

Data: $(s_1, a_1, \dots, s_H, a_H)$ $a_h \sim \pi_b$.

Goal: $J(\pi) = \mathbb{E}_\pi \left[\sum_{h=1}^H \gamma_h \right]$.

Estimator: Let $P_h := \frac{\pi(a_h | s_h)}{\pi_b(a_h | s_h)} \rightarrow P_{1:h} = \prod_{i=1}^h P_i$.

(step-wise) IS: $\sum_{h=1}^H P_{1:h} \gamma_h$.

WZS: $\sum_{h=1}^H \frac{1}{\Delta} P_{1:h}^i \gamma_h^i$

$\sum_{h=1}^H \frac{1}{\Delta} \sum P_{1:h}^i$

$$\underline{\underline{P_{1:H} \left(\sum_{h=1}^H \gamma_h \right)}}$$

$$\sum_{h=1}^H P_{1:h} \gamma_h$$

$$= P_1 \gamma_1 + P_{1:2} \gamma_2 + \dots + P_{1:H} \gamma_H$$

$$= P_1 (\gamma_1 + P_2 (\gamma_2 + P_3 (\gamma_3 + \dots + P_H \gamma_H)))$$

Define $v_0 := 0$.

$$v_{H-h+1} = P_h (r_h + v_{H-h})$$

$$= \frac{\pi(a_h | s_h)}{\pi_b(a_h | s_h)} \cdot \left(r_h + v_{H-h} \right)$$

$s_h, a_h \sim \pi_b$

$\rightarrow Q^\pi(s_h, a_h)$

$$= \sqrt{\pi}(s_{h+1}) \cdot \sqrt{\pi}(a_h)$$

$$\sqrt{\pi}(s_h)$$

$$r_h + \sqrt{\pi}(s_{h+1})$$

$$s_h, a_h \sim \pi_b$$

$$Q^\pi(s_n, a_n)$$

$$DR: V_{H-h+1} = \rho_h (V_n + V_{H-h} - \hat{Q}^\pi(s_n, a_n))$$

Policy Gradient $\Pi = \{ \pi_\theta : \theta \in \Theta \}$

$$\max_{\pi} J(\pi) \quad (\pi = \pi_\theta)$$

Idea: $\nabla J(\pi) = \nabla_{\theta} J(\pi_{\theta})$

Goal Given $\tau = (s_1, a_1, s_2, a_2, \dots, s_H, a_H)$

w/ $a_n \sim \pi$ estimate $\nabla J(\pi)$ from τ .

$$\begin{aligned} \nabla_{\theta} J(\pi_{\theta}) &= \mathbb{E}_{\tau \sim \pi} [f(\tau)] = \sum_{\tau} P^{\pi}(\tau) f(\tau) \\ &= \nabla \sum_{\tau} R(\tau) \underline{P^{\pi}(\tau)} \\ &= \sum R(\tau) \nabla P^{\pi}(\tau). \end{aligned}$$

SGD: $\theta \leftarrow \theta + \alpha \cdot f(\tau)$ ✓

$$= \sum_{\tau} R(\tau) \underbrace{P^{\pi}(\tau)} \underbrace{\nabla \log P^{\pi}(\tau)}$$

$$= \sum_{\tau} R(\tau) P^{\pi}(\tau) \nabla \log \left(\frac{d_0(s_1) \cdot \pi(a_1|s_1)}{P(s_2|s_1, a_1) \cdot \pi(a_2|s_2) \dots}$$

$$= \sum_{\tau} R(\tau) \underbrace{P^{\pi}(\tau)} \sum_{h=1}^H \left(\nabla \log P(s_h|s_{h-1}, a_{h-1}) + \nabla \log \pi(a_h|s_h) \right)$$

$$\Downarrow \mathbb{E}_{\tau \sim \pi} \left[R(\tau) \sum_{h=1}^H \nabla \log \pi(a_h|s_h) \right] \leftarrow \text{"RETURN GRAD"}$$

$$\Delta \frac{dJ(\pi_{\theta})}{d\theta} = \lim_{\Delta\theta \rightarrow 0} \frac{J(\pi_{\theta+\Delta\theta}) - J(\pi_{\theta})}{\Delta\theta}$$

$$\tau \sim \pi_{\theta}. \quad \mathbb{E}_{\tau} [R(\tau)] = J(\pi_{\theta})$$

$$\left(\frac{\prod_{h=1}^H \frac{\pi_{\theta+\Delta\theta}(a_h|s_h)}{\pi_{\theta}(a_h|s_h)} \right) R(\tau)$$

$$\left(\sum_{h=1}^H r_h \right) \quad \left(\sum_{h=1}^H \nabla \log \pi(a_h | s_h) \right)$$

discounted

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

$$\sum_{h=1}^H \nabla \log \pi(a_h | s_h) \quad \left(\sum_{h'=h}^H r_{h'} \right)$$

"Actor-critic"

$$\mathbb{E}_{(s,a) \sim \pi_{\theta}} \left[\nabla \log \pi(a|s) \cdot \hat{Q}^{\pi}(s,a) \right] \quad \text{FAE}$$

$\sum r_h$

π_{θ}

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

$$\forall f: S \rightarrow \mathbb{R}$$

$$\mathbb{E}_{a \sim \pi} \left[\nabla \log \pi(a|s) \right] f(s)$$

$$= \sum_a \cancel{\pi(a|s)} \cdot \frac{\nabla \pi(a|s)}{\cancel{\pi(a|s)}} = \vec{0}$$

PG: find $\hat{\pi}$, $\|\nabla J(\hat{\pi})\| = \varepsilon$.

$$\textcircled{I} = \mathbb{R}^d$$

$$J(\pi^*) - J(\hat{\pi})$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi^*}} \left[Q^{\hat{\pi}}(s, \pi^*) - Q^{\hat{\pi}}(s, \hat{\pi}) \right]$$

$$\left| \left(\nabla J(\hat{\pi}) \right)^T (\hat{\theta} - \theta) \right| \leq \varepsilon \quad \forall \theta$$

$$\leq \frac{1}{1-\gamma} \mathbb{E}_{d^{\pi^*}} \left[Q^{\hat{\pi}}(s, \pi_{Q^{\hat{\pi}}}) - Q^{\hat{\pi}}(s, \hat{\pi}) \right] \mathbb{E} \left[Q^{\hat{\pi}}(s, a) \right]$$

$$\leq \frac{1}{1-\gamma} \left\| \frac{d^{\pi^*}}{d^{\hat{\pi}}} \right\|_{\infty} \cdot \mathbb{E}_{d^{\hat{\pi}}} \left[Q^{\hat{\pi}}(s, \pi_{Q^{\hat{\pi}}}) - Q^{\hat{\pi}}(s, \hat{\pi}) \right]$$

policy parameters

II

$$\underline{\mathbb{H}_d^{\hat{\alpha}} \left[\mathbb{Q} \left(s, \frac{\hat{\alpha}}{\pi} \right) - \mathcal{V}_6 f_{\hat{\alpha}}(a|s) \right]}$$