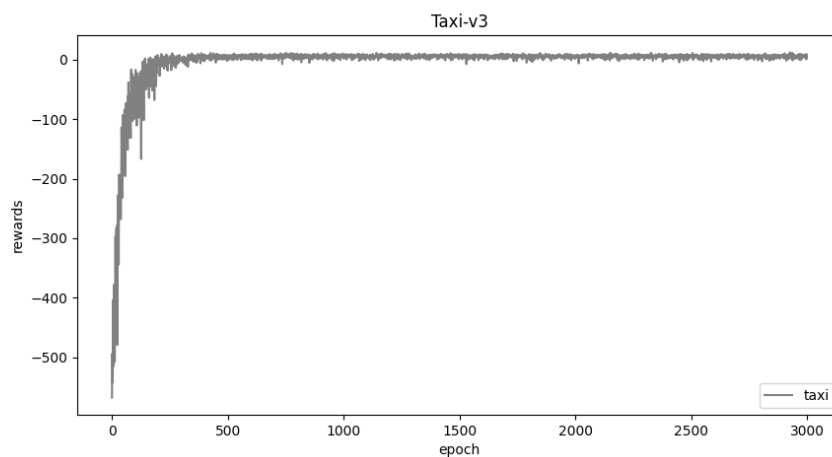
1 def choose_action(self, state):
2     """
3     - Implement the action-choosing function.
4     - Choose the best action with given state and epsilon
5     Parameters:
6         self: the agent itself.
7         state: the current state of the environment.
8         (Don't pass additional parameters to the function.)
9         (All you need have been initialized in the constructor.)
10    Returns:
11        action: the chosen action.
12    """
13    with torch.no_grad():
14        # Begin your code
15        # TODO
16        x = torch.tensor(state, dtype=torch.float)
17        if np.random.uniform(0, 1) > self.epsilon: # exploration
18            # Predict Q value from evaluate net
19            action_values = self.evaluate_net(x)
20            # Pick action that contains largest Q value
21            action = torch.argmax(action_values).item()
22        else: # exploitation
23            action = env.action_space.sample()
24        # End your code
25    return action

```

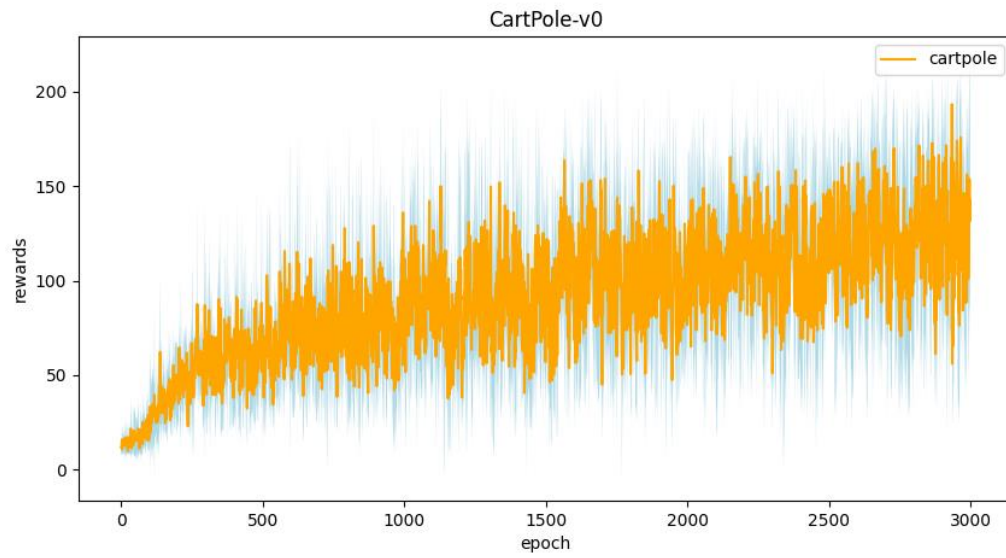
Part II. Experiment Results:

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

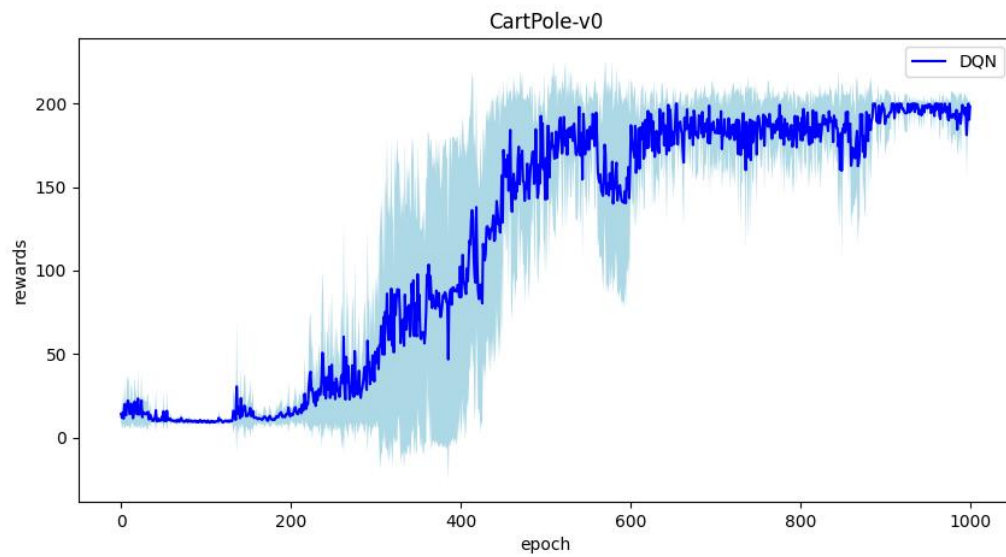
1. taxi.png:



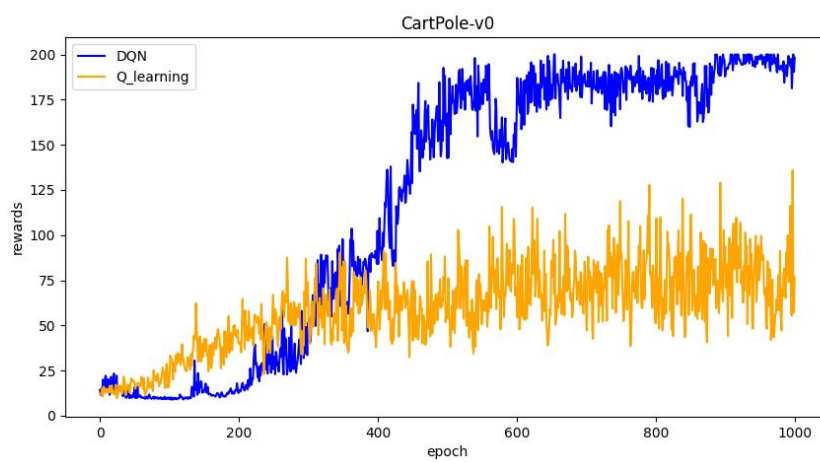
2. cartpole.png



3. DQN.png



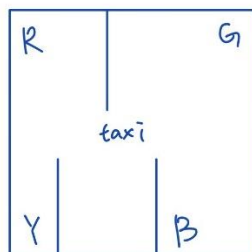
4. compare.png



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%)

```
average reward: 8.05
Initial state:
taxi at (2, 2), passenger at Y, destination at R
max Q:1.6226146699999995
```



passenger at Y, destination at R
gamma = 0.9

reward $\begin{cases} -1 & \text{per step unless other reward is triggered} \\ +20 & \text{delivering passenger} \\ -10 & \text{executing "pickup" and "drop-off" illegally.} \end{cases}$

5 steps to pickup the passenger
5 steps to drop-off the passenger

$$\begin{aligned} Q_{\text{opt}} &= -1 + 0.9 \cdot (-1) + (0.9)^2 \cdot (-1) + \dots + (0.9)^8 \cdot (-1) + (0.9)^9 \cdot 20 \\ &= (-1) \cdot \frac{1 - (0.9)^9}{1 - 0.9} + (0.9)^9 \cdot 20 \\ &= -10 \cdot (1 - (0.9)^9) + (0.9)^9 \cdot 20 \\ &= -10 + (0.9)^9 \cdot 30 \div 1.62261467 \end{aligned}$$

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%)

```
average reward: 92.48
max Q:29.309158145736305
```

Rewards: +1 for every step taken

Episode end: Truncation is executed for episode length greater than 200

gamma = 0.97

$$\begin{aligned} \Rightarrow Q_{\text{opt}} &= 1 + 1 \cdot (0.97)^1 + 1 \cdot (0.97)^2 + \dots + 1 \cdot (0.97)^{199} \\ &= \frac{1 - (0.97)^{200}}{1 - 0.97} = \frac{1 - (0.97)^{200}}{0.03} \div 33.2579586 \end{aligned}$$

3.

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

Because the interval is continuous, which is not easy to determine the state of the cartpole. As a result, we need to discretize it to get the state.

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

In my opinion, the performance will become better if we increase the num_bins, which represents the number of the states in the bounded interval. Because when we increase the number of the states, it implies that we have more states to approximate the continuous interval, which leads to the better performance.

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

Increasing "num_bins" can result in concerns such as longer time required to update and save the Q table due to the increased number of states, as well as the increased memory required to save the larger Q table.

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

DQN outperforms discretized Q learning in the Cartpole-v0 environment. The reason for this is that Q learning discretizes the continuous data into states, which can result in data loss. In contrast, DQN can use the continuous data and preserve more details, leading to better performance.

5.

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

The epsilon-greedy algorithm serves the purpose of balancing exploration and exploitation in action selection. It allows for the selection of the best-known action while also exploring new options, thus taking advantage of prior knowledge and discovering new possibilities.

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

If the epsilon-greedy algorithm is not used in the CartPole-v0 environment, there are two potential scenarios. If only exploration is used, there will be no way to choose the best-known action, as all actions will be random. If only exploitation is used, the algorithm can only rely on the known information and may miss unknown high-performance conditions without any randomness or 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%)

It may be possible to achieve the same performance without the epsilon-greedy algorithm in the CartPole-v0 environment if we can find another method to replace the learning rate in the algorithm while maintaining the same proportion of exploration/exploitation. With the same distribution, there is a possibility to achieve the same performance.

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

The epsilon-greedy algorithm is not needed during the testing section because it is used for exploration and exploitation during training. In the testing section, the goal is simply to find the best path and obtain the best reward, and therefore there is no need to use the algorithm to train the data.

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

The use of `"with torch.no_grad():"` inside the `"choose_action"` function in DQN disables gradient calculation for every tensor, meaning that `requires_grad` is set to `False`. This is because the function is used to choose an action and update the Q-values, but the gradients are not needed for these operations and can be disabled to improve computational efficiency.