

ECE433/COS435 Introduction to RL
Assignment 7: Actor-critic Algorithms for continuous
action space: SAC
Spring 2025

Fill me in

Your name here.

Due April 11, 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 implement the SAC algorithm, and then evaluate their performance on the Gym environment Pendulum-v1. We provide a sample implementation to TD3, which can be helpful as a reference for your implementation of SAC.

Question 1. Soft Actor-Critic (SAC)

Question 1.a

First, you want to build your Actor_SAC and Critic_SAC networks. Paste your class below:

Solution

```
1 class Actor_SAC(nn.Module):
2     def __init__(self, state_dim, action_dim, max_action):
3         super(Actor_SAC, self).__init__()
4         # [HINT] Construct a neural network as the actor. Return its
          value using forward You need to write down three linear layers.
5         # 1. l1: state_dim -> 256
6         # 2. l2: 256 -> 256
7         # 3. l3: 256 -> mean and log std of the action
8         #####
9         # YOUR IMPLEMENTATION HERE #
10        pass
11        #####
12        self.max_action = max_action
13
14    def forward(self, state):
15        # [HINT] Use the three linear layers to compute the mean and
          log std of the action
16        # Apply ReLU activation after layer l1 and l2
17        #####
18        # YOUR IMPLEMENTATION HERE #
19        pass
20        #####
21        log_std = torch.clamp(log_std, min=LOG_STD_MIN, max=
          LOG_STD_MAX)
22        return mean, log_std
23
24    def sample(self, state):
25        # [HINT] Use the forward method to compute the action, its
          log probability
26        # 1. Compute the mean and log std of the action
27        # 2. Compute the standard deviation of the action
28        # 3. Get the normal distribution of the action
29        # 4. Sample the action from the normal distribution
30        # 5. Apply tanh to the action and multiply by max_action to
          ensure the action is in the range of the action space
31        # 6. Compute the log probability of the action
32
33        #####
34        # YOUR IMPLEMENTATION HERE #
35        pass
36        #####
37        return action, log_prob

```

```
1 class Critic_SAC(nn.Module):
2     def __init__(self, state_dim, action_dim):
3         super(Critic_SAC, self).__init__()
4         # Q1 architecture
5         # [HINT] Construct a neural network as the first critic.
          Return its value using forward You need to write down three linear
          layers.
6         # 1. l1: state_dim+action_dim -> 256

```

```

7         # 2. l2: 256 -> 256
8         # 3. l3: 256 -> 1
9         #####
10        # YOUR IMPLEMENTATION HERE #
11        pass
12        #####
13
14        # Q2 architecture
15        # [HINT] Construct a neural network as the second critic.
Return its value using forward. You need to write down three
linear layers.
16        # 1. l4: state_dim+action_dim -> 256
17        # 2. l5: 256 -> 256
18        # 3. l6: 256 -> 1
19        #####
20        # YOUR IMPLEMENTATION HERE #
21        pass
22        #####
23
24
25    def forward(self, state, action):
26        sa = torch.cat([state, action], 1)
27        # [HINT] We use layers l1, l2, l3 to obtain q1
28        # 1. Apply ReLU activation after layer l1
29        # 2. Apply ReLU activation after layer l2
30        # 3. Return output as q1 from layer l3
31
32        # [HINT] We use layers l4, l5, l6 to obtain q2
33        # 1. Apply ReLU activation after layer l4
34        # 2. Apply ReLU activation after layer l5
35        # 3. Return output as q2 from layer l6
36
37        #####
38        # YOUR IMPLEMENTATION HERE #
39        pass
40        #####
41        return q1, q2
42
43
44    def Q1(self, state, action):
45        sa = torch.cat([state, action], 1)
46        # [HINT] only returns q1 for actor update using layers l1, l2
, l3
47        # 1. Apply ReLU activation after layer l1
48        # 2. Apply ReLU activation after layer l2
49        # 3. Return output as q1 from layer l3
50        #####
51        # YOUR IMPLEMENTATION HERE #
52        pass
53        #####
54        return q1

```

Question 1.b

Now we are ready to construct a SAC trainer! In the following function you will need to: (1) Calculate the TD value using target_Q network and update the critic; (2) Compute the actor loss using the current state and sampled action and update the actor; (3) Update the target networks.

Solution

```
1 def train(self, replay_buffer, batch_size=256):
2     state, action, next_state, reward, not_done = replay_buffer.
3     sample(batch_size)
4
5     # [HINT] compute the target Q value
6     # 1. Sample the next action and its log probability from the
7     actor with next_state
8     # 2. Compute the next Q values (Q1 and Q2) using the
9     critic_target with next_state and next_action
10    # 3. Min over the Q values: target_Q = min(Q1, Q2) - log_prob(a'|
11    s') * alpha
12    # 4. Compute the target Q value: target_Q = reward + not_done *
13    discount * target_Q
14
15    #####
16    # YOUR IMPLEMENTATION HERE #
17    pass
18    #####
19
20    self.critic_optimizer.zero_grad()
21    critic_loss.backward()
22    self.critic_optimizer.step()
23
24    # [HINT] compute the actor loss
25    # 1. Sample the action and its log probability from the actor
26    with state
27    # 2. Compute the Q values (Q1 and Q2) using the critic with state
28    and action
29    # 3. Min over the Q values: Q = min(Q1, Q2)
30    # 4. Compute the actor loss: actor_loss = alpha * log_prob(a|s) -
31    Q
32
33    #####
34    # YOUR IMPLEMENTATION HERE #
35    pass
36    #####
37
38    self.actor_optimizer.zero_grad()
39    actor_loss.backward()
40    self.actor_optimizer.step()
41
42    for param, target_param in zip(self.critic.parameters(), self.
43    critic_target.parameters()):
```

```
35         target_param.data.copy_(self.tau * param.data + (1 - self.tau
           ) * target_param.data)
```

Question 1.c

Time to see how it works. We would expect a reward that converges to negative values close to 0. 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 50 episodes.

Solution

Training curve: