

# ECE433/COS435 Introduction to RL

## Assignment 5: Proximal Policy Optimization (PPO)

### Spring 2025

Fill me in

Your name here.

Due March 28, 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.**

## Question 1.

In this assignment, we will implement Proximal Policy Optimization (PPO) for the CartPole-v1 environment from OpenAI Gym. Along the way, we will incorporate several important policy gradient concepts:

- REINFORCE (Monte Carlo Policy Gradient)
- Importance Weighting (IW) (For Off-Policy Learning)
- Generalized Advantage Estimation (GAE) (To Balance Bias and Variance)

### Question 1.a

First, you need to implement the policy network, which decides the actions to take, and the value network, which estimates the returns. Paste your class below:

## Solution

```
1 class PolicyNetwork(nn.Module):
2     def __init__(self, state_dim, action_dim):
3         super(PolicyNetwork, self).__init__()
4         self.fc1 = nn.Linear(state_dim, 128)
5         self.fc2 = nn.Linear(128, 128)
6         self.fc3 = nn.Linear(128, action_dim)
7
8     def forward(self, x):
9         x = torch.relu(self.fc1(x))
10        x = torch.relu(self.fc2(x))
11        x = self.fc3(x)
12        return Categorical(logits=x)
```

```
1 class ValueNetwork(nn.Module):
2     def __init__(self, state_dim):
3         super(ValueNetwork, self).__init__()
4         self.fc1 = nn.Linear(state_dim, 128)
5         self.fc2 = nn.Linear(128, 128)
6         self.fc3 = nn.Linear(128, 1)
7
8     def forward(self, x):
9         x = torch.relu(self.fc1(x))
10        x = torch.relu(self.fc2(x))
11        x = self.fc3(x)
12        return x
```

## Question 1.b

Below you need to implement GAE. The function should return the advantage estimates for each state-action pair. Paste below:

## Solution

```
1 def compute_advantages(next_value, rewards, masks, values, gamma
2     =0.99, lambda_gae=0.95):
3     values = values + [next_value] # Append bootstrap value for last
4     state
5     advantages = 0
6     returns = [] # Store computed returns
7
8     for step in reversed(range(len(rewards))): # Iterate in reverse
9         (backward pass)
10        delta = rewards[step] + gamma * values[step + 1] * masks[step]
11        - values[step]
12        advantages = delta + gamma * lambda_gae * masks[step] *
13        advantages
14        returns.insert(0, advantages + values[step])
```

```
11     return returns
```

## Question 1.c

Now let's move to the training parts. You need to implement the function below which applies the core PPO algorithm, uses the experiences collected from the environment to perform multiple epochs of updates on the policy and value networks. Paste below:

### Solution

```
1 def ppo_update(policy_net, value_net, optimizer, ppo_epochs,
2               mini_batch_size, states, actions, log_probs, returns, advantages,
3               clip_param=0.2):
4     for _ in range(ppo_epochs):
5         for state, action, old_log_probs, return_, advantage in
6             ppo_iter(mini_batch_size, states, actions, log_probs, returns,
7                     advantages):
8             dist = policy_net(state)
9             new_log_probs = dist.log_prob(action)
10
11             ##### Code implementation here #####
12             ratio = (new_log_probs-old_log_probs).exp()
13             clipped_ratio = torch.clamp(ratio, 1.0-clip_param, 1.0+
14                                         clip_param)
15             actor_loss = -torch.min(ratio*advantage, clipped_ratio*
16                                     advantage).mean()
17             value = value_net(state).squeeze(1)
18             critic_loss = (return_-value).pow(2).mean()
19             ##### Code implementation End #####
20
21             loss = 0.5 * critic_loss + actor_loss
22
23             optimizer.zero_grad()
24             loss.backward()
25             optimizer.step()
```

## Question 1.d

Finally, you should be able to run the main PPO training. You can freely adjust all hyperparameters for better performance. The provided hyperparameters, if implemented correctly, should be able to make rewards close to/higher than 400.

Note: Please try several times if you think your code is correct, the learning curves can have some variances over different runs. Present the best run you can get.

## Solution

