

# DLP Lab6: Deep Q-Network and Deep Deterministic Policy Gradient

Zhi-Yi Chin  
[joycenerd.cs09@nycu.edu.tw](mailto:joycenerd.cs09@nycu.edu.tw)

August 28, 2021

## 1 LunarLander-v2 tensorboard plot

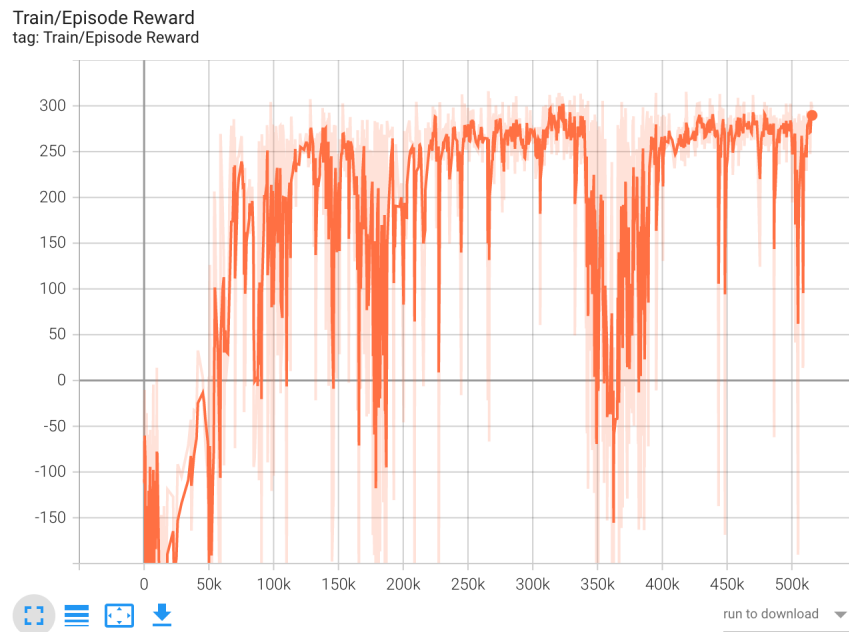


Figure 1: DQN episode reward with 2000 episodes.

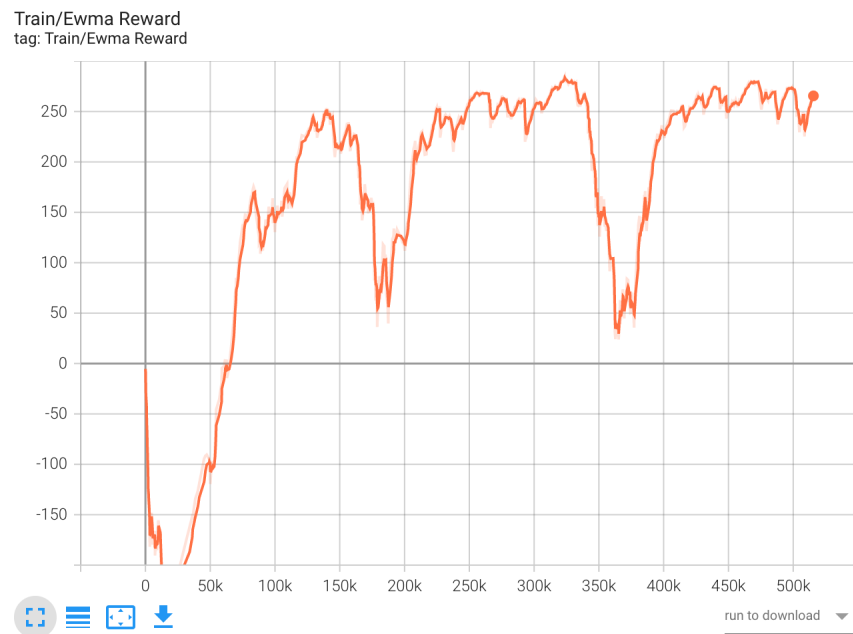


Figure 2: DQN ewma reward with 2000 episodes.

## 2 LunarLanderContinuous-v2 tensorboard plot

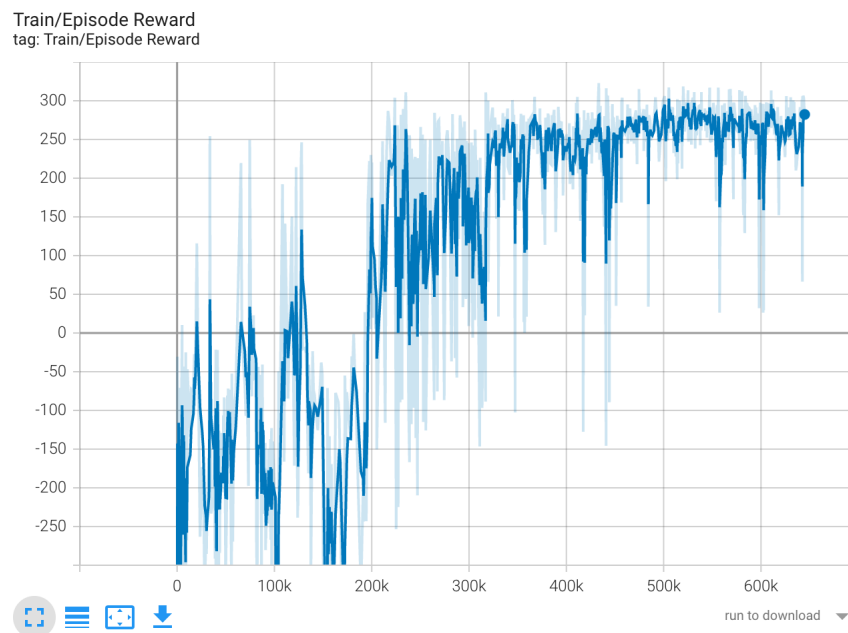


Figure 3: DDPG episode reward with 2000 episodes.

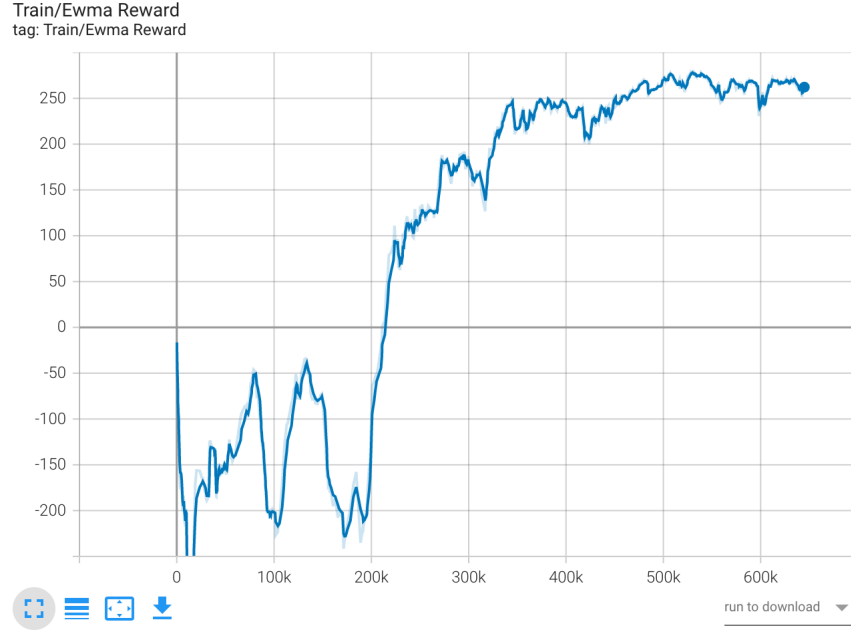


Figure 4: DDPG ewma reward with 2000 episodes.

### 3 Describe major implementation

#### 3.1 DQN

First, we build a network that can predict  $Q(s, a)$ . This network is a deep neural network (DNN); it takes an input state and outputs the Q-values of all possible actions for that state. We understand that the input layer of the DNN has the same size as a state size and that the output layer has the size of the number of actions that the agent can take. In our case, the number of actions is 4 (no-op, fire left agent, fire main engine, fire right engine). Also, we set the first hidden layer size to 384 and the second hidden layer size to 256. To summarize, when the agent is at a particular state, it forwards that state through the DNN and chooses the action with the highest Q-value.

```

1 class Net(nn.Module):
2     def __init__(self, state_dim=8, action_dim=4, hidden_dim=(384,256)):
3         super().__init__()
4         ## TODO ##
5         self.fc1=nn.Linear(state_dim,hidden_dim[0])
6         self.fc2=nn.Linear(hidden_dim[0],hidden_dim[1])
7         self.fc3=nn.Linear(hidden_dim[1],action_dim)
8         self.relu=nn.ReLU()
9
10    def forward(self, x):
11        ## TODO ##
12        x=self.fc1(x)
13        x=self.relu(x)
14        x=self.fc2(x)
15        x=self.relu(x)
16        out=self.fc3(x)
17        return out

```

Listing 1: DNN network to predict  $Q(s, a)$ .

We will do the update on the behavior network using gradient descent. We choose Adam as the optimizer with a learning rate of 0.0005.

```

1 self._optimizer = Adam(self._behavior_net.parameters(),lr=args.lr)

```

Listing 2: Optimizer for the behavior network.

In order for the agent to choose to explore or exploit, we use  $\epsilon$ -greedy. We randomly generate a number if the number is larger than  $\epsilon$  (initial to 1 and gradually decrease to 0.01), then we choose the action  $a_i$  with the largest  $Q(s, a_i)$  from the behavior network. Otherwise, we let the agent explore.

```

1 def select_action(self, state, epsilon, action_space):
2     '''epsilon-greedy based on behavior network'''
3     ## TODO ##
4     sample=random.random()
5     if sample>epsilon:
6         with torch.no_grad():
7             # t.max(1) will return target column value of each row
8             # second column on max result is index where max element was
9             # found, so we pick action with the larger expected reward
10            state=torch.from_numpy(state).view(1,-1).to(self.device)
11            return self._behavior_net(state).max(1)[1].item()
12     else:
13         return action_space.sample()

```

Listing 3: Select action based on  $\epsilon$ -greedy.

To update the behavior network, first, we sample a minibatch of transitions (state, action, reward, next state, done) from the replay memory. Then we get the current  $Q(s, a)$  output from the behavior network. And use target network to get  $Q(s_{t+1}, a')$  which finds action  $a'$  that yields the maximum  $Q$  in the next state. Then we get the target  $Q$ ,  $y_i$  by TD learning.

$$y_i = \begin{cases} R_T & \text{for terminal state } s_T \\ R_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') & \text{for non-terminal state } s_t \end{cases} \quad (1)$$

Next, we minimize the mean square error between our target  $Q$  value (according to the Bellman optimality equation) and our current  $Q$  output:

$$L(\theta_i) = \mathbb{E}_{s, a \sim p(\cdot)} [(y_i - Q(s, a; \theta_i))^2] \quad (2)$$

Optimally, we want the error to decrease, meaning that our current policy's outputs are becoming more similar to the true  $Q$  values. Therefore, with the loss function defined as above, we perform a gradient step on the loss function according to:

$$\nabla_{\theta_i} L(\theta_i) = \mathbb{E}_{s, a \sim p(\cdot), s' \sim \epsilon} [(R_{t+1} + \gamma \max_{a'} Q(s_{t+1}, a') - Q(s, a; \theta_i)) \nabla_{\theta_i} Q(s, a; \theta_i)] \quad (3)$$

```

1 def _update_behavior_network(self, gamma):
2     # sample a minibatch of transitions
3     state, action, reward, next_state, done = \
4         self._memory.sample(self.batch_size, self.device)
5
6     ## TODO ##
7     q_value = self._behavior_net(state).gather(dim=1, index=action.long())
8
9     with torch.no_grad():
10         q_next = self._target_net(next_state).max(dim=1)[0].view(-1,1)
11         q_target = reward + gamma*q_next*(1-done) # 1-done: end of trajectory
12
13     criterion = nn.MSELoss()
14     loss = criterion(q_value, q_target)
15
16     # optimize
17     self._optimizer.zero_grad()
18     loss.backward()
19     nn.utils.clip_grad_norm_(self._behavior_net.parameters(), 5)
20     self._optimizer.step()

```

Listing 4: Updating the behavior network using TD learning and mean square error.

Every 1000 steps, we update the target network by copying from the behavior network. We do this because the behavior network update so frequently that sometimes it may be unstable, so we maintain a stable target network aside from the behavior network.

```

1 def _update_target_network(self):
2     '''update target network by copying from behavior network'''
3     ## TODO ##
4     self._target_net.load_state_dict(self._behavior_net.state_dict())

```

Listing 5: Update target network by copying the behavior network.

## 3.2 DDPG

Both DQN and DDPG need a replay memory for exploitation. We sample from replay memory randomly each time we update our behavior network. Each time we sample a minibatch of transitions from the replay memory.

```

1 class ReplayMemory:
2     __slots__ = ['buffer']
3
4     def __init__(self, capacity):
5         self.buffer = deque(maxlen=capacity)
6
7     def __len__(self):
8         return len(self.buffer)
9
10    def append(self, *transition):
11        # (state, action, reward, next_state, done)
12        self.buffer.append(tuple(map(tuple, transition)))
13
14    def sample(self, batch_size, device):
15        '''sample a batch of transition tensors'''
16        ## TODO ##
17        transitions = random.sample(self.buffer, batch_size)
18        return (torch.tensor(x, dtype=torch.float, device=device)
19                for x in zip(*transitions))

```

Listing 6: Replay memory.

DDPG used four neural networks: a Q network, a deterministic policy network, a target Q network, and a target policy network.

- $\theta^Q$ : Q network.
- $\theta^\mu$ : deterministic policy function.
- $\theta^{Q'}$ : target Q network.
- $\theta^{\mu'}$ : target policy network.

The Q network and policy network are like simple Advantage Actor-Critic. However, in DDPG, the actor directly maps states to actions (the network’s output directly to the output) instead of outputting the probability distribution across a discrete action space. The target networks are the time-delayed copies of their original networks that slowly track the learned networks. Using these target value networks greatly improves stability in learning. This is because, in methods that do not use target networks, the update equations of the network are interdependent on the values calculated by the network itself, which makes it prone to divergence. For example, in Eq.4,  $Q(s', a')$  depends Q function itself (at the moment it is being optimized).

$$Q(s, a) \leftarrow Q(s, a) + \alpha[R(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a)] \quad (4)$$

So, we have the standard Actor and Critic architecture for the deterministic policy network and the Q network.

```

1 class ActorNet(nn.Module):
2     def __init__(self, state_dim=8, action_dim=2, hidden_dim=(400, 300)):
3         super().__init__()
4         ## TODO ##
5         self.fc1=nn.Linear(state_dim,hidden_dim[0])
6         self.fc2=nn.Linear(hidden_dim[0],hidden_dim[1])
7         self.fc3=nn.Linear(hidden_dim[1],action_dim)
8         self.relu=nn.ReLU()
9         self.tanh=nn.Tanh()
10
11     def forward(self, x):
12         ## TODO ##
13         x=self.fc1(x)
14         x=self.relu(x)
15         x=self.fc2(x)
16         x=self.relu(x)
17         out=self.fc3(x)
18         out=self.tanh(out)
19         return out
20
21 class CriticNet(nn.Module):
22     def __init__(self, state_dim=8, action_dim=2, hidden_dim=(400, 300)):
23         super().__init__()
24         h1, h2 = hidden_dim
25         self.critic_head = nn.Sequential(
26             nn.Linear(state_dim + action_dim, h1),
27             nn.ReLU(),
28         )
29         self.critic = nn.Sequential(
30             nn.Linear(h1, h2),
31             nn.ReLU(),
32             nn.Linear(h2, 1),
33         )
34
35     def forward(self, x, action):
36         x = self.critic_head(torch.cat([x, action], dim=1))
37         return self.critic(x)

```

Listing 7: Actor and Critic network of DDPG.

Since we have two networks, so we will need two optimizers to update the network. We use Adam as the optimizer and set the learning rate to 0.001 for both networks.

```

1 self._actor_opt = Adam(self._actor_net.parameters(),lr=args.lra)
2 self._critic_opt = Adam(self._critic_net.parameters(),lr=args.lrc)

```

Listing 8: Optimizer for the two network.



In reinforcement learning, for discrete action spaces, exploration is done via probabilistically selecting a random action (such as  $\epsilon$ -greedy or Boltzmann exploration), we have done in DQN as well. For continuous action spaces, exploration is done via adding noise to the action itself (there is also the parameter space noise). We use Gaussian noise in our implementation.

```

1 def select_action(self, state, noise=True):
2     '''based on the behavior (actor) network and exploration noise'''
3     ## TODO ##
4     with torch.no_grad():
5         state=torch.from_numpy(state).view(1,-1).to(self.device)
6         action=self._actor_net(state)
7
8         if noise:
9             N=torch.from_numpy(self._action_noise.sample())
10            N=N.view(1,-1).to(self.device)
11            action+=N
12            action=action.squeeze().cpu().numpy()
13
14    return action

```

Listing 9: Select action based on the behavior (actor) network and exploration noise.

The behavior network update splits into two parts. (1) Actor (policy) network update. (2) Critic (value) network update. The value network is updated as is done in Q-learning. The Bellman equation obtains the updated Q value:

$$y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1} | \theta^{\mu'}) | \theta^{Q'}) \quad (5)$$

However, in DDPG, the next-state Q values are calculated with the target value network and target policy network. In the implementation, we used the target actor-network to choose the next action and the target critic network to choose the next-state Q value. Then, we minimize the mean-squared loss between the updated Q value and the original Q value:

$$L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i | \theta^Q))^2 \quad (6)$$

Note that the original Q value is calculated with the value network, not the target value network.

For the policy function, our objective is to maximize the expected return:

$$J(\theta) = \mathbb{E}[Q(s, a) |_{s=s_t, a_t=\mu(s_t)}] \quad (7)$$

To calculate the policy loss, we take the derivative of the objective function concerning the policy parameter. Since the actor (policy) function is differentiable, so we have to apply the chain

rule.

$$\nabla_{\theta^\mu} J(\theta) \approx \nabla_a Q(s, a) \nabla_{\theta^\mu} \mu(s|\theta^\mu) \quad (8)$$

Nevertheless, since we are updating the policy in an off-policy way with batches of experience, we take the mean of the sum of gradients calculated from the mini-batch.

$$\nabla_{\theta^\mu} J(\theta) \approx \frac{1}{N} \sum_i [\nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s=s_i}] \quad (9)$$

We add a negative sign to the policy loss in the implementation since our objective is to maximize the expected return.

```

1 def _update_behavior_network(self, gamma):
2     actor_net, critic_net = self._actor_net, self._critic_net
3     target_actor_net, target_critic_net = self._target_actor_net,
4                                           self._target_critic_net
5     actor_opt, critic_opt = self._actor_opt, self._critic_opt
6
7     # sample a minibatch of transitions
8     state, action, reward, next_state, done = self._memory.sample(
9         self.batch_size, self.device)
10
11     ## update critic ##
12     # critic loss
13     ## TODO ##
14     q_value = self._critic_net(state, action)
15     with torch.no_grad():
16         a_next = self._target_actor_net(next_state)
17         q_next = self._target_critic_net(next_state, a_next)
18         q_target = reward + gamma * q_next * (1 - done)
19
20     criterion = nn.MSELoss()
21     critic_loss = criterion(q_value, q_target)
22
23     # optimize critic
24     actor_net.zero_grad()
25     critic_net.zero_grad()
26     critic_loss.backward()
27     critic_opt.step()
28
29     ## update actor ##
30     # actor loss
31     ## TODO ##
32     action = self._actor_net(state)
33     actor_loss = -self._critic_net(state, action).mean()
34

```

```

35     # optimize actor
36     actor_net.zero_grad()
37     critic_net.zero_grad()
38     actor_loss.backward()
39     actor_opt.step()

```

Listing 10: Update the actor and the critic network.

We make a copy of the target network parameters for the target network updates and have them slowly track those of the learned networks via "soft updates". We control the speed of the update by  $\tau$ , which we set to 0.005.

$$\begin{aligned}
 \theta^{Q'} &\leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'} \\
 \theta^{\mu'} &\leftarrow \tau \theta^\mu + (1 - \tau) \theta^{\mu'} \\
 &\text{where } \tau \ll 1
 \end{aligned} \tag{10}$$

```

1  @staticmethod
2  def _update_target_network(target_net, net, tau):
3      '''update target network by _soft_ copying from behavior network'''
4      for target, behavior in zip(target_net.parameters(), net.parameters()):
5          ## TODO ##
6          target.data.copy_((1.0-tau)*target.data+tau*behavior.data)

```

Listing 11: Update the target network by soft copying from the behavior network.

For testing, we run ten episodes, which is equal to 10 trajectories. In each episode, we are given an initial state from the environment. We repeat selecting an action from our well-trained actor-network without exploitation (noise) term and execute the action until the end of the trajectory. We print out the total rewards we get in each episode and the average return of 10 episodes.

```

1 def test(args, env, agent, writer):
2     print('Start Testing')
3     seeds = (args.seed + i for i in range(10))
4     rewards = []
5     for n_episode, seed in enumerate(seeds):
6         total_reward = 0
7         env.seed(seed)
8         state = env.reset()
9         ## TODO ##
10        for t in itertools.count(start=1):
11            # env.render()
12
13            # select action
14            action = agent.select_action(state, noise=False)
15
16            # execute action
17            next_state, reward, done, _ = env.step(action)
18
19            state = next_state
20            total_reward += reward
21
22            if done:
23                writer.add_scalar('Test/Episode Reward', total_reward, n_episode)
24                print(f"Total reward: {total_reward:.4f}")
25                rewards.append(total_reward)
26                break
27
28    print('Average Reward', np.mean(rewards))
29    env.close()

```

## 4 Difference between your implementation and algorithms

Both DQN and DDPG have a warmup time in our implementation. This happens at the beginning of training. In this short time, we will not update the parameters of the network; we let the agent explore and save the transitions in the replay memory.

In DQN, we do not update the behavior network in every iteration; instead, we update every four iterations.

The original paper of DDPG uses the Ornstein-Uhlenbeck Process to generate the exploration noise, which is correlated with the previous noise, to prevent the noise from canceling out or freezing the overall dynamics. However, in our implementation, we use a simple Gaussian noise instead, which yields good results as well.

## 5 Implementation and the gradient of actor updating

In policy (actor) learning, We want to learn a deterministic policy  $\mu_\theta(s)$  which gives the action that maximizes  $Q(s, a)$ . Because the action space is continuous, and we assume the Q-function is differentiable with respect to the action, we can perform gradient ascent (with respect to policy parameters only) to solve:

$$\max_{\theta} \mathbb{E}_{s \sim D}[Q(s, \mu_\theta(s))] \quad (11)$$

Note that the Q-function parameters are treated as constants here, which means that the critic is fixed when we update the actor. The calculation of policy (actor) loss and the update of the actor are mentioned in Section 3.2. We copy the context again below.

To calculate the policy loss, we take the derivative of the objective function concerning the policy parameter. Since the actor (policy) function is differentiable, so we have to apply the chain rule.

$$\nabla_{\theta^\mu} J(\theta) \approx \nabla_a Q(s, a) \nabla_{\theta^\mu} \mu(s|\theta^\mu)$$

However, since we are updating the policy in an off-policy way with batches of experience, we take the mean of the sum of gradients calculated from the mini-batch.

$$\nabla_{\theta^\mu} J(\theta) \approx \frac{1}{N} \sum_i [\nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s=s_i}]$$

We add a negative sign to the loss in the implementation since our objective is to maximize the expected return.

```
1  ## update actor ##
2  # actor loss
3  ## TODO ##
4  action=self._actor_net(state)
5  actor_loss = -self._critic_net(state,action).mean()
6
7  # optimize actor
8  actor_net.zero_grad()
9  critic_net.zero_grad()
10 actor_loss.backward()
11 actor_opt.step()
```

Listing 12: Update actor by maximizing the expected return.

## 6 Implementation and the gradient of critic updating

The Bellman equation describing the optimal action-value function,  $Q^*(s, a)$  is given by:

$$Q^*(s, a) = \mathbb{E}_{s' \sim P}[r(s, a) + \gamma \max_{a'} Q^*(s', a')] \quad (12)$$

where  $s' \sim P$  is shorthand for saying that the next state  $s'$  is sampled by the environment from a distribution  $P(\cdot|s, a)$ . This Bellman equation is the starting point for learning an approximator to  $Q^*(s, a)$ . Suppose the approximator is a neural network  $Q(s, a)$  with parameters  $\theta^Q$ , and that we have collected a set  $\mathcal{D}$  of transitions  $(s, a, r, s', d)$ . We can set a mean-squared Bellman error (MSBE) function, which tells us roughly how closely  $Q$  comes to satisfying the Bellman equation:

$$L = \mathbb{E}_{(s, a, r, s', d) \sim \mathcal{D}}[(Q(s, a) - (r + \gamma(1 - d) \max_{a'} Q(s', a')))^2] \quad (13)$$

Q-learning algorithms for function approximators, such as DQN (and all its variants) and DDPG, primarily minimize this MSBE loss function. All of them employ two main tricks. (1) Replay buffer, in which we will draw a batch of transitions to update the network. (2) Target network, which is the updated (next-state) Q value.

$$y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$$

In the implementation, we use the target actor-network to get the next action and the target critic network to get the next Q value and applied the equation above to get the target Q value.

We can rewrite Eq.13 as follow:

$$L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$$

where the original Q value is calculated with the value network, not the target value network.

Also the gradient of the critic loss is

$$\nabla_{\theta^Q} L = \frac{2}{N} (Q(s, a|\theta^Q) - y) \quad (14)$$

```

1  ## update critic ##
2  # critic loss
3  ## TODO ##
4  q_value = self._critic_net(state,action)
5  with torch.no_grad():
6      a_next = self._target_actor_net(next_state)
7      q_next = self._target_critic_net(next_state,a_next)
8      q_target = reward+gamma*q_next*(1-done)
9
10 criterion = nn.MSELoss()
11 critic_loss = criterion(q_value, q_target)
12
13 # optimize critic
14 actor_net.zero_grad()
15 critic_net.zero_grad()
16 critic_loss.backward()
17 critic_opt.step()

```

Listing 13: Update the critic with mean-squared Bellman loss

## 7 Effects of the discount factor

The discount factor essentially determines how much the reinforcement learning agents care about rewards in the distant future relative to those in the immediate future. If  $\gamma = 0$ , the agent will be wholly myopic and only learn about actions that produce an immediate reward. If  $\gamma = 1$ , the agent will evaluate each of its actions based on the total of all its future rewards. Picking a particular value of  $\gamma$  rewrite an agent's discounted reward  $G$  as

$$G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \dots \quad (15)$$

We concern more about the future near the current state and concern less about the far future. Also, if a trajectory has no end, then without discount factor might get an infinite value on the expected return.

## 8 Benefits of $\epsilon$ -greedy in comparison to greedy action

When the agent explores, it can improve its current knowledge and gain better rewards in the long run. However, it gets more reward immediately when it exploits, even if it is a sub-optimal behavior. As the agent cannot do both at the same time, there is a trade-off.

A greedy agent can get stuck in a sub-optimal state since it never lets itself explore new possibilities, or there might be changes in the environment as time passes. In  $\epsilon$ -greedy action selection, the agent uses both exploitations to take advantage of prior knowledge and exploration to look for new options.

## 9 The necessity of the target network

Take DQN as example, if we didn't use the target network, the update target was  $(r_t + \max_a Q(s_{t+1}, a; \theta^Q) - Q(s_t, a_t; \theta^Q))^2$ , then learning would become unstable because the target  $r_t + \max_a Q(s_{t+1}, a; \theta^Q)$  and the prediction  $Q(s_t, a_t; \theta^Q)$  are not independent, as they both rely on  $\theta^Q$ .

We use the target network that does not update as frequently as the behavior network to stabilize the training process.

## 10 Explain the effect of replay buffer size in case of too large or too small

If the replay buffer is too small, the replay buffer serves little to no purpose. A small buffer might force the network to only care about what it saw recently, and it might overfit and break.

Usually, we want our experience replay buffer to be as large as possible. The larger the experience replay, the less likely we will sample correlated elements; hence the more stable the network's training will be. However, a large experience replay also requires much memory, and it might slow down the training. So, there is a trade-off between training stability and memory requirements.

## 11 Double-DQN

Double DQN handles the problem of the overestimation of Q-values. From Eq.1, we can see that, by calculating the TD target, we face a simple problem: how are we sure that the best action for the next state is the action with the highest Q-values? We know that the accuracy of Q-values depends on what action we tried and what neighboring states we explored. Consequently, we do not have enough information about the best action to take at the beginning of the training. Therefore, taking the maximum Q-value (which is noisy) as the best action to take can lead to false positives. If non-optimal actions are regularly given a higher Q-value than the optimal best action, the learning will be complicated. Double DQN provides the solution. When computing the Q target, we use two networks to decouple the action selected from the target Q-value generation. We use our DQN network to select the best action to take for the next state (the action with the highest Q-value). Then we use our target network to calculate the target Q-value of taking that action at the next state.

$$Q(s, a) \leftarrow r_t + \gamma Q'(s_{t+1}, \arg \max_a Q(s_t, a_t)) \quad (16)$$

where  $Q(s, a)$  is the behavior network and  $Q'(s, a)$  is the target network. Therefore, double DQN helps us reduce the overestimation of Q values and, as a consequence, helps us train faster and have more stable training.

The only difference between DQN and double DQN is when updating the behavior network.



```

1 def _update_behavior_network(self, gamma):
2     # sample a minibatch of transitions
3     state, action, reward, next_state, done = self._memory.sample(self.batch_size,
4                                                                    self.device)
5
6     ## TODO ##
7     q_value = self._behavior_net(state).gather(dim=1, index=action.long())
8
9     with torch.no_grad():
10         action=self._behavior_net(next_state).max(dim=1)[1].view(-1,1)
11         q_next = self._target_net(next_state).gather(dim=1, index=action.long())
12         q_target = reward + gamma*q_next*(1-done)
13
14     criterion = nn.MSELoss()
15     loss = criterion(q_value, q_target)
16
17     # optimize
18     self._optimizer.zero_grad()
19     loss.backward()
20     nn.utils.clip_grad_norm_(self._behavior_net.parameters(), 5)
21     self._optimizer.step()

```

Listing 14: Update the behavior network of DDQN, which takes the action from behavior network and generate the next-state Q value from the target network.

Figure 5 and 6 show the episode and the ewma reward of DDQN respectively.

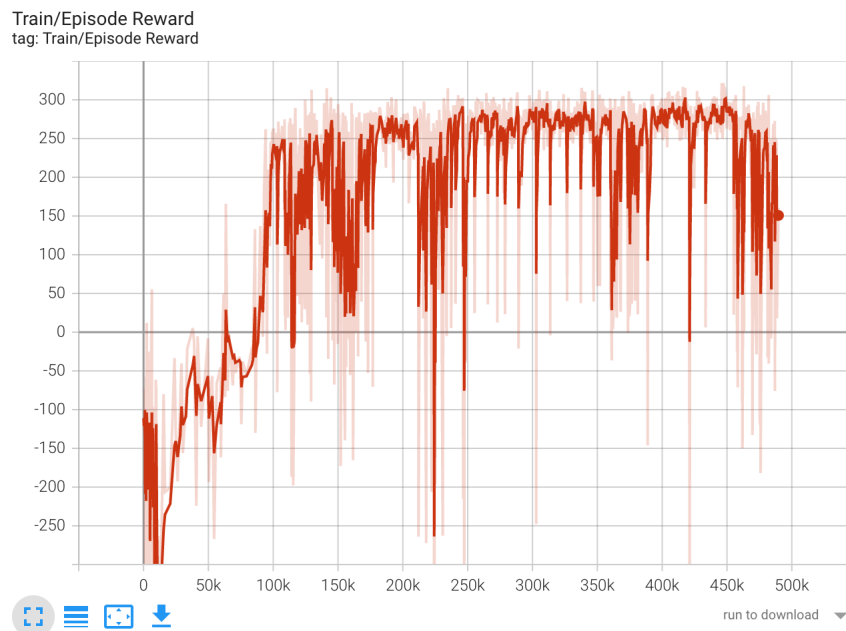


Figure 5: DDQN episode reward with 2000 episodes

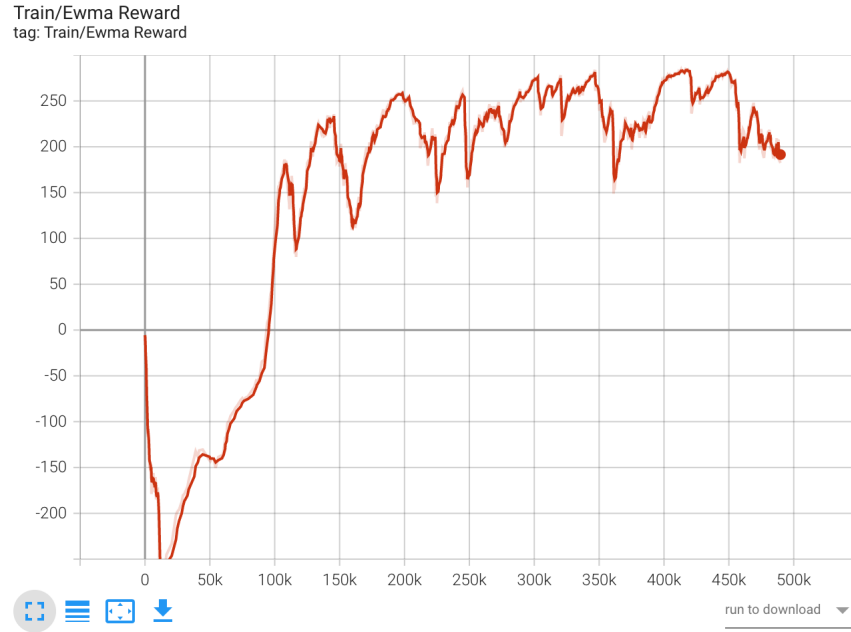


Figure 6: DDQN episode reward with 2000 episodes

## 12 Results

For all the results, we train for 2000 episodes.

```
Start Testing
Total reward: 247.2968
Total reward: 282.2878
Total reward: 250.9439
Total reward: 262.0189
Total reward: 316.1820
Total reward: 265.7174
Total reward: 278.1898
Total reward: 284.7278
Total reward: 294.1539
Total reward: 297.6763
Average Reward 277.91946829695576
(dqn) zchin@eva-System-Product-Name:~/Deep_Learning_Practice_labs/lab6$
```

Figure 7: DQN testing result. The average reward is 277.

```
Total reward: 256.9653
Total reward: 291.9388
Total reward: 281.0775
Total reward: 285.6077
Total reward: 293.4235
Total reward: 270.9022
Total reward: 287.2607
Total reward: 300.5433
Total reward: 310.0303
Total reward: 152.8534
Average Reward 273.06028514742735
(dqn) zchin@eva-System-Product-Name:~/Deep_Learning_Practice_labs/lab6$
```

Figure 8: DDPG testing result. The average reward is 273.

```
Start Testing
Total reward: 248.7711
Total reward: 270.5899
Total reward: 280.4892
Total reward: 275.7324
Total reward: 28.7430
Total reward: 262.5992
Total reward: 293.9000
Total reward: 296.4455
Total reward: 306.9310
Total reward: 301.6734
Average Reward 256.58746746801523
(dqn) zchin@eva-System-Product-Name:~/Deep_Learning_Practice_labs/lab6$
```

Figure 9: DDQN testing result. The average reward is 256.