

# Homework 4

4771 Reinforcement Learning (Spring 2026)

Submission deadline: 11:59pm, April 12

The latex template can be found [here](#). The starter code is [here](#). This starter code merges the implementations of DQN and PPO. We will focus on the PPO part, while the DQN part serves as a solution for HW3.

## 1 PPO for CartPole

In this homework, we implement PPO with the Monte Carlo estimator to play CartPole. The algorithm is shown in [Algorithm 1](#). It has a different policy update rule compared to the version discussed in the class. Define  $\varrho_\theta(s, a) = \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}$  and  $\hat{A}(s, a) = \hat{G}(s, a) - V_\phi(s)$ . **Instead of** updating the policy with

$$\theta \leftarrow \theta + \alpha \nabla_\theta \left( \varrho_\theta(s, a) \hat{A}(s, a) - \beta \text{KL}(\pi_\theta(\cdot|s), \pi_{\theta_k}(\cdot|s)) \right) \quad (1)$$

as in Page 37 of [this slide](#) (omitting batching), we update it with

$$\theta \leftarrow \theta + \alpha \nabla_\theta \left( \min\{\varrho_\theta(s, a) \hat{A}(s, a), \varrho_\theta^{\text{clip}}(s, a) \hat{A}(s, a)\} \right) \quad (2)$$

where  $\varrho_\theta^{\text{clip}}(s, a)$  is a projection of  $\varrho_\theta(s, a)$  to the range  $[1 - \epsilon, 1 + \epsilon]$  for some  $\epsilon$  usually chosen as 0.1 or 0.2. See [Algorithm 1](#) for more detailed pseudocode.

The update rules (1) and (2) are two variants of PPO (see Equation (7) and (8) of the [PPO paper](#)), and are usually called PPO-Penalty and PPO-Clip respectively ([spinning up webpage](#)). The KL penalty and the clipping operation both introduce a more stable policy updates. The former imposes a soft penalty, while the latter simply makes the gradient being zero if the policy change  $\varrho_\theta(s, a)$  is outside a certain range.

In the CartPole environment, we find that PPO-Clip is significantly more stable, so this will be the version you implement in this homework. Besides the policy update, please verify yourself the equivalence between [Algorithm 1](#) and the algorithm on in Page 37 of [this slide](#).

(a) (40%) Implement PPO-Clip to play CartPole, and paste the output plot here. Your score for this part will be:

$$\min \left\{ \frac{\text{maximum running average within 2000 episodes}}{475}, 1 \right\} \times 100\%.$$

Recommended choices of the hyperparameters are provided in the starter code.

(b) (5%) List the hyperparameters you use in (a):

Hyperparameter	Value
N ( $N$ in <a href="#">Algorithm 1</a> )	
M ( $M$ in <a href="#">Algorithm 1</a> )	
MINIBATCH_SIZE ( $B$ in <a href="#">Algorithm 1</a> )	
LR ( $\alpha$ in <a href="#">Algorithm 1</a> )	
EPS ( $\epsilon$ in <a href="#">Algorithm 1</a> )	

If you use any trick beyond what is described in [Algorithm 1](#), please briefly describe it here. Also, list any additional hyperparameters you introduced in the table above.

---

### Algorithm 1 PPO-Clip

---

- 1 **Parameters:**  $N, M, B, \alpha$  (learning rate),  $\lambda$  (GAE parameter),  $\beta$  (KL penalty coefficient),  $\epsilon$  (clipping parameter).
- 2 Initialize the policy network  $\pi_\theta$  with weights  $\theta$
- 3 Initialize the value network  $V_\phi$  with weights  $\phi$ .
- 4 Sample  $s_{\text{tmp}} \sim \rho$ , where  $\rho$  is the initial state distribution.
- 5 **for**  $k = 1, \dots$  **do**
- 6      $s_1 \leftarrow s_{\text{tmp}}$
- 7     **for**  $i = 1, 2, \dots, N$  **do**
- 8         Select action  $a_i \sim \pi_\theta(\cdot | s_i)$ .
- 9         Observe reward  $r_i$  and the next state  $s'_i$ .
- 10         If  $s'_i$  is a terminal state, sample  $s_{i+1} \sim \rho$ ; otherwise, set  $s_{i+1} = s'_i$ .
- 11      $s_{\text{tmp}} \leftarrow s_{N+1}$
- 12     Compute  $\hat{G}(s_i, a_i)$  for all  $i = 1, \dots, N$  by calling [Algorithm 2](#) with inputs  $\{s_i, a_i, r_i, s'_i\}_{i=1}^N$  and  $V_\phi$ .
- 13     Define  $\hat{\pi}(a_i | s_i) = \pi_\theta(a_i | s_i)$  for all  $i = 1, \dots, N$ .
- 14     Define  $\hat{A}(s_i, a_i) = \hat{G}(s_i, a_i) - V_\phi(s_i)$  for all  $i = 1, \dots, N$ .
- 15     Define  $\hat{V}(s_i) = r_i + \gamma V_\phi(s'_i) \mathbb{I}\{s'_i \text{ is not a terminal state}\}$  for all  $i = 1, \dots, N$ .
- 16      $\hat{G}(s_i, a_i), \hat{\pi}(a_i | s_i), \hat{A}(s_i, a_i), \hat{V}(s_i)$  should all be viewed as constant below—call `tensor.detach()` for them.
- 17     **for**  $j = 1, \dots, M$  **do**
- 18         Randomly sample a minibatch  $\mathcal{B} \subset \{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$  with size  $|\mathcal{B}| = B$ .
- 19         Update the policy network:
 
$$\theta \leftarrow \theta + \alpha \nabla_\theta \frac{1}{B} \sum_{(s,a,r,s') \in \mathcal{B}} \min \{ \varrho_\theta(s, a) \hat{A}(s, a), \varrho_\theta^{\text{clip}}(s, a) \hat{A}(s, a) \},$$
- 20         where  $\varrho_\theta(s, a) \triangleq \frac{\pi_\theta(a|s)}{\hat{\pi}(a|s)}$  and  $\varrho_\theta^{\text{clip}}(s, a) \triangleq \max\{1 - \epsilon, \min\{1 + \epsilon, \varrho_\theta(s, a)\}\}$ .
- 21         Update the value network:
 
$$\phi \leftarrow \phi - \alpha \nabla_\phi \frac{1}{B} \sum_{(s,a,r,s') \in \mathcal{B}} (V_\phi(s) - \hat{V}(s))^2.$$

---



---

### Algorithm 2 Monte Carlo Estimator

---

**Input:**  $\{(s_i, a_i, r_i, s'_i)\}_{i=1}^N$  and  $V_\phi$ .

$\hat{G}(s_N, a_N) \leftarrow r_N + \gamma V_\phi(s'_N) \mathbb{I}\{s'_N \text{ is not a terminal state}\}$

**for**  $i = N - 1, \dots, 1$  **do**

$\hat{G}(s_i, a_i) \leftarrow r_i + \gamma \hat{G}(s_{i+1}, a_{i+1}) \mathbb{I}\{s'_i \text{ is not a terminal state}\}$

**return**  $\{\hat{G}(s_i, a_i)\}_{i=1}^N$ .

---

## 2 Survey

(5%) Leave any feedback for the course.