

$$D = \{s, a, s'\} \quad s, a \sim \rho, \quad s' \sim P(\cdot | s, a)$$

Fitted Q-Iteration:

$$\text{Given } \mathcal{F} \subseteq S \times A \mapsto [0, \frac{1}{1-\gamma}]$$

Assumptions:

$$\textcircled{1} \text{ Bell-completeness: } \forall f \in \mathcal{F}: \quad \overset{\text{Bell-operator}}{\mathcal{T}} f \in \mathcal{F}$$

$$\textcircled{2} \max_{\pi} \max_{s, a} \frac{d^{\pi}(s, a)}{\rho(s, a)} \leq C < \infty$$

Min-Max Approach For PE

$$\textcircled{1} \text{ Coverage condition: } \max_{s, a} \frac{d^{\pi}(s, a)}{\rho(s, a)} \leq C < \infty$$

$$\text{Function classes: } \mathcal{F} \subseteq S \times A \mapsto [0, \frac{1}{1-\gamma}], \quad (Q^{\pi} \in \mathcal{F})$$

$$\mathcal{W} \subseteq S \times A \mapsto [-c, c] \quad \left(\frac{d^{\pi}}{\rho} \in \mathcal{W} \right)$$

$$\textcircled{2} \text{ Realizability: } Q^{\pi} \in \mathcal{F}, \quad \underline{\frac{d^{\pi}}{\rho}} \in \mathcal{W}$$

$$\textcircled{3} \text{ symmetric: } w \in \mathcal{W}, \Rightarrow -w \in \mathcal{W}$$

Algorithm:

$$\text{define: } \hat{l}(w, Q) := \frac{1}{N} \sum_{i=1}^N w(s_i, a_i) \left(Q(s_i, a_i) - r(s_i, a_i) - \gamma Q(s'_i, \pi(s'_i)) \right)$$

$$l(w, Q) = \mathbb{E}_{\text{sample}} \left[w(s, a) \left(Q(s, a) - r(s, a) - \gamma \mathbb{E}_{s' \sim P(s, \cdot)} Q(s', \pi(s')) \right) \right]$$

Min-Max PE:

$$\hat{Q} = \underset{Q \in \mathcal{F}}{\text{argmin}} \max_{w \in \mathcal{W}} \hat{l}(w, Q)$$

$$\hat{V}^{\pi} := \hat{Q}(s_0, \pi(s_0)) \quad (s_0 \text{ is fixed})$$

$$\text{Goal: } \left| \hat{V}^\pi - V^\pi \right| \leq \frac{1}{\sqrt{n}} \cdot \ln(|W||F|/\delta) \cdot \text{poly}\left(\frac{1}{1-\gamma}, C\right)$$

Analysis:

① Uniform convergence: (Hoeffding + Union Bound)

$\forall w \in W, Q \in \mathcal{F}$:

$$\left| \hat{l}(w, Q) - l(w, Q) \right| \leq \underbrace{\sqrt{\frac{\ln(|W||F|/\delta)}{N}}}_{\delta} \cdot \text{poly}\left(\frac{1}{1-\gamma}, C\right)$$

w.h.p.

② min-max

$$\min_Q \left[\max_w \hat{l}(Q, w) \right]$$

loss for Q

$$\forall Q: \left| \max_w \hat{l}(Q, w) - \max_w l(Q, w) \right| \leq \max_w \left| \hat{l}(Q, w) - l(Q, w) \right| \leq \delta$$

Recall: $\hat{Q} = \arg \min_Q \max_w \hat{l}(Q, w)$ (*)

$$\max_w l(w, Q^\pi) = 0 \quad (*, *)$$

$$\begin{aligned} \max_w l(w, \hat{Q}) &\leq \max_w \hat{l}(w, \hat{Q}) + \delta \\ &\leq \max_w \hat{l}(w, Q^\pi) + \delta \quad (*) \\ &= \max_w l(w, Q^\pi) + \delta \end{aligned}$$

$$= 2\delta \quad (**)$$

$$(3) \max_W \ell(v, \hat{Q}) \leq 2\delta \Rightarrow |\hat{V}^\pi - V^\pi|$$

$$\max_{W \in \mathcal{W}} E_{\text{samp}} \left(W(s_a) (\hat{Q}(s_a) - r(s_a) - \gamma E_{S'|S_a} \hat{Q}(s', \pi(s')) \right) \leq 2\delta \quad (a)$$

By symmetric condition on W :

$$\max_{W \in \mathcal{W}} E_{\text{samp}} \left[W(s_a) (-\hat{Q}(s_a) + r(s_a) + \gamma E_{S'|S_a} \hat{Q}(s', \pi(s')) \right] \leq 2\delta \quad (b)$$

By $W^\pi := \frac{d^\pi}{P} \in \mathcal{W}$, & (a)+(b)

$$\left| E_{\text{samp}} \left[W^\pi(s_a) (\hat{Q}(s_a) - r(s_a) - \gamma E_{S'|S_a} \hat{Q}(s', \pi(s')) \right] \right| \leq 2\delta$$

$$\min_Q \sum_{i=1}^n \left(Q(s_a) - r(s_a) - \gamma E_{S'|S_a} Q(s', \pi(s')) \right)^2 / n \leq \checkmark$$

is not unbiased wrt:

$$E_{\text{samp}} \left[\left(Q(s_a) - r(s_a) - \gamma E_{S'|S_a} Q(s', \pi(s')) \right)^2 \right] \checkmark$$

$$\left| E_{S_a, d^\pi} (\hat{Q}(s_a) - r(s_a) - \gamma E_{S'|S_a} \hat{Q}(s', \pi(s')) \right| \leq 2\delta \quad (**)$$

$$d^\pi(s_a) = (1-\gamma) \delta(s_a) \mathbb{1}(a = \pi(s_a))$$

$$+ \gamma E_{S_a \sim d^\pi} P(S|S_a) \mathbb{1}(a = \pi(S)) \quad (*)$$

$$E_{\pi} \hat{Q}(s_a) - \gamma E E \hat{Q}(s', \pi(s'))$$

$$\begin{aligned}
 \text{Sand''} & \quad \text{Sand'' Simpson} \\
 &= (1-\delta) \widehat{\Theta}(s_0, \pi(s_0)) \quad (*) \\
 &= (1-\delta) \widehat{V}^\pi
 \end{aligned}$$

Together w/ (**)

$$\Rightarrow \left| (1-\delta) \widehat{V}^\pi - (1-\delta) V^\pi \right| \leq \delta$$

$$\Rightarrow E_{\text{Sand''}} \Gamma(s, \alpha)$$

□

Summary:

$$\left| \widehat{V}^\pi - V^\pi \right| \leq \frac{1}{\sqrt{n}} \cdot \sqrt{\ln(M) / \delta} \cdot \text{poly}\left(\frac{1}{1-\delta}, c\right)$$