

# Homework III

Deadline: 2025-01-05

1. (10 pts) Recall the definition of state visitation measure

$$d_{\mu}^{\pi}(s) = \mathbb{E}_{s_0 \sim \mu} [d_{s_0}^{\pi}(s)] = \mathbb{E}_{s_0 \sim \mu} \left[ (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}[s_t = s | s_0, \pi] \right],$$

where  $(s_0, a_0, s_1, a_1, \dots)$  is trajectory starting from initial distribution  $\mu$  and then following policy  $\pi$ . Let  $T$  obey the geometric distribution, i.e.,  $\mathbb{P}[T = t] = \gamma^t(1 - \gamma)$ ,  $t = 0, 1, \dots$ . Show that

$$\mathbb{P}[s_T = s] = d_{\mu}^{\pi}(s).$$

Then suggest a way to sample from  $d_{\mu}^{\pi}$ .

2. (20 pts) Implement and test the Projected Policy Gradient method and the Softmax Policy Gradient method in Lecture 7 for the Gridworld problem in Homework I (Question 7, use  $\gamma = 0.9$  and uniform distribution for  $\mu$ ). The action/advantage values and visitation measure in the policy gradient should be evaluated exactly based on the transition model. Display the convergence plots ( $V^*(\mu) - V^k(\mu)$  vs # of iterations) of the two algorithms in a figure. Can you observe the finite iteration convergence of the Projected Policy Gradient method?
3. (5 pts) Let  $V_{\tau}^{\pi}$  and  $Q_{\tau}^{\pi}$  be the value functions under the entropy regularization, and recall the definition of the Bellman optimality operator  $\mathcal{T}_{\tau}$  in this case. Show that

$$\mathcal{T}_{\tau} V_{\tau}^{\pi}(s) - V_{\tau}^{\pi}(s) = \tau \text{KL}(\pi(\cdot|s) \|\widehat{\pi}(\cdot|s)),$$

where  $\widehat{\pi}(\cdot|s) \propto \exp(Q_{\tau}^{\pi}(\cdot|s)/\tau)$ .

4. Consider the soft policy iteration algorithm in Lecture 8 (page 28).

- (10 pts) Show the policy improvement property of the algorithm:

$$V_{\lambda}^{\pi^{k+1}}(s) \geq V_{\lambda}^{\pi^k}(s), \quad \forall s.$$

- (10 pts) Show the  $\gamma$ -rate convergence of the algorithm:

$$\|V_{\lambda}^* - V_{\lambda}^{\pi^k}\|_{\infty} \leq \gamma^k \|V_{\lambda}^* - V_{\lambda}^{\pi^0}\|_{\infty}.$$

5. (20 pts) Reproduce the figure on page 26 of Lecture 9 for comparing different bandit algorithms.