

# Homework 3: Policy Gradient, Policy Evaluation, & Exploration

CS 6789: Foundations of Reinforcement Learning

## 1 Exploration in Absorbing MDPs (30 Points)

Let's first recall Problem 3 in HW2 where we constructed absorbing MDPs and have shown that the absorbing MDP can be used to provide optimism. We complete the rest exercise regarding absorbing MDPs in this section.

Specifically, we will see the Explore-or-Terminate phenomena, where the algorithm either has already found a near optimal policy—hence we terminate, or the algorithm can guarantee to explore new state-action pairs outside the current known set.

### 1.1 Exploration or Termination

Note that we know  $\hat{P}^\dagger$ , as it uses our learned model  $\hat{P}$  and it has deterministic absorbing structure at unknown state-action pairs. So we can plan in this model using  $r^\dagger$ . Let's define  $\hat{\pi}^*$  as the optimal policy of  $\hat{P}^\dagger$  and  $r^\dagger$  (imagine we run PI under  $\hat{P}^\dagger$  and  $r^\dagger$  until convergence).

We will consider two cases.

In case one, when we execute  $\hat{\pi}^*$  under the real model  $P$ ,  $\hat{\pi}^*$  does not escape the known set  $\mathcal{K}$  with big probability, i.e.,  $\sum_{s,a \in (\mathcal{S} \times \mathcal{A}) \setminus \mathcal{K}} d^{\hat{\pi}^*}(s, a) \leq \epsilon$  for some small  $\epsilon \in \mathbb{R}^+$ .

**Q (a) (15 points):** Let us prove the following inequality:

$$V^{\pi^*} \leq V^{\hat{\pi}^*} + \mathcal{O} \left( \frac{1}{(1-\gamma)^2} \epsilon + \frac{1}{(1-\gamma)^2} \sqrt{\frac{S \ln(SA/\delta)}{k}} \right).$$

This says that if  $k$  is big enough, and the escape probability is small, then we find a near-optimal policy  $\hat{\pi}^*$ .

**(Hint:** here we need to consider the difference between  $\hat{V}^{\dagger, \pi}$  and  $V^{\dagger, \pi}$  using simulation lemma)

**Q (b) (15 points):** On the other hand, we encounter the case where  $\hat{\pi}^*$  escapes, i.e.,  $\sum_{s,a \in (\mathcal{S} \times \mathcal{A}) \setminus \mathcal{K}} d^{\hat{\pi}^*}(s, a) \geq \epsilon$ . Let us assume that we have sampling oracle  $\mathcal{O}(\pi)$  that samples an i.i.d state-action pair from  $d^\pi$ , i.e.,  $(s, a) \sim d^\pi$

Prove the following conclusion. *With probability at least  $1 - \delta$ , after  $M = \Omega\left(\frac{kSA}{\epsilon} + \frac{\ln(1/\delta)}{\epsilon^2}\right)$  many sampling oracle  $\mathcal{O}(\hat{\pi}^*)$  calls, at least one state-action pair from the unknown set  $(\mathcal{S} \times \mathcal{A}) \setminus \mathcal{K}$  becomes known, i.e., it will be sampled more than  $k$  times.*

**(Hint:** First let's ask ourselves that after  $M$  many calls of  $\mathcal{O}(\hat{\pi}^*)$ , how many  $(s, a)$  samples out of  $M$  samples will fall into  $(\mathcal{S} \times \mathcal{A}) \setminus \mathcal{K}$ ? Now, if I have  $N$  many  $(s, a)$  samples falling into  $(\mathcal{S} \times \mathcal{A}) \setminus \mathcal{K}$ , can I say anything useful here using the Pigeon hole principle here?)

## 1.2 Summary (no questions)

To this end, we have seen that using  $\hat{P}^\dagger$  and  $r^\dagger$ , we plan a policy  $\hat{\pi}^*$  that either is near optimal if it stays inside the known set with big probability, or we guarantee to explore, i.e., making at least one previous unknown state-action pair to be known, by executing  $\hat{\pi}^*$  in the real MDP for polynomially many times. With the new known set, we can repeat the above whole process again which will either terminate with a near optimal policy or will guarantee to identify another known state-action pair. Note that the total number of state-action pairs that can be made to the known set is at most  $SA$ . Thus, we can expect that the algorithm will find a near optimal policy in polynomial sample complexity.

## 2 Policy Gradient Theorem (30 Points)

We will focus on infinite horizon discounted MDP here.

### 2.1 Baselines

**Q (15 points):** For any function  $f(s)$  (that is only a function of the states) show that:

$$\nabla V^{\pi_\theta} = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[ (Q^{\pi_\theta}(s, a) - f(s)) \nabla \log \pi_\theta(a|s) \right]$$

### 2.2 Action-dependent Baselines for Factored Policies

We consider the following factored policy representation. We assume action  $a \in \mathbb{R}^d$ . We denote  $a[i]$  as the  $i$ -th component of the  $d$ -dim vector  $a$ . Also we denote  $a[-i]$  as all the components of  $a$  except the  $i$ -th component.

Following the most common open-source implementations in the RL community, we represent our parameterized policy as follows  $\pi_\theta(a|s) = \prod_{i=1}^d \pi_{i;\theta}(a[i]|s)$ , where each  $\pi_{i;\theta}$  maps state  $s$  to a distribution in  $\mathbb{R}$ . Note that one of the most common examples of this in practice is in continuous control where  $\pi_\theta(\cdot|s)$  is a Gaussian distribution with a diagonal covariance matrix (i.e., each coordinate of the action is independent).

Now let us consider  $d$  many baselines,  $f_1, \dots, f_d$ , where each  $f_i(s, a[-i]) \in \mathbb{R}$ . Namely, each  $f_i$  takes both state and action information, but  $f_i$  ignores the  $i$ -th component of the action vector  $a$ .

**Q (15 points):** Prove the following:

$$\nabla V^{\pi_\theta} = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_\theta}} \mathbb{E}_{a \sim \pi_\theta(\cdot|s)} \left[ \sum_{i=1}^d \nabla_\theta \ln \pi_{i;\theta}(a[i]|s) (Q^{\pi_\theta}(s, a) - f_i(s, a[-i])) \right].$$

Namely, we indeed can use action-dependent baselines if our policy is factorized over dimensions.

### 2.3 Best Baseline (Bonus)

**Q** Find the  $f(\cdot)$  which leads to the minimum variance estimate of our gradient, assuming that our gradient estimate will be:

$$\widehat{\nabla V^{\pi_\theta}} = \frac{1}{1-\gamma} (Q^{\pi_\theta}(s, a) - f(s)) \nabla \log \pi_\theta(a|s)$$

where  $s \sim d^{\pi_\theta}$ ,  $a \sim \pi_\theta(\cdot|s)$ .

## 3 Batch Policy Evaluation via Regression (40 Points)

Often in practice, instead of using an unbiased estimate of policy gradient, one would like to use a biased estimate, i.e., one would like to perform bias-variance tradeoff for policy gradient. In this homework problem, let us derive a specific algorithm that performs policy evaluation in a dynamic programming fashion.

Let's take a look at the setup below. We consider finite horizon setting  $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, H, r, P, \rho\}$  with  $\rho$  being the initial state distribution. We consider offline setting where we have a  $H$  many datasets  $\mathcal{D}_h = \{s_h^i, a_h^i, \bar{s}_h^i, r_h^i\}_{i=1}^N$ , where  $s_h^i, a_h^i \sim \nu_h$ , and  $r_h^i = r(s_h^i, a_h^i)$ , and  $\bar{s}_h^i \sim P(\cdot|s_h^i, a_h^i)$ ; We consider linear function approximation with feature mapping  $\phi : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}^d$ .

To make sure the problem is tractable, we make the following two key assumptions:

- Assumption 1: Realizability, i.e., for any  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , we have for all  $h$ ,  $Q_h^\pi(s, a) = \theta_h^\top \phi(s, a)$ , for some  $w$  which has bounded norm  $\|\theta_h\|_2 \leq W \in \mathbb{R}^+$ ; In other words,  $Q_h^\pi$  is a linear function with respect to feature  $\phi$ , for all  $h$ ;
- Assumption 2: Coverage: i.e., for data distribution  $\nu_h$  at any  $h$ , we have  $\sigma_{\min}(\mathbb{E}_{s, a \sim \nu_h} \phi(s, a) \phi(s, a)^\top) \geq \kappa \in \mathbb{R}^+$ . In other words, the data distribution  $\nu_h$ 's induced feature covariance matrix is full rank and has lower bounded eigenvalues. This assumption ensures that with reasonably enough data, our data matrix, i.e.,  $\Phi_h = [\phi(s_h^1, a_h^1), \dots, \phi(s_h^N, a_h^N)]^\top \in \mathbb{R}^{N \times d}$  will be full rank.

We fix a policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ , and we would like to estimate  $Q_0^\pi$  and its parameter  $\theta_0$  using the offline datasets  $\mathcal{D}_h$  for  $h = 0, \dots, H-1$ .

### 3.1 Least square regression at $h = H - 1$

From a dynamic programming perspective, it is obvious what we would like to do at  $h = H - 1$ , i.e., the last time step. We use  $\mathcal{D}_{H-1}$  to perform regression:

$$\hat{w}_{H-1} = \arg \min_w \sum_{i=1}^N (\phi(s_{H-1}^i, a_{H-1}^i)^\top w - r_{H-1}^i)^2.$$

Let's write  $\Phi_{H-1} \in \mathbb{R}^{N \times d}$ , where the  $i$ -th row of  $\Phi_{H-1}$  is  $\phi(s_{H-1}^i, a_{H-1}^i)^\top$ , and let's write  $r_{H-1} \in \mathbb{R}^N$  where the  $i$ -th element of  $r_{H-1}$  is  $r_{H-1}^i$ .

Assume  $\Phi_{H-1}$  is full rank (i.e., assume  $N > d$ , and rank of  $\Phi_{H-1}$  is  $d$ )

**Q (10 points):** Derive that:

$$\hat{w}_{H-1} = (\Phi_{H-1}^\top \Phi_{H-1})^{-1} \Phi_{H-1}^\top r_{H-1}.$$

### 3.2 Linear square regression at $h$

Now let's do dynamic programming. Assume that we have learned  $\hat{w}_{h+1}$  already. We use  $\hat{w}_{h+1}$  to learn  $\hat{w}_h$ . Specifically, we consider the following least square regression using the dataset  $\mathcal{D}_h$ :

$$\hat{w}_h = \arg \min_w \sum_{i=0}^N (\phi(s_h^i, a_h^i)^\top w - (r_h^i + \hat{w}_{h+1}^\top [\phi(\bar{s}_h^i, \pi(\bar{s}_h^i))]))^2.$$

Note that if we believed that  $\hat{w}_{h+1}$  is a good estimator, then we would believe that  $r_h^i + \hat{w}_{h+1}^\top \phi(\bar{s}_h^i, \pi(\bar{s}_h^i)) \approx Q_h^\pi(s_h^i, a_h^i)$ , since  $s_{h+1}^i \sim P(\cdot | s_h^i, a_h^i)$ . From an induction perspective, we started at  $H - 1$  where we should expect that  $\hat{w}_{H-1}$  is a good fit because the reward itself is a linear function with respect to  $\phi$ , and if we were doing linear regression correctly, we should expect  $\hat{w}_{h+1}$  is a good estimator as well, as  $Q_{h+1}^\pi$  is a linear function with respect to  $\phi$ .

Let's denote  $\Phi_h \in \mathbb{R}^{N \times d}$  as a matrix where the  $i$ -th row is  $\phi(s_h^i, a_h^i)^\top$ , and  $\bar{\Phi}_h \in \mathbb{R}^{N \times d}$  as a matrix where the  $i$ -th row is  $\phi(\bar{s}_h^i, \pi(\bar{s}_h^i))$ , and  $r_h \in \mathbb{R}^N$  where the  $i$ -th element of  $r$  is  $r_h^i$ .

**Q (10 points):** Derive that:

$$\hat{w}_h = (\Phi_h^\top \Phi_h)^{-1} \Phi_h^\top r_h + (\Phi_h^\top \Phi_h)^{-1} \Phi_h^\top \bar{\Phi}_h \hat{w}_{h+1}.$$

### 3.3 Solution at $h = 0$

We can repeat the same process above for all  $h = H - 1$  to  $h = 0$ .

**Q (20 Points):** Show that for  $\hat{w}_0$ :

$$\hat{w}_0 = \sum_{h=0}^{H-1} \left( \prod_{\tau=0}^{h-1} (\Phi_\tau^\top \Phi_\tau)^{-1} \Phi_\tau^\top \bar{\Phi}_\tau \right) (\Phi_h^\top \Phi_h)^{-1} \Phi_h^\top r_h$$

### 3.4 Quality of the Solution (Bonus)

We further simplify the setting by assuming that  $r_h^i = 0$  for all  $h = 0 \dots H - 2$ , i.e., we have nonzero reward at  $H - 1$  but zero reward in all previous time steps (i.e., we get reward 1 or 0 at the end of the game). In this case, the solution  $\hat{w}_0$  has a very simple expression:

$$\hat{w}_0 = \left( \prod_{\tau=0}^{H-2} (\Phi_\tau^\top \Phi_\tau)^{-1} \Phi_\tau^\top \bar{\Phi}_\tau \right) (\Phi_{H-1}^\top \Phi_{H-1})^{-1} \Phi_{H-1}^\top r_{H-1}$$

Given  $\phi(s_0, \pi(s_0))$ , we are interested in the prediction error  $|Q_0^\pi(s_0, \pi(s_0)) - \hat{w}_0 \phi(s_0, \pi(s_0))|$  at the initial state  $s_0$ .

Recall our assumption where for all  $h = 0, \dots, H - 1$ , we have that for all  $s, a$ ,  $Q_h^\pi(s, a) = \theta_h \cdot \phi(s, a)$  for some  $\theta_h$ .

Let us denote  $\delta_h^i = \theta_{h+1} \cdot \phi(\bar{s}_h^i, \pi(\bar{s}_h^i)) - \mathbb{E}_{s' \sim P(\cdot | s_h^i, a_h^i)} \theta_{h+1}^\top \phi(s', \pi(s'))$ . Since  $\bar{s}_h^i \sim P(\cdot | s_h^i, a_h^i)$ , we have that  $\mathbb{E}[\delta_h^i | s_h^i, a_h^i] = 0$ . Denote  $\delta_h \in \mathbb{R}^N$  as a vector whose  $i$ -th element is  $\delta_h^i$ .

**Q** Prove that:

$$\hat{w}_0 - \theta_0 = \sum_{h=0}^{H-1} \left[ \left( \prod_{\tau=0}^{h-1} (\Phi_\tau^\top \Phi_\tau)^{-1} \Phi_\tau^\top \bar{\Phi}_\tau \right) (\Phi_h^\top \Phi_h)^{-1} \Phi_h^\top \delta_h \right]$$

**Remark (no question)** As we can see above,  $\hat{w}_0$  is indeed an unbiased estimate of  $\theta_0$ ! This is because for any  $h$ , we have:

$$\mathbb{E} \left[ \left( \prod_{\tau=0}^{h-1} (\Phi_\tau^\top \Phi_\tau)^{-1} \Phi_\tau^\top \bar{\Phi}_\tau \right) (\Phi_h^\top \Phi_h)^{-1} \Phi_h^\top \delta_h \mid \mathcal{D}_0, \dots, \mathcal{D}_h \right] = 0,$$

as  $\mathbb{E}[\delta_h^i | s_h^i, a_h^i] = 0$ . However, the variance of our estimator  $\hat{w}_0$  could be exponential in  $H$  in worst case (notice the term  $\prod_{\tau=0}^{h-1} (\Phi_\tau^\top \Phi_\tau)^{-1} \Phi_\tau^\top \bar{\Phi}_\tau$ ). For people who are interested in the potential exponentially large variance worst case, we encourage you to think about a lower bound construction.

## 4 Linear Quadratic Regulator (Bonus)

Recall the definition of infinite horizon LQR:

$$\begin{aligned} x_{t+1} &= Ax_t + Bu_t, \\ c(x, u) &= x^\top Qx + u^\top Ru, \end{aligned}$$

where  $x \in \mathbb{R}^d$  and  $u \in \mathbb{R}^k$ .

We consider a specific linear policy:  $\pi(x) = Kx$  where  $K \in \mathbb{R}^{k \times d}$  (assume this control is stable, i.e., the average expected cost under this control is finite). From the lectures, we know that

the value function of the linear policy is quadratic. Let us denote  $P_K \in \mathbb{R}^{d \times d}$ , and define the value function  $V_K(x)$  as:

$$V_K(x) := \sum_{t=0}^{\infty} x_t^\top Q x_t + u_t^\top R u_t = x^\top P_K x,$$

where  $x_0 = x$ ,  $u_t = K x_t$ , and Q function  $Q_K(x, u)$  as:

$$Q_K(x, u) = x^\top Q x + u^\top R u + V_K(Ax + Bu).$$

## 4.1 Policy Iteration

Recall the Policy Iteration algorithm from our second lecture: Given a policy  $\pi : \mathbb{R}^d \mapsto \mathbb{R}^k$ , PI computes a new policy  $\pi'$ :

$$\pi'(x) := \arg \min_{u \in \mathbb{R}^k} Q^\pi(x, u), \forall x \in \mathbb{R}^d.$$

**Q:** Now use the fact that we are considering an LQR system with  $V_K(x) = x^\top P_K x$ , derive the closed-form solution of  $\pi'$

**(Hint:** first show that  $\pi'$  again will be a linear policy).