

ECE433/COS435 Introduction to RL
Assignment 6: Actor-critic Algorithms for continuous
action space: DDPG & TD3
Spring 2025

Fill me in

Your name here.

Due April 4, 2025

Collaborators

Fill me in

Please fill in the names and NetIDs of your collaborators in this section.

Instructions

Writeups should be typesetted in Latex and submitted as PDFs. You can work with whatever tool you like for the code, but **please submit the asked-for snippet and answer in the solutions box as part of your writeup. We will only be grading your write-up.**

In this assignment, we will introduce two modern RL algorithms that aim to tackle MDPs with continuous action space: Deep Deterministic Policy Gradient (DDPG) and Twin Delayed DDPG (TD3).

The workflow for the rest of this assignment is as follows: we will initially implement DDPG algorithms, and then evaluate their performance on the Gym environment MountainCarContinuous. The implementation for TD3 is directly related to DDPG. By studying these two algorithms, you will gain insight into reinforcement learning in continuous action spaces and become familiar with various techniques in RL.

Question 1. Deep Deterministic Policy Gradient(DDPG)

Question 1.a

First, you want to build your Actor-DDPG and Critic-DDPG networks. Paste your class below:

Solution

```
1 class Actor-DDPG(nn.Module):
2     def __init__(self, state_dim: int, action_dim: int, max_action:
      torch.tensor):
3         super(Actor-DDPG, self).__init__()
4         #####
5         # YOUR IMPLEMENTATION HERE #
6         pass
7         #####
8
9         self.max_action = max_action
10
11     def forward(self, state):
12         #####
13         # YOUR IMPLEMENTATION HERE #
14         pass
15         #####
16
17 class Critic-DDPG(nn.Module):
18     def __init__(self, state_dim: int, action_dim: int):
19         super(Critic-DDPG, self).__init__()
20         #####
21         # YOUR IMPLEMENTATION HERE #
22         pass
23         #####
24
25     def forward(self, state: torch.Tensor, action: torch.Tensor)->torch
      .Tensor:
26         sa = torch.cat([state, action], 1)
27         #####
28         # YOUR IMPLEMENTATION HERE #
29         pass
30         #####
```

Question 1.b

Now we are ready to construct a DDPG trainer! In the following function you will need to: (1) Calculate the TD value using target_Q network and update the critic; (2) Calculate the deterministic policy gradient and update the actor; (3) Update the target networks.

Paste below:

Solution

```
1 def train(self, replay_buffer, batch_size=64):
2     # Sample from replay buffer
3     state, action, next_state, reward, not_done = replay_buffer.
        sample(batch_size)
4
5     # Compute the target_Q value with critic_target for the batch
6
7     #####
8     # YOUR IMPLEMENTATION HERE #
9     pass
10    #####
11
12    # Get current Q estimate
13    current_Q = self.critic(state, action)
14
15    # Compute critic loss
16    #####
17    # YOUR IMPLEMENTATION HERE #
18    critic_loss = None
19    #####
20
21    # Optimize the critic
22    self.critic_optimizer.zero_grad()
23    critic_loss.backward()
24    self.critic_optimizer.step()
25
26    # Compute actor loss
27    #####
28    # YOUR IMPLEMENTATION HERE #
29    actor_loss = None
30    #####
31
32    # Optimize the actor
33    self.actor_optimizer.zero_grad()
34    actor_loss.backward()
35    self.actor_optimizer.step()
36
37    # Update the target models
38    for param, target_param in zip(self.critic.parameters(), self.
        critic_target.parameters()):
39        #####
40        # YOUR IMPLEMENTATION HERE #
41        new_target_params = None
42        target_param.data.copy_(new_target_params)
43        #####
44
45    for param, target_param in zip(self.actor.parameters(), self.
        actor_target.parameters()):
46        #####
47        # YOUR IMPLEMENTATION HERE #
```

```

48     new_target_params = None
49     target_param.data.copy_(new_target_params)
50     #####

```

Question 1.c

Time to see how it works. We would expect a reward that converges to around 90. The estimated wall time for running the whole process is around 10 – 20 minutes, and you should be able to see a large positive reward at around 8×10^4 timesteps. If the initialization is unsuccessful, which could result in the reward being stuck at around 0, try restarting the main function or debugging your trainer.

Note: Even if your implementation has no errors, its performance may not be perfect every time; this is acceptable because it involves randomness. Paste your code for drawing training curve below, also provide the best training curve you can get.

Solution

1. Paste the code for plotting the training curve below:

```

1 import matplotlib.pyplot as plt
2
3 #####
4 # YOUR IMPLEMENTATION HERE #
5 pass
6 #####

```

2. Training curve below:

Question 2. Twin Delayed DDPG (TD3)

Question 2.a

First, we will implement the actor and critic network for TD3. For the actor and every critic (we need to maintain an additional critic), please make sure it has the same structure as the one in the previous DDPG question so that we can conduct an ablation study.

Paste your class below:

Solution

```

1 class Actor_TD3(nn.Module):
2     def __init__(self, state_dim: int, action_dim: int, max_action:
3         float):
4         super(Actor_TD3, self).__init__()
5         #####
6         # YOUR IMPLEMENTATION HERE #
7         #####
8         self.max_action = max_action

```

```

8
9
10 def forward(self, state):
11     #####
12     # YOUR IMPLEMENTATION HERE #
13     #####

```

Solution

```

1 class Critic_TD3(nn.Module):
2     def __init__(self, state_dim: int, action_dim: int):
3         super(Critic_TD3, self).__init__()
4         # Q1 architecture
5         #####
6         # YOUR IMPLEMENTATION HERE #
7         pass
8         #####
9
10        # Q2 architecture
11        #####
12        # YOUR IMPLEMENTATION HERE #
13        pass
14        #####
15
16    def forward(self, state: torch.Tensor, action: torch.Tensor) ->
17        Tuple[torch.Tensor, torch.Tensor]:
18        sa = torch.cat([state, action], 1)
19
20        #####
21        # YOUR IMPLEMENTATION HERE #
22        pass
23        #####
24        return q1, q2
25
26    # Implement a function that returns only Q1, which is helpful when
27    # calculating actor loss
28    def Q1(self, state: torch.Tensor, action: torch.Tensor) -> torch.
29        Tensor:
30        sa = torch.cat([state, action], 1)
31        #####
32        # YOUR IMPLEMENTATION HERE #
33        pass
34        #####
35        return q1

```

Question 2.b

Now let's implement the TD3 trainer!

Solution

```
1 def train(self, replay_buffer, batch_size=256):
2     self.total_it += 1
3
4     # Sample replay buffer
5     state, action, next_state, reward, not_done = replay_buffer.
6     sample(batch_size)
7
8     with torch.no_grad():
9         # 1. Select action according to policy and add clipped noise.
10        #####
11        # YOUR IMPLEMENTATION HERE #
12        pass
13        #####
14
15        # Compute the target Q value
16        target_Q1, target_Q2 = self.critic_target(next_state,
17        next_action)
18
19        # 2. Compute the target_Q here
20        #####
21        # YOUR IMPLEMENTATION HERE #
22        pass
23        #####
24
25        # Get current Q estimates
26        current_Q1, current_Q2 = self.critic(state, action)
27
28        # Compute critic loss
29        critic_loss = F.mse_loss(current_Q1, target_Q) + F.mse_loss(
30        current_Q2, target_Q)
31
32        # Optimize the critic
33        self.critic_optimizer.zero_grad()
34        critic_loss.backward()
35        self.critic_optimizer.step()
36
37        # Delayed policy updates
38        if self.total_it % self.policy_freq == 0:
39
40            # Compute actor loss
41            #####
42            # YOUR IMPLEMENTATION HERE #
43            pass
44            #####
45
46            # Optimize the actor
47            self.actor_optimizer.zero_grad()
48            actor_loss.backward()
49            self.actor_optimizer.step()
```

```

48
49     # Update the frozen target models using weighted mean
50     for param, target_param in zip(self.critic.parameters(), self.
critic_target.parameters()):
51         #####
52         # YOUR IMPLEMENTATION HERE #
53         new_target_params = None
54         target_param.data.copy_(new_target_params)
55         #####
56
57     for param, target_param in zip(self.actor.parameters(), self.
actor_target.parameters()):
58         #####
59         # YOUR IMPLEMENTATION HERE #
60         new_target_params = None
61         target_param.data.copy_(new_target_params)
62         #####

```

Question 2.c

Time to see how it works. Plot your reward v.s. training.episode curve.

Note: Even if your implementation has no errors, its performance may not be perfect every time; this is acceptable because it involves randomness. Paste your code for drawing training curve below, also provide the best training curve you can get.

Solution

1. Paste the code for plotting the training curve below:

```

1 import matplotlib.pyplot as plt
2
3 #####
4 # YOUR IMPLEMENTATION HERE #
5 pass
6 #####

```

2. Training curve below: