

# Assignment 3

EE675: Introduction to Reinforcement Learning

April 10, 2024

## Instructions

- Kindly name your submission files as 'RollNo\_Name\_A3.ipynb', Marks will be deducted for all submissions that do not follow the naming guidelines.
- You are required to work out your answers and submit only the iPython Notebook. The code should be well commented and easy to understand as there are marks for this.
- Submissions are to be made through HelloIITK portal. Submissions made through mail will not be graded.
- Answers to the theory questions, if any, should be included in the notebook itself. While using special symbols use the  $\text{\LaTeX}$  mode
- Make sure your plots are clear and have title, legends and clear lines, etc.
- Plagiarism of any form will not be tolerated. If your solutions are found to match with other students or from other uncited sources, there will be heavy penalties and the incident will be reported to the disciplinary authorities.
- In case you have any doubts, feel free to reach out to TAs for help.

## Part-A (Deadline - 20th Apr 2024)

**(Cart Pole Balancing using Policy Gradient: REINFORCE)** [20 Marks] Through this Cart Pole balancing exercise we will learn policy gradient algorithms. Consider the cart-pole problem where a pole is attached by an un-actuated joint to a cart, which moves along a frictionless track. The pendulum is placed upright on the cart and the goal is to balance the pole by applying forces in the left and right direction on the cart. For information about the observation\_space, action\_space and rewards of the environment refer the [cart pole documentation](#).

For this part implement the policy gradient algorithm REINFORCE as shown below in the figure. You are required to use a **linear policy** (say parameterized by  $\theta = [\theta_1 \ \theta_2]^\top$ , a 2x4 matrix for the cart pole problem) of state  $s$  passing through a softmax i.e.

$$\pi(a|s, \theta) = \text{softmax}(\theta \cdot s) = [\exp\{\theta_1^\top s\} \ \exp\{\theta_2^\top s\}]^\top / \exp\{\theta_1^\top s\} + \exp\{\theta_2^\top s\}$$

**REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for  $\pi_*$** 

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$   
 Algorithm parameter: step size  $\alpha > 0$   
 Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):  
   Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$   
   Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :  
      $G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$  ( $G_t$ )  
      $\theta \leftarrow \theta + \alpha \gamma^t G \nabla \ln \pi(A_t|S_t, \theta)$

Figure 1: REINFORCE

1. Write the expression for the gradient  $\nabla \ln \pi(A_t|S_t, \theta)$  for the policy shown above (write in  $\text{\LaTeX}$  in the jupyter-notebook submission) [2 Marks]
2. Implement REINFORCE algorithm with appropriate choice of algorithm parameters [5 Marks]
3. Plot the training rewards over 1000 episodes [2 Marks]
4. Test the trained policy and compute the average reward over 5 episodes [1 Marks]
5. Implement REINFORCE algorithm with baseline as shown in Figure 2 with appropriate choice of algorithm parameters. Take the state value function as a linear function of the state [7 Marks]

$$\hat{v}(s, \mathbf{w}) = \mathbf{w}^\top s$$

6. Compare and plot the training performance for both the algorithms [3 Marks]

**REINFORCE with Baseline (episodic), for estimating  $\pi_\theta \approx \pi_*$** 

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$   
 Input: a differentiable state-value function parameterization  $\hat{v}(s, \mathbf{w})$   
 Algorithm parameters: step sizes  $\alpha^\theta > 0$ ,  $\alpha^\mathbf{w} > 0$   
 Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$  (e.g., to  $\mathbf{0}$ )

Loop forever (for each episode):  
   Generate an episode  $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$ , following  $\pi(\cdot|\cdot, \theta)$   
   Loop for each step of the episode  $t = 0, 1, \dots, T - 1$ :  
      $G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k$  ( $G_t$ )  
      $\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$   
      $\mathbf{w} \leftarrow \mathbf{w} + \alpha^\mathbf{w} \delta \nabla \hat{v}(S_t, \mathbf{w})$   
      $\theta \leftarrow \theta + \alpha^\theta \gamma^t \delta \nabla \ln \pi(A_t|S_t, \theta)$

Figure 2: REINFORCE with baseline

**Note:** Take each episode to have a maximum of 500 steps

## Part-B (Optional)

**(Cart Pole Balancing using Policy Gradient: Actor Critic)** [10 Marks] Implement the Actor-Critic algorithm shown in figure 3 for the Cart Pole environment. Compare and plot the training performance vs methods in Part A.

### One-step Actor–Critic (episodic), for estimating $\pi_{\theta} \approx \pi_*$

Input: a differentiable policy parameterization  $\pi(a|s, \theta)$   
Input: a differentiable state-value function parameterization  $\hat{v}(s, \mathbf{w})$   
Parameters: step sizes  $\alpha^{\theta} > 0$ ,  $\alpha^{\mathbf{w}} > 0$   
Initialize policy parameter  $\theta \in \mathbb{R}^{d'}$  and state-value weights  $\mathbf{w} \in \mathbb{R}^d$  (e.g., to  $\mathbf{0}$ )  
Loop forever (for each episode):  
    Initialize  $S$  (first state of episode)  
     $I \leftarrow 1$   
    Loop while  $S$  is not terminal (for each time step):  
         $A \sim \pi(\cdot|S, \theta)$   
        Take action  $A$ , observe  $S', R$   
         $\delta \leftarrow R + \gamma \hat{v}(S', \mathbf{w}) - \hat{v}(S, \mathbf{w})$       (if  $S'$  is terminal, then  $\hat{v}(S', \mathbf{w}) \doteq 0$ )  
         $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S, \mathbf{w})$   
         $\theta \leftarrow \theta + \alpha^{\theta} I \delta \nabla \ln \pi(A|S, \theta)$   
         $I \leftarrow \gamma I$   
         $S \leftarrow S'$

Figure 3: Actor-Critic algorithm