

# CS234 Problem Session Solutions

Week 6: Feb 17

## 1) [CA Session] Conservative Policy Iteration

Let us consider an MDP with a fixed start state  $s_0$ .

Let us consider the conservative policy update rule:

$$\pi_{new}(s, a) = (1 - \alpha)\pi(s, a) + \alpha\pi'(s, a)$$

for some  $\alpha \in [0, 1]$ .

(a) What is  $\pi_{new}(s, a)$  when  $\alpha = 1$ ?

**Solution**  $\pi_{new}(s, a) = \pi'(s, a)$

Recall that  $A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s)$ .

Let  $P(s_t; \pi)$  be the distribution over states at time  $t$  while following  $\pi$  from the start state  $s_0$ . Recall that the discounted stationary state distribution of a policy  $\pi$  is  $d^\pi(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s; \pi)$ . We now define the policy advantage of some policy  $\pi'$  with respect to a policy  $\pi$  as  $\mathbb{A}^\pi(\pi') = \mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi'(s)} [A^\pi(s, a)]]$ . Recall Lemma 1 from assignment 2.

**Lemma 1:** For all policies  $\pi', \pi$ , we have that  $V^{\pi'}(s_0) - V^\pi(s_0) = \frac{1}{1 - \gamma} \mathbb{E}_{s \sim d^{\pi'}} [\mathbb{E}_{a \sim \pi'(s)} [A^\pi(s, a)]]$

(b) How does  $V^{\pi'}(s_0) - V^\pi(s_0)$  differ from the policy advantage  $\mathbb{A}^\pi(\pi')$ ? A high-level description in words will suffice.

**Solution** The main difference is that the expectation over states is over a different distribution  $d^\pi$  in each case. The policy advantage is also not normalized by  $\frac{1}{1 - \gamma}$

(c) Compute a simplified expression for  $\mathbb{A}^\pi(\pi_{new})$  in terms of the policy advantage of  $\pi'$ .

**Solution**

$$\begin{aligned} \mathbb{A}^\pi(\pi_{new}) &= \mathbb{E}_{s \sim d^\pi} [\mathbb{E}_{a \sim \pi_{new}(s)} [A^\pi(s, a)]] \\ &= \mathbb{E}_{s \sim d^\pi} [(1 - \alpha)\mathbb{E}_{a \sim \pi(s)} [A^\pi(s, a)] + \alpha\mathbb{E}_{a \sim \pi'(s)} [A^\pi(s, a)]] \\ &= \mathbb{E}_{s \sim d^\pi} [\alpha\mathbb{E}_{a \sim \pi'(s)} [A^\pi(s, a)]] \\ &= \alpha\mathbb{A}^\pi(\pi') \end{aligned}$$

With  $\pi_{new}$ , at any given timestep, the probability that we select an action according to  $\pi'$  is  $\alpha$ . Let us define the random variable  $c_t$  as the number of actions chosen from  $\pi'$  before time  $t$ .

(d) Let us denote  $\rho_t = Pr(c_t \geq 1)$ . Compute an expression for  $\rho_t$  in terms of  $\alpha$  and  $t$ .

**Solution** We can see  $Pr(c_t = 0) = (1 - \alpha)^t$ . Thus,  $\rho_t = Pr(c_t \geq 1) = 1 - (1 - \alpha)^t$ .

Now let  $\epsilon = \max_s |\mathbb{E}_{a \sim \pi'(s)} [A^\pi(s, a)]|$ .

(e) Prove that  $\mathbb{E}_{s \sim P(s_t; \pi_{new})} [\sum_a \pi_{new}(s, a) A^\pi(s, a)] \geq \alpha \mathbb{E}_{s \sim P(s_t; \pi)} [\sum_a \pi'(s, a) A^\pi(s, a)] - 2\alpha\rho_t\epsilon$ .

**Solution**

$$\begin{aligned}
& \mathbb{E}_{s \sim P(s_t; \pi_{new})} [\sum_a \pi_{new}(s, a) A^\pi(s, a)] \\
&= \alpha \mathbb{E}_{s \sim P(s_t; \pi_{new})} [\sum_a \pi'(s, a) A^\pi(s, a)] \\
&= \alpha(1 - \rho_t) \mathbb{E}_{s \sim P(s_t | c_t = 0; \pi_{new})} [\sum_a \pi'(s, a) A^\pi(s, a)] + \alpha\rho_t \mathbb{E}_{s \sim P(s_t | c_t \geq 1; \pi_{new})} [\sum_a \pi'(s, a) A^\pi(s, a)] \\
&= \alpha(1 - \rho_t) \mathbb{E}_{s \sim P(s_t | c_t = 0; \pi_{new})} [\sum_a \pi'(s, a) A^\pi(s, a)] + \alpha\rho_t \mathbb{E}_{s \sim P(s_t | c_t \geq 1; \pi_{new})} [\sum_a \pi'(s, a) A^\pi(s, a)] \\
&\geq \alpha \mathbb{E}_{s \sim P(s_t | c_t = 0; \pi_{new})} [\sum_a \pi'(s, a) A^\pi(s, a)] - 2\alpha\rho_t\epsilon \\
&= \alpha \mathbb{E}_{s \sim P(s_t; \pi)} [\sum_a \pi'(s, a) A^\pi(s, a)] - 2\alpha\rho_t\epsilon
\end{aligned}$$

Notice that the last line holds because  $P(s_t | c_t = 0; \pi_{new}) = P(s_t; \pi)$ .

(f) Now let us lower bound the improvement of our policy. Please prove that the following equation holds:

$$V^{\pi_{new}}(s_0) - V^\pi(s_0) \geq \frac{\alpha}{1 - \gamma} (\mathbb{A}^\pi(\pi') - \frac{2\alpha\gamma\epsilon}{1 - \gamma(1 - \alpha)})$$

## Solution

$$\begin{aligned}
& V^{\pi_{new}}(s_0) - V^{\pi}(s_0) \\
&= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi_{new}}} [\mathbb{E}_{a \sim \pi_{new}(s)} [A^{\pi}(s, a)]] \\
&= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s \sim P(s_t; \pi_{new})} \left[ \sum_a \pi_{new}(s, a) A^{\pi}(s, a) \right] \\
&\geq \sum_{t=0}^{\infty} \gamma^t [\alpha \mathbb{E}_{s \sim P(s_t; \pi)} \left[ \sum_a \pi'(s, a) A^{\pi}(s, a) \right] - 2\alpha\rho_t\epsilon] \\
&= \alpha \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s \sim P(s_t; \pi)} \left[ \sum_a \pi'(s, a) A^{\pi}(s, a) \right] - 2\alpha\epsilon \sum_{t=0}^{\infty} \gamma^t (1 - (1-\alpha)^t) \\
&= \frac{\alpha}{1-\gamma} \mathbb{E}_{s \sim d^{\pi}} \left[ \sum_a \pi'(s, a) A^{\pi}(s, a) \right] - 2\alpha\epsilon \left( \frac{1}{1-\gamma} - \frac{1}{1-\gamma(1-\alpha)} \right) \\
&= \frac{\alpha}{1-\gamma} \left[ \mathbb{A}^{\pi}(\pi') - \frac{2\alpha\gamma\epsilon}{1-\gamma(1-\alpha)} \right]
\end{aligned}$$

## 2) [Breakout Rooms] Trajectory Likelihoods

Suppose  $\pi_1$  and  $\pi_2$  are two different stochastic policies. We now observe a trajectory  $H = (S_0, A_0, R_0, S_1, \dots, S_{T-1}, A_{T-1}, R_{T-1})$ . Assume the rewards are finite and denote  $R(s, a, s', r) = \Pr(R_t = r | S_t = s, A_t = a, S_{t+1} = s')$ .

(a) Simplify  $\frac{\Pr(H|\pi_1)}{\Pr(H|\pi_2)}$  using terms from the MDP definition. Your final answer should be able to be computed without needing to know the transition function, the reward function, or the reward distribution.

### Solution

$$\begin{aligned} \frac{\Pr(H|\pi_1)}{\Pr(H|\pi_2)} &= \frac{\Pr(S_0)\pi_1(S_0, A_0)P(S_0, A_0, S_1)R(S_0, A_0, S_1, R_0)\pi_1(S_1, A_1)P(S_1, A_1, S_2)\dots}{\Pr(S_0)\pi_2(S_0, A_0)P(S_0, A_0, S_1)R(S_0, A_0, S_1, R_0)\pi_2(S_1, A_1)P(S_1, A_1, S_2)\dots} \\ &= \frac{\pi_1(S_0, A_0)\pi_1(S_1, A_1)\pi_1(S_2, A_2)\dots}{\pi_2(S_0, A_0)\pi_2(S_1, A_1)\pi_2(S_2, A_2)\dots} \\ &= \prod_{t=0}^T \frac{\pi_1(S_t, A_t)}{\pi_2(S_t, A_t)} \end{aligned}$$

### 3) [Breakout Rooms] Off Policy Actor Critic Policy Gradients

We will derive an expression for the policy gradient for a new objective function,  $J'$ . This new objective is similar to one used in off-policy actor-critics. Assume there is a fixed policy  $\pi_b$ . Let

$$d'(s) = \sum_{t=0}^{L-1} Pr(S_t = s | \pi_b)$$

The objective function  $J'$  is defined as

$$J'(\theta) = \sum_{s \in S} d'(s) E[R_t | S_t = s, \theta]$$

Derive an expression for the policy gradient for this objective. The terms in your answer should only be terms used in defining an MDP (including the reward function defined as  $R(s, a)$ ). Note that  $\theta$  are not the parameters of  $\pi_b$ , but the parameters of another policy  $\pi$ .

#### Solution

$$\begin{aligned} \frac{\partial J'(\theta)}{\partial \theta} &= \frac{\partial}{\partial \theta} \sum_{s \in S} d'(s) \sum_{a \in A} \pi(s, a, \theta) \mathbb{E}[R_t | S_t = s, A_t = a, \theta] \\ &= \sum_{s \in S} \sum_{t=0}^{L-1} Pr(S_t = s | \pi_b) \sum_{a \in A} \frac{\partial}{\partial \theta} \pi(s, a, \theta) R(s, a) \\ &= \sum_{t=0}^{L-1} \sum_{s \in S} Pr(S_t = s | \pi_b) \sum_{a \in A} R(s, a) \frac{\partial}{\partial \theta} \pi(s, a, \theta) \\ &= \sum_{t=0}^{L-1} \mathbb{E} \left[ \sum_{a \in A} R(S_t, a) \frac{\partial}{\partial \theta} \pi(S_t, a, \theta) | \pi_b \right] \\ &= \mathbb{E} \left[ \sum_{t=0}^{L-1} \sum_{a \in A} R(S_t, a) \frac{\partial}{\partial \theta} \pi(S_t, a, \theta) | \pi_b \right] \end{aligned}$$

Problem is borrowed from <sup>1</sup>

---

<sup>1</sup>[https://people.cs.umass.edu/~pthomas/courses/CMPSCI\\_687\\_Fall2018/687\\_F18\\_main.pdf](https://people.cs.umass.edu/~pthomas/courses/CMPSCI_687_Fall2018/687_F18_main.pdf)