

Assignment 4

余北辰 519030910245

1 Experimental Purpose

- Implement the DQN algorithm and its improved algorithm DDQN
- Play with them in the classical RL control environment MountainCar
- Compare the performance of DQN and DDQN

2 Experimental Background

In this section, we present the methods of DQN and DDQN separately. In addition, we present the background of the environment MountainCar.

2.1 DQN

DQN, or Deep Q-Network, approximates a state-value function in a Q-Learning framework with a neural network.

The DQN algorithm use two same network: policy network Q and target network \hat{Q} . In each episode, the agent interact with the environment using the given state s_t and an action a_t . The algorithm stores the reward r_t and next state s_{t+1} together with s_t and a_t into the buffer as (s_t, a_t, r_t, s_{t+1}) . Then, randomly sample a batch from the buffer, and estimate the target as $y = r + \gamma \max_a \hat{Q}(s, a)$. After that, optimize the mean square error between the estimation and Q value as $\min(y - Q(s, a))^2$. Finally, update \hat{Q} by setting $\hat{Q} := Q$ every C steps.

The pseudo-code for DQN is shown below:

Algorithm 1: deep Q-learning with experience replay.

Initialize replay memory D to capacity N

Initialize action-value function Q with random weights θ

Initialize target action-value function \hat{Q} with weights $\theta^- = \theta$

For episode = 1, M **do**

 Initialize sequence $s_1 = \{x_1\}$ and preprocessed sequence $\phi_1 = \phi(s_1)$

For $t = 1, T$ **do**

 With probability ϵ select a random action a_t

 otherwise select $a_t = \arg\max_a Q(\phi(s_t), a; \theta)$

 Execute action a_t in emulator and observe reward r_t and image x_{t+1}

 Set $s_{t+1} = s_t, a_t, x_{t+1}$ and preprocess $\phi_{t+1} = \phi(s_{t+1})$

 Store transition $(\phi_t, a_t, r_t, \phi_{t+1})$ in D

 Sample random minibatch of transitions $(\phi_j, a_j, r_j, \phi_{j+1})$ from D

 Set $y_j = \begin{cases} r_j & \text{if episode terminates at step } j+1 \\ r_j + \gamma \max_{a'} \hat{Q}(\phi_{j+1}, a'; \theta^-) & \text{otherwise} \end{cases}$

 Perform a gradient descent step on $(y_j - Q(\phi_j, a_j; \theta))^2$ with respect to the network parameters θ

 Every C steps reset $\hat{Q} = Q$

End For

End For

2.2 DDQN

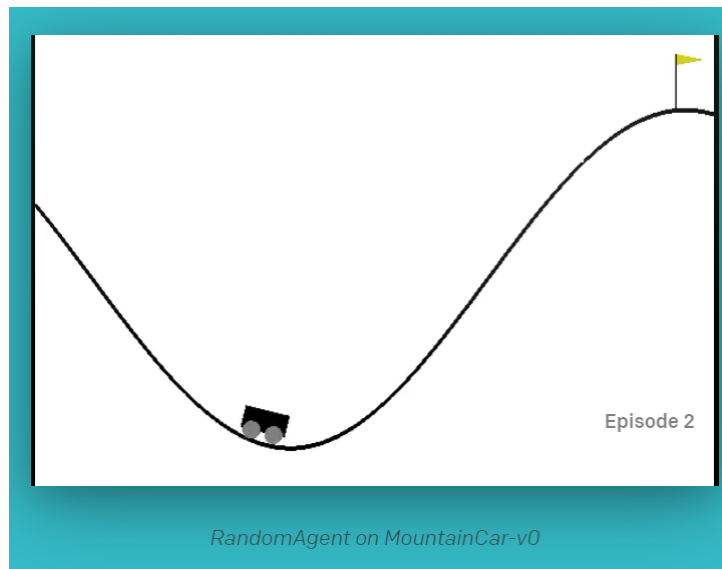
The DDQN(double DQN) is an improvement of DQN. To reduce overestimating of Q value, DDQN modified the method of estimating TD target as $y = r + \gamma \hat{Q}(s, \arg \max_a Q(s, a))$. If Q overestimates a , then \hat{Q} will give it a proper value; If \hat{Q} overestimates a , then a will not be selected. Thus, overestimation is avoided.

Algorithm 1 : Double Q-learning (Hasselt et al., 2015)

```
Initialize primary network  $Q_\theta$ , target network  $Q_{\theta'}$ , replay buffer  $\mathcal{D}$ ,  $\tau \ll 1$ 
for each iteration do
    for each environment step do
        Observe state  $s_t$  and select  $a_t \sim \pi(a_t, s_t)$ 
        Execute  $a_t$  and observe next state  $s_{t+1}$  and reward  $r_t = R(s_t, a_t)$ 
        Store  $(s_t, a_t, r_t, s_{t+1})$  in replay buffer  $\mathcal{D}$ 
    for each update step do
        sample  $e_t = (s_t, a_t, r_t, s_{t+1}) \sim \mathcal{D}$ 
        Compute target Q value:
             $Q^*(s_t, a_t) \approx r_t + \gamma Q_{\theta'}(s_{t+1}, \arg \max_{a'} Q_{\theta'}(s_{t+1}, a'))$ 
        Perform gradient descent step on  $(Q^*(s_t, a_t) - Q_\theta(s_t, a_t))^2$ 
        Update target network parameters:
             $\theta' \leftarrow \tau * \theta + (1 - \tau) * \theta'$ 
```

2.3 MountainCar

MountainCar is a classical gym environment. In MountainCar, A car is on a one-dimensional track, positioned between two "mountains". The goal is to drive up the mountain on the right; however, the car's engine is not strong enough to scale the mountain in a single pass.



There are three possible actions. 0 means push left, 1 means do nothing, and 2 means push right. The observation of the car is a vector of two dimension: the first dimension represents the position of the car and the second one represents the velocity.

3 Experimental Procedure

The source code for this experiment includes `DQN.py`, `DDQN.py` and `main.py`. The implement of DQN references the [blog](#) and is modified based on our experimental scenario.

3.1 DQN.py

DQN.py includes two classes: a class of network and a class of DQN agent.

Network

We will utilize `torch.nn` to construct our neural network for DQN.

```
1 class Net(nn.Module):
2     def __init__(self, ):
3         super(Net, self).__init__()
4         self.fc1 = nn.Linear(N_STATES, 50)
5         self.fc1.weight.data.normal_(0, 0.1)
6         self.out = nn.Linear(50, N_ACTIONS)
7         self.out.weight.data.normal_(0, 0.1)
8
9     def forward(self, x):
10        x = self.fc1(x)
11        x = F.relu(x)
12        actions_value = self.out(x)
13        return actions_value
```

DQN agent

The initialization of DQN agent is seen below.

```
1 class DQN(object):
2     def __init__(self):
3         self.eval_net, self.target_net = Net(), Net()
4
5         self.learn_step_counter = 0
6         self.memory_counter = 0
7         self.memory = np.zeros((MEMORY_CAPACITY, N_STATES * 2 + 2))
8         self.optimizer = torch.optim.Adam(self.eval_net.parameters(), lr=LR)
9         self.loss_func = nn.MSELoss()
```

As seen, we use two same network to fit policy and target separately.

```
1     def choose_action(self, x):
2         x = torch.unsqueeze(torch.FloatTensor(x), 0)
3         if np.random.uniform() < EPSILON:
4             actions_value = self.eval_net.forward(x)
5             action = torch.max(actions_value, 1)[1].data.numpy()
6             action = action[0] if ENV_A_SHAPE == 0 else action.reshape(
7                 ENV_A_SHAPE)
8         else:
9             action = np.random.randint(0, N_ACTIONS)
10            action = action if ENV_A_SHAPE == 0 else action.reshape(
11                ENV_A_SHAPE)
12            return action
```

We use a ϵ -greedy policy to choose action. If the random number is less than ϵ , the greedy action is chosen. Otherwise, a random action is chosen.

```

1 def store_transition(self, s, a, r, s_):
2     transition = np.hstack((s, [a, r], s_))
3     index = self.memory_counter % MEMORY_CAPACITY
4     self.memory[index, :] = transition
5     self.memory_counter += 1

```

We use a buffer to store the transitions. If the memory buffer is full, then the old transitions are filled.

```

1 def learn(self):
2     if self.learn_step_counter % TARGET_REPLACE_ITER == 0:
3         self.target_net.load_state_dict(self.eval_net.state_dict())
4     self.learn_step_counter += 1
5
6     sample_index = np.random.choice(MEMORY_CAPACITY, BATCH_SIZE)
7     b_memory = self.memory[sample_index, :]
8     b_s = torch.FloatTensor(b_memory[:, :N_STATES])
9     b_a = torch.LongTensor(b_memory[:, N_STATES:N_STATES+1].astype(int))
10    b_r = torch.FloatTensor(b_memory[:, N_STATES+1:N_STATES+2])
11    b_s_ = torch.FloatTensor(b_memory[:, -N_STATES:])
12
13    q_eval = self.eval_net(b_s).gather(1, b_a)
14    q_next = self.target_net(b_s_).detach()
15    q_target = b_r + GAMMA * q_next.max(1)[0].view(BATCH_SIZE, 1)
16    loss = self.loss_func(q_eval, q_target)
17
18    self.optimizer.zero_grad()
19    loss.backward()
20    self.optimizer.step()

```

`learn()` is the key function. A transition (s_t, a_t, r_t, s_{t+1}) is randomly sampled from the buffer. The loss is calculated as introduced in previous, and the parameters of the networks are updated.

3.2 DDQN.py

The most parts of `DDQN.py` is the same as `DQN.py`. The only difference is the `learn()` function.

```

1 def learn(self):
2     if self.learn_step_counter % TARGET_REPLACE_ITER == 0:
3         self.target_net.load_state_dict(self.eval_net.state_dict())
4     self.learn_step_counter += 1
5
6     sample_index = np.random.choice(MEMORY_CAPACITY, BATCH_SIZE)
7     b_memory = self.memory[sample_index, :]
8     b_s = torch.FloatTensor(b_memory[:, :N_STATES])
9     b_a = torch.LongTensor(b_memory[:, N_STATES:N_STATES+1].astype(int))
10    b_r = torch.FloatTensor(b_memory[:, N_STATES+1:N_STATES+2])
11    b_s_ = torch.FloatTensor(b_memory[:, -N_STATES:])
12
13    q_eval = self.eval_net(b_s).gather(1, b_a)
14    q_next = self.target_net(b_s_)
15    q_next_cur = self.eval_net(b_s_)
16    q_target = b_r + GAMMA * \
17        q_next.gather(1, torch.max(q_next_cur, 1)
18            [1].unsqueeze(1)).squeeze(1)

```

```

19         loss = self.loss_func(q_eval, q_target.detach())
20
21         self.optimizer.zero_grad()
22         loss.backward()
23         self.optimizer.step()

```

The code of the `learn()` function is modified based on the equation in the section 2.2.

3.3 main.py

```

1  def main():
2      parser = argparse.ArgumentParser()
3      parser.add_argument('-n', '--episode_num', type=int, default=500)
4      parser.add_argument('-a', '--action', type=str, default='plot')
5
6      args = parser.parse_args()
7      episode_num = args.episode_num
8      action = args.action
9
10     if action == 'plot':
11         plot(episode_num)
12     elif action == 'train':
13         train(episode_num)
14
15
16 if __name__ == '__main__':
17     main()

```

`main()` use parser to get the parameters entered on the command line, and depending on these parameters, training the models or drawing curves is chosen.

training

```

1  def train_net(net, episode_num):
2      return_list = []
3      step_list = []
4      score_list = []
5      for i in range(episode_num):
6          s = env.reset()
7          ep_r = 0
8          step = 0
9          score = 0
10         discount = 1.0
11         while True:
12             step += 1
13             # env.render()
14             a = net.choose_action(s)
15             s_, r, done, info = env.step(a)
16
17             if s_[0] <= -0.5:
18                 r = 100 * abs(s_[1])
19             elif s_[0] > -0.5 and s_[0] <= 0.5:
20                 r = math.pow(2, 5*(s_[0] + 1)) + (100 * abs(s_[1])) ** 2
21             else:
22                 r = 1000
23             net.store_transition(s, a, r, s_)

```

```

24
25         score += discount * r
26         discount *= GAMMA
27
28         ep_r += r
29         if net.memory_counter > MEMORY_CAPACITY:
30             net.learn()
31             if done:
32                 return_list.append(ep_r)
33                 step_list.append(step)
34                 score_list.append(score)
35             if done:
36                 break
37         s = s_
38     return return_list, step_list, score_list
39
40 def train(episode_num):
41     dqn = DQN()
42     ddqn = DDQN()
43     dqn_return_list, dqn_step_list, dqn_score_list = train_net(
44         dqn, episode_num)
45     ddqn_return_list, ddqn_step_list, ddqn_score_list = train_net(
46         ddqn, episode_num)
47     np.save('dqn_return_list', dqn_return_list)
48     np.save('dqn_step_list', dqn_step_list)
49     np.save('dqn_score_list', dqn_score_list)
50     np.save('ddqn_return_list', ddqn_return_list)
51     np.save('ddqn_step_list', ddqn_step_list)
52     np.save('ddqn_score_list', ddqn_score_list)

```

In the training part, The return, step number and score of each episode are stored in three lists. The lists are saved as the `.npy` files and will be used to plot the curves.

The origin reward is not good enough to train the model, so we modified the reward function. The chosen of the reward function references this [blog](#).

$$r_t = \begin{cases} 1000, & \text{if } x_t \geq 0.5 \\ 2^{5(x_t+1)} + (100|v_t|)^2, & \text{if } -0.5 < x_t < 0.5 \\ 0 + 100|v_t|, & \text{if } x_t \leq -0.5 \end{cases}$$

When the car coordinates $x < -0.5$, although the position is not good, in order to accelerate to the right, we still need to make sure that the absolute value of velocity is large.

When the car coordinate $x > -0.5$, at this time, the more rightward the position, the bigger the reward; and the bigger the absolute value of speed, the reward should be steeper.

plotting

```

1 def plot(episode_num):
2     dqn_return_list = np.load('dqn_return_list.npy', allow_pickle=True)
3     dqn_step_list = np.load('dqn_step_list.npy', allow_pickle=True)
4     dqn_score_list = np.load('dqn_score_list.npy', allow_pickle=True)
5     ddqn_return_list = np.load('ddqn_return_list.npy', allow_pickle=True)
6     ddqn_step_list = np.load('ddqn_step_list.npy', allow_pickle=True)
7     ddqn_score_list = np.load('ddqn_score_list.npy', allow_pickle=True)
8

```

```

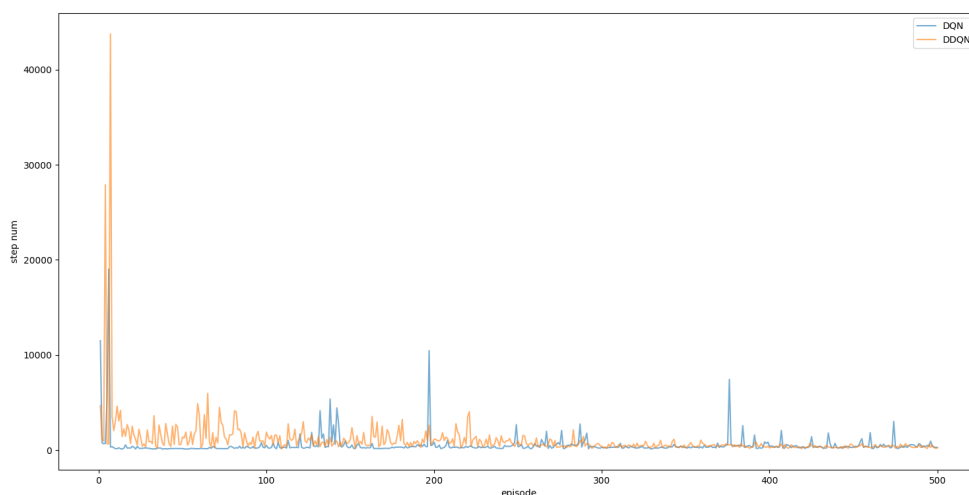
9
10     ep_range = [i + 1 for i in range(episode_num)]
11
12     plt.figure(figsize=(18, 9))
13     plt.plot(ep_range, dqn_step_list, label='DQN', alpha=0.6)
14     plt.plot(ep_range, ddqn_step_list, label='DDQN', alpha=0.6)
15     plt.xlabel('episode')
16     plt.ylabel('step num')
17     plt.legend()
18     plt.savefig('step.png')
19
20     plt.figure(figsize=(18, 9))
21     plt.plot(ep_range, dqn_return_list, label='DQN', alpha=0.6)
22     plt.plot(ep_range, ddqn_return_list, label='DDQN', alpha=0.6)
23     plt.xlabel('episode')
24     plt.ylabel('return')
25     plt.legend()
26     plt.savefig('return.png')
27
28     plt.figure(figsize=(18, 9))
29     plt.plot(ep_range, dqn_score_list, label='DQN', alpha=0.6)
30     plt.plot(ep_range, ddqn_score_list, label='DDQN', alpha=0.6)
31     plt.xlabel('episode')
32     plt.ylabel('score')
33     plt.legend()
34     plt.savefig('score.png')

```

4 Experimental Result

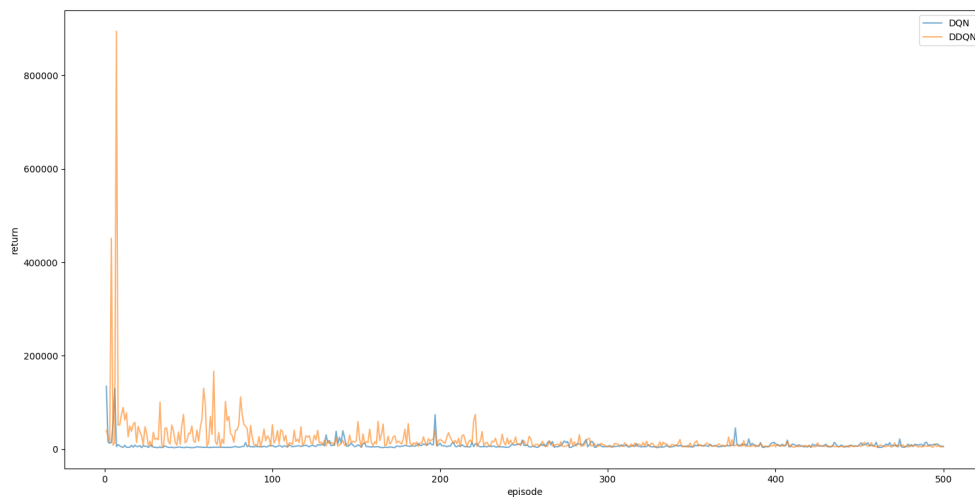
We can enable `env.render()` to observe the training of the car. We can find that both in DQN and DDQN, the car can easily get into the position of flag in just several episodes.

The curve of the step number:

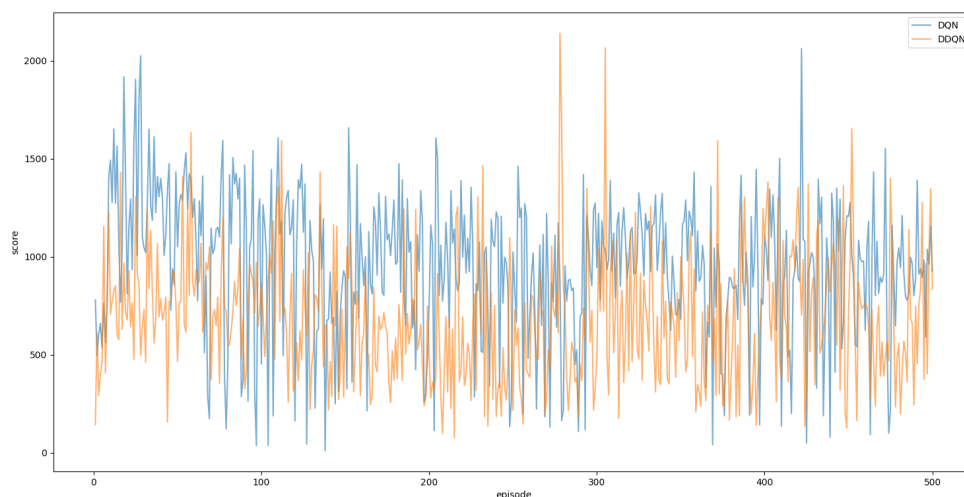


DQN converges quickly while DDQN converges in around 250 episodes.

The curve of the total return:



The curve of the score:



From the above figures, we can conclude that DDQN is indeed more stable than DQN.

In fact, the differences between the result of DQN and the results of DDQN is not so large. I think the biggest difference between the two is that DDQN takes longer to converge, but is more stable than DQN after convergence. In some ways, such as score, DQN is even higher than DDQN in average. I think the reason for this phenomenon is that the hyper parameters chosen during training are not ideal, resulting in poor training of DDQN.

5 Experimental Summary

- Systematically learn and understand DQN and DDQN
- Learn to use pytorch to build neural networks
- Implement and successfully run the algorithm of DQN and DDQN based on pseudo-code
- Use matplotlib to plot curves