# Actor-Critic for Continuing Tasks

# Objectives

- [ ] Derive a sample-based estimate for the gradient

- [ ] Understand the Actor-Critic algorithm

# Estimating the Policy Gradient

☐ Definition:

- ■ Estimating the policy gradient is a crucial step in policy gradient algorithms

- ■ It allows us to update the policy parameters in the direction that maximizes the expected return.

# Estimating the Policy Gradient

☐ The gradient of the average reward:

$$\nabla r(\pi) = \sum_s \mu(s) \sum_a \nabla \pi(a \mid s, \boldsymbol{\theta}) q_\pi(s, a)$$

☐ We simply make updates from states we observe while following policy pi.

$$S_0, A_0, R_1, S_1, A_1, \ldots, S_t, A_t, R_{t+1}, \ldots$$

Actor-Critic for Continuing Tasks

4

# Estimating the Policy Gradient

☐ This gradient from state St provides an approximation to the gradient of the average reward.

☐ The stochastic gradient descent update looks like for the

$$\nabla r(\pi) = \sum_s \mu(s) \sum_a \nabla \pi(a \mid s, \boldsymbol{\theta}) q_\pi(s, a)$$

$$\boldsymbol{\theta}_{t+1} \doteq \boldsymbol{\theta}_t + \alpha \sum_a \nabla \pi(a \mid S_t, \boldsymbol{\theta}_t) q_\pi(S_t, a)$$

$$S_0, A_0, R_1, S_1, A_1, \ldots, S_t, A_t, R_{t+1}, \ldots$$

Actor-Critic for Continuing Tasks

5

# Estimating the Policy Gradient

☐ Estimating the policy gradient methods

    ☐ Finite Difference Methods: These methods approximate the gradient by perturbing the policy parameters and observing the change in the expected return.

    ☐ Score Function Methods: Score function methods, directly estimate the gradient of the expected return using samples obtained from interactions with the environment.

# Estimating the Policy Gradient

☐ Estimating the policy gradient methods

☐ Likelihood Ratio Methods: Likelihood ratio methods, further improve the efficiency of gradient estimation by subtracting a baseline function from the returns before computing the gradient.

☐ Actor-Critic Methods: Actor-critic methods combine policy gradient estimation with value function estimation.

# Estimating the Policy Gradient

❑ Estimating the policy gradient methods

    ❑ Natural Policy Gradient: Natural policy gradient methods incorporate the geometry of the parameter space into gradient estimation to improve convergence properties and stability.

# Actor-Critic Algorithm

- ☐ The Actor-Critic algorithm combines elements of both value-based and policy-based methods.
- ☐ It consists of two main components: the actor and the critic:
    - ☐ **Actor**: The actor is responsible for learning the policy
    - ☐ **Critic**: The critic is responsible for evaluating the policy.
- ☐ The Actor-Critic algorithm operates by iteratively updating the actor and critic networks based on observed experiences in the environment.

# Actor-Critic Algorithm

☐ Actor- Critic Algorithm:

  ☐ Step 1- Initialization: Initialize the actor and critic networks with random parameters.

  ☐ Step 2- Interact with Environment: Sample trajectories by following the current policy in the environment.

  ☐ Step 3- Compute Returns: Compute the returns for each time step in the trajectory.

  ☐ Step 4- Update Critic: Use the returns to update the parameters of the critic network

# Actor-Critic Algorithm

☐ Actor- Critic Algorithm:

    ☐ Step 5- Compute Advantages: Compute advantages for each time step by subtracting the estimated value from the observed return.

    ☐ Step 6- Update Actor: Update the parameters of the actor network using policy gradients

    ☐ Step 7- Repeat: Repeat steps 2-6 for multiple episodes or until convergence criteria are met.

# Actor-Critic Algorithm

☐ The Actor-Critic algorithm in a simple grid world environment

```python
1 # 3.11 Actor-Critic Algorithm- HoaDNt@fe.edu.vn
2 import numpy as np
3
4 class GridWorld:
5     def __init__(self):
6         self.grid_size = (3, 3)
7         self.num_actions = 4  # Up, Down, Left, Right
8         self.start_state = (0, 0)
9         self.goal_state = (2, 2)
10
11    def step(self, state, action):
12        # Define the dynamics of the environment
13        row, col = state
14        if action == 0:  # Up
15            row = max(0, row - 1)
16        elif action == 1:  # Down
17            row = min(self.grid_size[0] - 1, row + 1)
18        elif action == 2:  # Left
19            col = max(0, col - 1)
20        elif action == 3:  # Right
21            col = min(self.grid_size[1] - 1, col + 1)
22        next_state = (row, col)
23        reward = 0
24        if next_state == self.goal_state:
25            reward = 1  # Reward of +1 upon reaching the goal state
26        return next_state, reward
27
```

# Actor-Critic Algorithm

```python
28  class ActorCritic:
29      def __init__(self, num_actions, alpha_actor, alpha_critic, gamma):
30          self.num_actions = num_actions
31          self.alpha_actor = alpha_actor
32          self.alpha_critic = alpha_critic
33          self.gamma = gamma
34          self.actor_params = np.zeros((3, 3, num_actions))  # Tabular actor parameters
35          self.critic_values = np.zeros((3, 3))  # Tabular critic values
36
37      def select_action(self, state):
38          # Select action probabilistically based on actor parameters
39          action_probs = self.softmax(self.actor_params[state])
40          action = np.random.choice(self.num_actions, p=action_probs)
41          return action
42
43      def update(self, state, action, reward, next_state):
44          # Compute TD error (advantage)
45          td_error = reward + self.gamma * self.critic_values[next_state] - self.critic_values[state]
46
47          # Update critic values
48          self.critic_values[state] += self.alpha_critic * td_error
49
50          # Update actor parameters
51          self.actor_params[state][action] += self.alpha_actor * td_error
52
53      def softmax(self, x):
54          e_x = np.exp(x - np.max(x))
55          return e_x / e_x.sum(axis=0)
```

Actor-Critic for Continuing Tasks

# Actor-Critic Algorithm

```python
57 # Create a grid world environment
58 grid_world = GridWorld()
59
60 # Create an Actor-Critic agent
61 num_actions = 4  # Up, Down, Left, Right
62 alpha_actor = 0.1
63 alpha_critic = 0.1
64 gamma = 0.9
65 actor_critic_agent = ActorCritic(num_actions, alpha_actor, alpha_critic, gamma)
66
67 # Train the Actor-Critic agent
68 num_episodes = 1000
69 for _ in range(num_episodes):
70     state = grid_world.start_state
71     while state != grid_world.goal_state:
72         action = actor_critic_agent.select_action(state)
73         next_state, reward = grid_world.step(state, action)
74         actor_critic_agent.update(state, action, reward, next_state)
75         state = next_state
76
77 # Evaluate the learned policy
78 total_reward = 0
79 state = grid_world.start_state
80 while state != grid_world.goal_state:
81     action = actor_critic_agent.select_action(state)
82     next_state, reward = grid_world.step(state, action)
83     total_reward += reward
84     state = next_state
85
86 print("Total reward obtained by learned policy:", total_reward)
```

Actor-Critic for Continuing Tasks

14

# Actor-Critic Algorithm

☐ Run that code and show a total reward obtained by learned policy

# Summary

☐ Derive a sample-based estimate for the gradient

☐ Understand the Actor-Critic algorithm

# Q & A