

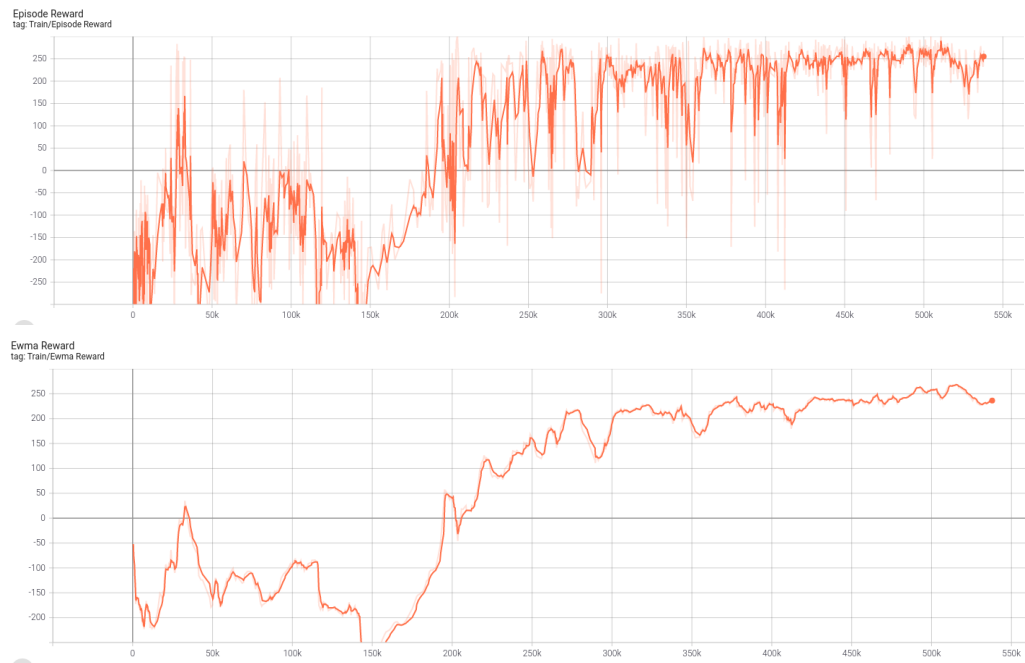
# LAB6 Deep Q-Network and Deep Deterministic Policy Gradient

## 1. Tensorboard

### LunarLander-v2 (DQN)



### LunarLanderContinuous-v2 (DDPG)



## 2. Implement detail

### 1. Describe your major implementation of both algorithms in detail. (TODO)

DQN :

Create a network to generate distribution of four action .

```
41 class Net(nn.Module):
42     def __init__(self, state_dim=8, action_dim=4, hidden_dim=32):
43         super().__init__()
44         ## TODO ##
45         self.fc1 = nn.Linear(8, 64)
46         self.fc2 = nn.Linear(64, 64)
47         self.fc3 = nn.Linear(64, 32)
48         self.out = nn.Linear(32, 4)
49
50     def forward(self, x):
51         ## TODO ##
52         x = torch.relu(self.fc1(x))
53         x = torch.relu(self.fc2(x))
54         x = torch.relu(self.fc3(x))
55         action_distribution = self.out(x)
56
57         return action_distribution
```

Set the optimizer and loss function .

```
self._optimizer = torch.optim.Adam(self._behavior_net.parameters(), lr = 0.0005)
self.criterion = nn.MSELoss()
```

When training , use the action of the highest possibility  $Q(S, a_i)$  or randomly choose action with the epsilon .

With probability  $\varepsilon$  select a random action  $a_t$   
otherwise select  $a_t = \operatorname{argmax}_a Q(\phi(s_t), a; \theta)$

```
78 def select_action(self, state, epsilon, action_space):
79     '''epsilon-greedy based on behavior network'''
80     x = torch.unsqueeze(torch.FloatTensor(state), 0).cuda()
81     if random.random() < epsilon:
82         action_ = action_space.sample()
83     else:
84         with torch.no_grad():
85             action_value = self._behavior_net(x)
86             action_ = torch.argmax(action_value).cpu().numpy()
87
88     return action_
```

Update the behavior network with the algorithm :

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  $\theta$

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

End For

End For

Deep Learning and Practice, Spring 2021,

```

102 def _update_behavior_network(self, gamma):
103     # sample a minibatch of transitions
104     state, action, reward, next_state, done = self._memory.sample(
105         self.batch_size, self.device)
106
107     self._optimizer.zero_grad()
108     q_value = self._behavior_net(state).gather(1, action.long())
109     with torch.no_grad():
110         q_next = self._target_net(next_state)
111         q_target = reward + (1-done)*gamma*torch.max(q_next, dim=1)[0].view(self.batch_size, 1)
112
113     loss = self.criterion(q_value, q_target)
114     # optimize
115     self._optimizer.zero_grad()
116     loss.backward()
117     nn.utils.clip_grad_norm_(self._behavior_net.parameters(), 5)
118     self._optimizer.step()

```

Update the target network with behavior network every 4 steps .

```

120 def _update_target_network(self):
121     '''update target network by copying from behavior network'''
122     self._target_net.load_state_dict(self._behavior_net.state_dict())

```

My DQN training parameters are using default .

DDPG :

Create a actor network to generate value of “Main engine” and “Left-Right engine” , so the num of the output is two .

```

48 class ActorNet(nn.Module):
49     def __init__(self, state_dim=8, action_dim=2, hidden_dim=(400, 300)):
50         super().__init__()
51         ## TODO ##
52         h1, h2 = hidden_dim
53         self.fc1 = nn.Linear(state_dim, h1)
54         self.fc2 = nn.Linear(h1, h2)
55         self.out = nn.Linear(h2, 2)
56
57     def forward(self, x):
58         x = torch.relu(self.fc1(x))
59         x = torch.relu(self.fc2(x))
60         action_value = torch.tanh(self.out(x))
61
62         return action_value

```

Set the optimizer and loss function .

```

97     self._actor_opt = torch.optim.Adam(self._actor_net.parameters(), lr=args.lra)
98     self._critic_opt = torch.optim.Adam(self._critic_net.parameters(), lr=args.lrc)
99     self.criterion = nn.MSELoss()

```

When training , generate action values and add the noise .

Select action  $a_t = \mu(s_t|\theta^\mu) + N_t$  according to the current policy and exploration noise

```

111     def select_action(self, state, noise=True):
112         '''based on the behavior (actor) network and exploration noise'''
113         ## TODO ##
114         with torch.no_grad():
115             x = torch.unsqueeze(torch.FloatTensor(state),0).cuda()
116             action = self._actor_net(x)
117             action = action.detach().cpu().numpy()[0]
118             if noise:
119                 action = action + self._action_noise.sample()
120
121         return action

```

Update the behavior network with the algorithm :

Sample random minibatch of  $N$  transitions  $(s_i, a_i, r_i, s_{i+1})$  from  $R$

Set  $y_i = r_i + \gamma Q'(s_{t+1}, \mu'(s_{t+1}|\theta^{\mu'})|\theta^{Q'})$

Update critic by minimizing the loss:  $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$

Update the actor policy using the sampled gradient:

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

```

139     def update_behavior_network(self, gamma):
140         actor_net, critic_net, target_actor_net, target_critic_net = self._actor_net, self._critic_net, self._target_actor_net, self._target_critic_net
141         actor_opt, critic_opt = self._actor_opt, self._critic_opt
142
143         # sample a minibatch of transitions
144         state, action, reward, next_state, done = self.memory.sample(
145             self.batch_size, self.device)
146
147         ## update critic ##
148         # critic loss
149         ## TODO ##
150         q_value = critic_net(state, action)
151         with torch.no_grad():
152             a_next = target_actor_net(next_state)
153             q_next = target_critic_net(next_state, a_next)
154             q_target = reward + gamma*q_next*(1 - done)
155
156         critic_loss = self.criterion(q_value, q_target)
157         # optimize critic
158         actor_net.zero_grad()
159         critic_net.zero_grad()
160         critic_loss.backward()
161         critic_opt.step()
162
163         ## update actor ##
164         # actor loss
165         ## TODO ##
166         action = actor_net(state)
167         actor_loss = (-1.0)*critic_net(state, action).mean()
168         # optimize actor
169         actor_net.zero_grad()
170         critic_net.zero_grad()
171         actor_loss.backward()
172         actor_opt.step()

```

Update the target networks .

Update the target networks:

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

```
174 @staticmethod
175 def _update_target_network(target_net, net, tau):
176     '''update target network by _soft_ copying from behavior network'''
177     for target, behavior in zip(target_net.parameters(), net.parameters()):
178         ## TODO ##
179         target.data.copy_(tau*behavior.data + (1.0 -tau)*target.data)
```

## 2. Describe differences between your implementation and algorithms.

When start training , in order to have enough sample for training , we don't update the network within warmup step , just randomly choose action to play the game and store the result in the replay memory . And in the DQN , we update the target network every few iteration , it can reduce the correlation between target and behavior network .

## 3. Describe your implementation and the gradient of actor updating.

```
166 action = actor_net(state)
167 actor_loss = (-1.0)*critic_net(state, action).mean()
168 # optimize actor
169 actor_net.zero_grad()
170 critic_net.zero_grad()
171 actor_loss.backward()
172 actor_opt.step()
```

for actor :

expect large  $Q(s, a)$  ,

so make  $J(\theta^\mu) = E[Q(s, a)]$

$\Rightarrow \nabla_{\theta^\mu} J(\theta^\mu) = E\left[\frac{\partial Q(s, a|\theta^\mu)}{\partial \theta^\mu}\right]$ , with  $a = \mu(s|\theta^\mu)$

$= E\left[\frac{\partial Q(s, a|\theta^\mu)}{\partial a} \cdot \frac{\mu(s|\theta^\mu)}{\partial \theta^\mu}\right]$

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

So we can define loss function as :

Actor loss  $= - \nabla_{\theta^\mu} J(\theta^\mu)$

We only update the actor network .

**4. Describe your implementation and the gradient of critic updating.**

Compute the MSELoss with Q\_target from target network and Q\_value from behavior network .

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

```
150 q_value = critic_net(state, action)
151 with torch.no_grad():
152     a_next = target_actor_net(next_state)
153     q_next = target_critic_net(next_state, a_next)
154     q_target = reward + gamma*q_next*(1 - done)
155
156 critic_loss = self.criterion(q_value, q_target)
157 # optimize critic
158 actor_net.zero_grad()
159 critic_net.zero_grad()
160 critic_loss.backward()
161 critic_opt.step()
```

**5. Explain effects of the discount factor.**

As the training step is longer , the sum of the reward will become infinity , so the discount factor can reduce the reward in the later step , in other words , it reduce the correlation of the later step . In this lab , only using the one next step , but we still give the 0.99 as discount factor .

**6. Explain benefits of epsilon-greedy in comparison to greedy action selection.**

The epsilon-greedy make additional opportunity to select the best action , because the action selected by the network may not be the best selection , and with the episode growing , we increase the epsilon . In conclusion , it can improve the exploratory of the network in the early step .

**7. Explain the necessity of the target network.**

We use the periodically update network (target network) to improve the stability of training .

**8. Explain the effect of replay buffer size in case of too large or too small.**

The higher replay buffer size can make training more stable but need more time to train . And if too small , it is easier to occur overfitting .

### 3. Bonus

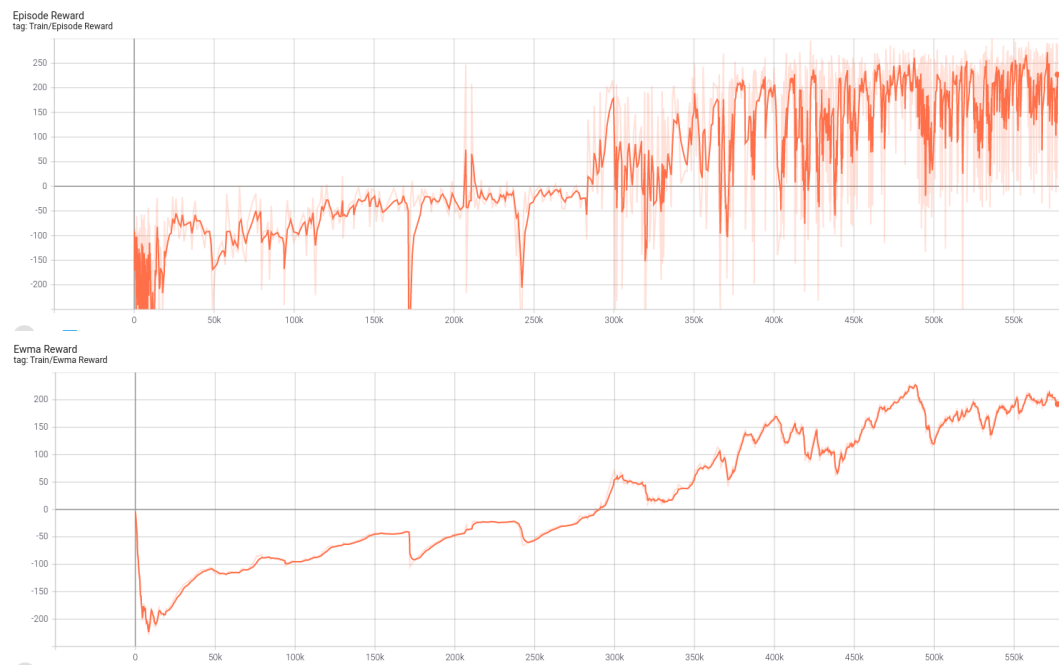
#### DDQN

The difference between DQN and DDQN is that in DDQN case , when compute the  $q\_target$  , it don't use the max of  $Q'(s_i, a_i)$  , but use the index of the max value in  $Q(s_i, a_i)$  as the index of  $Q'(s_i, a_i)$  , that mean in DDQN , the action selection and evaluation are generate from the different function .

```
q_value = self._behavior_net(state).gather(1, action.long())
with torch.no_grad():
    next_ = self._behavior_net(next_state).max(dim=1)[1].view(-1,1)
    q_next = self._target_net(next_state).gather(1, next_.long())
    q_target = reward + gamma*q_next

loss = self.criterion(q_value, q_target)
```

Tensorboard :



### 4. Performance ( Average reward of 10 games )

DQN ( train 1200 episode ) :

Average Reward 274.02778984625076

DDPG ( train 1500 episode ) :

Average Reward 230.17661800337765

DDQN ( train 1200 episode ) :

Average Reward 240.21417588222707