

Lecture 7: Advanced PG

So far, we saw:

Policy gradient methods

(policy, maybe on policy $V(s)$)
on-policy

Q-learning methods

(no policy, just a value fn.)

off-policy

Can we build a hybrid of these methods?

Should be possible!

Why?

How can we get started in a principled manner??

Performance Difference Lemma

For any two policies $\pi(a|s)$ & $\mu(a|s)$, the following holds:

$$\underbrace{J(\pi) - J(\mu)}_{\text{Return of policy } \pi} = \mathbb{E}_{\tau \sim P_\pi(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^{\mu}(s_t, a_t) \right]$$

Return of
policy π

Proof:

$$\begin{aligned} J(\pi) &= \mathbb{E}_{\tau \sim P_\pi(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \right] \\ &= \mathbb{E}_{\tau \sim P_\pi(\tau)} \left[V^\pi(s_0) \right] = \mathbb{E}_{\tau \sim P_\mu(\tau)} \left[V^\pi(s_0) \right] \\ &= \mathbb{E}_{\tau \sim P_\mu(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t V^\pi(s_t) - \sum_{t=1}^{\infty} \gamma^t V^\pi(s_t) \right] \\ &= \mathbb{E}_{\tau \sim P_\mu(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t (V^\pi(s_t) - \gamma V^\pi(s_{t+1})) \right] \end{aligned}$$

What does the performance difference lemma give us?

$$A^{\mu} \iff \pi$$

Policy iteration

$$\begin{aligned} \max_{\pi} J(\pi) &= \max_{\pi} (J(\pi) - J(\mu)) \\ &= \max_{\pi} E_{s \sim d^{\pi}} [A^{\mu}(s, a)] \end{aligned}$$

Why is this useful?

A practical off-policy policy gradient algorithm

$$J(\pi) - J(\mu) = \mathbb{E}_{\tau \sim P_\pi(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t A^\mu(s_t, a_t) \right]$$

data collection policy

estimated via a critic

Option 1: use states from μ , A^μ instead

Option 2: use states from π , A^π instead

Which one should we use?

An approximation

$$\max_{\pi} J(\pi) - J(\mu)$$

$$E_{\tau \sim p_{\pi}(s)} \left[\sum_{t=0}^{\infty} \gamma^t A^{\mu}(s_t, a_t) \right]$$

$$= \sum_{t=0}^{\infty} E_{\substack{s_t \sim p_{\pi}(s_t) \\ a_t \sim \pi(\cdot | s_t)}} \left[\gamma^t A^{\mu}(s_t, a_t) \right]$$

ignore utilize

$$E_{\substack{s_t \sim p_{\mu}(s_t) \\ a_t \sim \pi(\cdot | s_t)}} \left[\gamma^t A^{\mu}(s_t, a_t) \right]$$

only term involved in maximization

Different ways to solve this maximization

when π & μ are close!!

Generic optimization problem:

$$\max_{\pi} \mathbb{E}_{s \sim d^{\mu}(s)} \left[\mathbb{E}_{a \sim \pi(\cdot | s)} \left[A^{\mu}(s, a) \right] \right]$$

$$-\alpha D(\pi, \mu)$$

what is this precisely?

We want $P_{\pi}(s_t) \& P_{\mu}(s_t)$ to be close!

$$P_{\pi}(s_t) - P_{\mu}(s_t) = ??$$

say π & μ are deterministic & agree with probability $1-\varepsilon$, then:

$$P_{\pi}(s_t) = (1-\varepsilon)^t P_{\mu}(s_t) + (1-(1-\varepsilon)^t) \text{Paiff}(s_t)$$

$$\Rightarrow |P_{\pi}(s_t) - P_{\mu}(s_t)| = (1 - (1-\varepsilon)^t) |P_{\mu}(s_t) - \text{Paiff}(s_t)| \leq 2 \cdot \varepsilon \cdot t$$

This means that:

$$D(\pi, \mu) = \underline{D_{TV}(\pi(-|s), \mu(-|s))}$$

total variations distance.

$$D_{TV}(P, Q) \leq \sqrt{\frac{1}{2} D_{KL}(P \| Q)}$$

$$\& D_{TV}(P, Q) \leq \sqrt{\frac{1}{2} D_{KL}(Q \| P)}$$

Pinsker's inequality

Each of these choices gives a different algorithm:

① AWR $\rightarrow D_{KL}(\pi \| \mu)$

② PPO / TRPO $\rightarrow D_{KL}(\mu \| \pi)$

Advantage-weighted regression (AWR)

[Peng, Kumar, Liu, Levine, 2019]

$$\max_{\pi} \mathbb{E}_{s \sim d^{\mu}(s)} \left[\mathbb{E}_{a \sim \pi(\cdot | s)} [A^{\mu}(s, a)] - \alpha D_{KL}(\pi || \pi_0) \right]$$

Easier problem:

$$\max_{p(x)} \mathbb{E}_{x \sim p(x)} [f(x)] - \alpha D_{KL}(p || q)$$

$$\Rightarrow p^*(x) \propto q(x) \exp\left(\frac{f(x)}{\alpha}\right)$$

$$\text{So, } \pi^*(a|s) \propto \pi(a|s) \exp\left[\frac{A^{\mu}(s, a)}{\alpha}\right]$$

How do I use this oracle solution??