

Lec 7: Policy gradient methods for tabular MDP (Based on Chi Jin's lecture)

MDP: $\{P_h\}_{h \in [H]}, \{r_h\}_{h \in [H]}$

- objective function:

$$J(\pi) = V_1^\pi(s_1) = \mathbb{E}_\pi \left[\sum_{h=1}^H r_h(s_h, a_h) \mid s_1 \right]$$

where $(s_1, a_1, s_2, a_2, \dots, s_H, a_H) \sim (P, \pi)$

$$a_1 \sim \pi_1(\cdot \mid s_1)$$

$$s_2 \sim P_1(\cdot \mid s_1, a_1)$$

\vdots

$$a_h \sim \pi_h(\cdot \mid s_h)$$

$$s_{h+1} \sim P_h(\cdot \mid s_h, a_h)$$

- parameterization

$$\pi = (\pi_1, \dots, \pi_H) \quad \text{where } \pi_h(\cdot \mid s) \in \Delta(A)$$

Consider the natural policy gradient method (mirror descent)

- initialize $\pi_h^1(\cdot|s)$ to be uniform over A for all s, h

- for $t=1, 2, \dots, T$:

$$\pi^{t+1} = \text{proj}_{\Pi} \left(\pi^t + \frac{\eta}{2} \Delta J(\pi^t) \right)$$

$$= \arg \max_{\pi \in \Pi} \langle \Delta J(\pi^t), \pi - \pi^t \rangle - \frac{1}{\eta} \underbrace{D_{\Phi}(\pi, \pi^t)}$$

- return: π^1, \dots, π^T

$$\mathbb{E}_{\pi^t} \left[\sum_{h=1}^H \text{KL}[\pi_h(\cdot|s_h) \parallel \pi_h^t(\cdot|s_h)] \right]$$

- Compute the gradient

$$\frac{\partial J(\pi)}{\partial \pi_h(a|s)}$$

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{h'=1}^H r_{h'}(s_{h'}, a_{h'}) \right] = \underbrace{\mathbb{E}_{\pi} \left[\sum_{h'=1}^{h-1} r_{h'}(s_{h'}, a_{h'}) \right]}_{\text{independent of } \pi_h} + \underbrace{\mathbb{E}_{\pi} \left[\sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) \right]}_{(I)}$$

$$(I) = \mathbb{E} \left[Q_h^{\pi}(s, a) \right]$$

$s \sim \Pr(s_h = s | \pi)$
 $a \sim \pi_h(\cdot | s)$

$$\Rightarrow \frac{\partial J(\pi)}{\partial \pi_h(a|s)} = \underbrace{\Pr(s_h = s | \pi)}_{d_h^{\pi}(s)} Q_h^{\pi}(s, a)$$

Solve:

$$\arg \max_{\pi \in \Pi} \langle \Delta J(\pi^t), \pi - \pi^t \rangle - \frac{1}{\eta} D_{\Phi}(\pi, \pi^t)$$

Lagrangian multiplier:

$$L(\pi, \lambda) = \sum_{h,s,a} d_h^{\pi^t}(s,a) Q_h^{\pi^t}(s,a) (\pi_h(a|s) - \pi_h^t(a|s)) - \frac{1}{\eta} \sum_{h,s,a} d_h^{\pi}(s,a) \cdot \pi_h(a|s) \log \frac{\pi_h(a|s)}{\pi_h^t(a|s)} \\ - \sum_{h,s} \lambda_{h,s} \left(\sum_a \pi_h(a|s) - 1 \right)$$

$$\frac{\partial L(\pi, \lambda)}{\partial \pi_h(a|s)} = d_h^{\pi^t}(s,a) Q_h^{\pi^t}(s,a) - \frac{1}{\eta} d_h^{\pi}(s,a) \cdot \left(\log \frac{\pi_h(a|s)}{\pi_h^t(a|s)} + 1 \right) - \lambda_{h,s} = 0$$

$$\Rightarrow \pi_h(a|s) = \pi_h^t(a|s) e^{d_h^{\pi^t}(s,a) / Z_t}$$

Consider the simplest case: $\forall \pi \xrightarrow{\text{oracle}} Q_h^{\pi} \quad \forall h$

Theorem $\text{Regret}(T) = \sum_{t=1}^T [V_1^{\pi^*}(s_1) - V_1^{\pi^t}(s_1)] \leq \frac{H \log A}{\gamma} + H^2$

(when $\gamma \rightarrow \infty$, $\text{regret}(T) \rightarrow H^2 \Rightarrow$ policy iteration)

Lemma: (Performance difference Lemma) \forall policies π, π' :

$$V_1^{\pi}(s_1) - V_1^{\pi'}(s_1) = \mathbb{E}_{\pi} \left[\sum_{h=1}^H \langle Q_h^{\pi'}(s_h, \cdot), \pi_h(\cdot | s_h) - \pi'_h(\cdot | s_h) \rangle | s_1 \right]$$

Proof: Let $\pi^{(0)} = \{\pi_1', \dots, \pi_H'\}$, $\pi^{(1)} = \{\pi_1, \pi_2', \dots, \pi_H'\}$, \dots , $\pi^{(H)} = \{\pi_1, \dots, \pi_H\}$

($\pi^{(t)}$: follows π for the first t steps and π' for the remaining $H-t$ steps)

Then: $V_1^{\pi}(s_1) - V_1^{\pi'}(s_1) = V_1^{\pi^{(H)}}(s_1) - V_1^{\pi^{(0)}}(s_1) = \sum_{h=1}^H V_1^{\pi^{(h)}}(s_1) - V_1^{\pi^{(h-1)}}(s_1)$

But: $V_1^{\pi^{(h)}}(s_1) - V_1^{\pi^{(h-1)}}(s_1) = \sum_{s_h} d_h^{\pi}(s_h) \left(\sum_a Q_h^{\pi'}(s_h, a) \pi_h(a | s_h) - V_h^{\pi'}(s_h) \right)$

$$= \mathbb{E}_{\pi} \left[\langle Q_h^{\pi'}(s_h, \cdot), \pi_h(\cdot | s_h) - \pi'_h(\cdot | s_h) \rangle \right]$$

$$\pi_h^{tH}(a|s) = \pi_h^t(a|s) e^{y Q_h^{\pi^t}(s,a)} / Z_{tH,s}$$

$$\Rightarrow Q_h^{\pi^t}(s,a) = \frac{1}{y} \log \frac{\pi_h^{tH}(a|s)}{\pi_h^t(a|s)} + \frac{1}{y} \log Z_{tH,s}$$

$$\Rightarrow \langle Q_h^{\pi^t}(s, \cdot), \pi_h^{tH}(\cdot|s) \rangle = \frac{1}{y} \text{KL} [\pi_h^{tH}(\cdot|s) \| \pi_h^t(\cdot|s)] + \frac{1}{y} \log Z_{tH,s}$$

$$\begin{aligned} - \langle Q_h^{\pi^t}(s, \cdot), \pi_h^*(\cdot|s) \rangle &= \frac{1}{y} \left[\text{KL} [\pi^*(\cdot|s) \| \pi_h^t(\cdot|s)] - \text{KL} [\pi^*(\cdot|s) \| \pi_h^{tH}(\cdot|s)] \right] \\ &\quad + \frac{1}{y} \log Z_{tH,s} \end{aligned}$$

$$\begin{aligned} \langle Q_h^{\pi^t}(s, \cdot), \pi_h^{tH}(\cdot|s) - \pi_h^*(\cdot|s) \rangle &= \frac{1}{y} \text{KL}_S (\pi_h^{tH}, \pi_h^t) \\ &\quad - \frac{1}{y} \text{KL}_S (\pi_h^*, \pi_h^t) + \frac{1}{y} \text{KL}_S (\pi_h^*, \pi_h^{tH}) \\ &\geq \frac{1}{y} \text{KL}_S (\pi_h^*, \pi_h^{tH}) - \frac{1}{y} \text{KL}_S (\pi_h^*, \pi_h^t) \end{aligned}$$

(replace π^* by π^t)

$$\langle Q_h^{\pi^t}(s, \cdot), \pi_h^{tH}(\cdot|s) - \pi_h^t(\cdot|s) \rangle \geq 0$$

$$\Rightarrow \langle Q_h^{\pi^t}(s, \cdot), \pi_h^t(\cdot | s) - \pi_h^* (\cdot | s) \rangle \cancel{\geq} \langle Q_h^{\pi^t}(s, \cdot), \pi_h^t(\cdot | s) - \pi_h^{t+1}(\cdot | s) \rangle$$

$$+ \frac{1}{\gamma} KL_S(\pi_h^*, \pi_h^{t+1}) - \frac{1}{\gamma} KL_S(\pi_h^*, \pi_h^t)$$

We have:

$$V_1^{\pi^*}(s_1) - V_1^{\pi^t}(s_1) = \mathbb{E}_{\pi^*} \left[\sum_{h=1}^H \langle Q_h^{\pi^t}(\cdot | s_h), \pi_h^* (\cdot | s_h) - \pi_h^t (\cdot | s_h) \rangle \right]$$

$$\leq \frac{1}{\gamma} \sum_{h=1}^H \left[KL_{S_h}(\pi_h^*, \pi_h^t) - KL_{S_h}(\pi_h^*, \pi_h^{t+1}) \right]$$

$$+ \mathbb{E}_{\pi^*} \left[\sum_{h=1}^H \langle Q_h^{\pi^t}(s_h, \cdot), \pi_h^{t+1}(\cdot | s_h) - \pi_h^t(\cdot | s_h) \rangle \right]$$

Note that:

$$V_h^{\pi^{t+1}}(s_h) - V_h^{\pi^t}(s_h) = \mathbb{E}_{\pi^{t+1}} \left[\sum_{h'=h}^H \langle Q_{h'}^{\pi^t}(\cdot | s_{h'}), \pi_{h'}^{t+1}(\cdot | s_{h'}) - \pi_{h'}^t(\cdot | s_{h'}) \rangle \right]$$

$$\geq \langle Q_h^{\pi^t}(\cdot | s_h), \pi_h^{t+1}(\cdot | s_h) - \pi_h^t(\cdot | s_h) \rangle$$

$$V_1^{\pi^*}(\xi_1) - V_1^{\pi^t}(\xi_1) \leq \frac{1}{\eta} \mathbb{E}_{\pi^*} \sum_{h=1}^H \left[K_{\xi_h}(\pi_h^*, \pi_h^t) - K_{\xi_h}(\pi_h^*, \pi_h^{t+1}) \right] \\ + \mathbb{E}_{\pi^*} \left[\sum_{h=1}^{H+1} V_h^{\pi^{t+1}}(\xi_h) - V_h^{\pi^t}(\xi_h) \right]$$

$$\sum_{t=1}^T \text{RHS} \leq \frac{1}{\eta} H \log A + H(H)$$

Summary (unknown r and p)

reference	Algorithm	Setting	regret
[Azar et al. '12]	VI-UCB	stochastic reward	$O(\sqrt{\text{poly}(H) SAT})$
[Efroni et al. '20]	NPG	stochastic reward	$O(\sqrt{\text{poly}(H) S^2 AT})$
[Jin et al. '19]	Upper Occupancy Band	adversarial reward	$O(\sqrt{\text{poly}(H) S^2 AT})$
[Efroni et al. '20]	NPG	adversary reward	$O(T^{2/3})$