

Tabular Method

Data  $\{(s, a, r, s')\}$ .

→ Assume: each  $(s, a)$  appears  $n$  times.

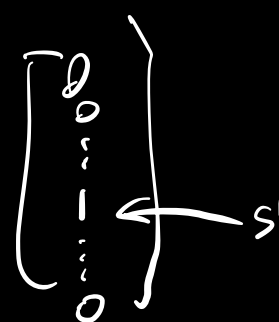
$$\underline{r} \sim \underline{R}(\cdot | s, a), \quad \underline{s}' \sim P(\cdot | s, a)$$

Let  $D_{s,a}$  be the  $n$   $(r, s')$  pairs from  $(s, a)$ .

Alg: "certainty-equivalence", estimate  $\hat{M}$ .

$$\hat{R}(s, a) = \frac{1}{n} \sum_{(r, s') \in D_{s,a}} r$$

$$\hat{P}(\cdot | s, a) = \frac{1}{n} \sum_{(r, s') \in D_{s,a}} e_{s'}$$



Output:  $\pi_{\hat{M}}^*$  and  $V_M^{\pi_{\hat{M}}^*}$ .

$$\text{Analysis: } \|V_M^* - V_M^{\pi_{\hat{M}}^*}\|_{\infty} \leq O\left(\frac{1}{\sqrt{n}}\right)$$

w.p.  $1 - \delta$ .

Define:  $\epsilon_R := \max_{s,a} \left| \hat{R}(s,a) - \frac{R(s,a)}{\Delta} \right|$

$\epsilon_P := \max_{s,a} \| \hat{P}(\cdot | s, a) - P(\cdot | s, a) \|_1$

# Simulation Lemma

$$\forall \pi: S \rightarrow A \quad (S \rightarrow \Delta(A))$$

$$\rightarrow \|V_M^\pi - V_{\hat{M}}^\pi\|_\infty \leq \frac{\epsilon_R}{1-\gamma} + \frac{\delta \epsilon_p V_{\max}}{2(1-\delta)} \quad \frac{R_{\max}}{1-\delta}$$

- ①  $\|V_M^* - V_M^{\pi_M^*}\|_\infty \leq \boxed{\times 2}$  ✓ ✓
- ② Prove Simu Lemma.
- ③ Bound  $\epsilon_R, \epsilon_p$  as a fn of  $n$ .

$\forall s$

$$V_M^*(s) - V_M^{\pi_M^*}(s) = \underbrace{V_M^{\pi_M^*}(s) - V_{\hat{M}}^{\pi_M^*}(s)}_{\text{"decision loss"}} + \underbrace{V_{\hat{M}}^{\pi_M^*}(s) - V_M^{\pi_M^*}(s)}_{\leq V_{\hat{M}}^{\pi_M^*}} \leq 2 \cdot \max_{\pi: S \rightarrow A} \underbrace{\|V_M^\pi - V_{\hat{M}}^\pi\|_\infty}_{\text{"evaluation error"}}$$

$\forall f \in \mathbb{R}^S \quad \forall s_0 \quad \forall \pi: S \rightarrow \Delta(A)$

$$\frac{f(s_0) - V_M^\pi(s_0)}{1-\gamma} = \frac{f(s) - \gamma f(s')}{1-\gamma}$$

$s \sim d^{\pi, s_0}$   
 $a \sim a(\cdot|s)$   
 $v \sim R(\cdot|sa)$

$$s' \sim p(\cdot | s, a)$$

$$= \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^{\pi, s_0}} [f(s) - (R(s, \pi) + \gamma \cdot \mathbb{E}_{s' \sim p(s, \pi)} [f(s')])]$$

$$f - T_M^\pi f$$

$$(T_M^\pi f)(s)$$

Bellman error/residual.

$$(f - T_M^\pi f)$$

$$\|f - V_M^\pi\|_\infty \leq \frac{1}{1-\gamma} \|f - T_M^\pi f\|_\infty$$

$$\|f - V_M^\pi\|_\infty \leq \|f - T_M^\pi f\|_\infty + \|T_M^\pi f - T_M^\pi V_M^\pi\|_\infty \leq \gamma \cdot \|f - V_M^\pi\|_\infty$$

$$\frac{1}{1-\gamma} \mathbb{E} \left[ \underbrace{s \sim d^{\pi, s_0}}_{\substack{a \sim a(\cdot|s), \\ r \sim R(\cdot|s,a), \\ s' \sim p(\cdot|s,a)}} \left[ \underbrace{f(s)} - \gamma - \gamma \underbrace{f(s')} \right] \right]$$

$$= V_M^{\pi}(s_0)$$

$$\mathbb{E}_{s \sim d_1}^{\pi, s_0} \left[ f(s) - \gamma f(s') \right]$$

$$+ \gamma \mathbb{E}_{s \sim d_2}^{\pi, s_0} \left[ f(s) - \gamma f(s') \right]$$

$$+ \gamma^2 \mathbb{E}_{s \sim d_3}^{\pi, s_0} \left[ f(s) - \gamma f(s') \right]$$

$$+ \dots$$

$$= f(s_0)$$

$$\mathbb{E}_{\phi} [f] = \mathbb{E}_{\theta} [f] \quad \forall f$$

$$\left| V_M^\pi(s_0) - V_M^\pi(s_0) \right| = \left| \frac{1}{1-\gamma} \mathbb{E}_{(s,a) \sim d^{\pi,s_0}} \left[ \frac{V_M^\pi(s) - \gamma V_M^\pi(s')}{r + \gamma V_M^\pi(s')} \right] \right|$$

$(s,a) \sim d^{\pi,s_0}$   
 $r \sim R(\cdot|s,a)$   
 $s' \sim P(\cdot|s,a)$

$$\leq \frac{1}{1-\gamma} \max_s \left| V_M^\pi(s) - R(s,\pi) - \gamma \mathbb{E}_{s' \sim P(\cdot|s,\pi)} [V_M^\pi(s')] \right|$$

$$\leq \frac{1}{1-\gamma} \max_s \left| \hat{R}(s,\pi) + \gamma \mathbb{E}_{s' \sim \hat{P}(\cdot|s,a)} [V_M^\pi(s')] - R(s,\pi) - \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [V_M^\pi(s')] \right|$$

$\epsilon_R \triangleright$

$$= \left| \gamma \langle \hat{P}(\cdot|s,a), V_M^\pi \rangle - \gamma \langle P(\cdot|s,a), V_M^\pi \rangle \right|$$

$\gamma \epsilon_P \cdot V_{max}$

$[0, V_{max}]$

$$= \gamma \left| \langle \hat{P}(\cdot|s,a) - P(\cdot|s,a), V_M^\pi \rangle \right|$$

$$\leq \gamma \underbrace{\| \hat{P}(\cdot|s,a) - P(\cdot|s,a) \|}_{\leq \epsilon_P} \cdot \| V_M^\pi \|_\infty$$

$$\downarrow$$

$$V_{max} = \frac{R_{max}}{1-\gamma}$$

Hölder's ineq:

$$u, v \in \mathbb{R}^n$$

$$|u^T v| \leq \|u\| \cdot \|v\|_*$$

for any  $\|\cdot\|, \|\cdot\|_*$ .

Special case:

$$= \frac{1}{1-\gamma} (\epsilon_R + \delta \epsilon_P \frac{V_{max}}{2\gamma})$$

$$\|\cdot\|_p \xleftrightarrow{\text{dual}} \|\cdot\|_q$$

for  $1/p + 1/q = 1$ .

$$\sqrt{u^T A u} = \|u\|_A$$

$$\|\cdot\|_A, \|\cdot\|_{A^{-1}}$$

$$u^T v = |u^T A^{1/2} A^{-1/2} v|$$

$$[0, V_{max}] \leq \|u^T A^{1/2}\|_2 \cdot \|\cdot\|_2$$

$$|\langle \hat{P} - P, V \rangle| \leq |\langle \hat{P} - P, V - \frac{V_{max}}{2} \mathbb{1} \rangle|$$

$$|\hat{R}(s,a) - R(s,a)|$$

↑

$$\frac{1}{n} \sum_{v \in D_{s,a}} v$$

$$R_{max} \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}$$

- w.p.  $\geq 1-\delta$

$$V - \frac{V_{max}}{2} \mathbb{1}$$

$$[-\frac{V_{max}}{2}, \frac{V_{max}}{2}]$$

$$\|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1$$

$$= \sum_{\tilde{s}} |\hat{P}(\tilde{s} | s, a) - P(\tilde{s} | s, a)|.$$

$$\frac{1}{n} \sum_{s' \in \mathcal{S}, a} \mathbb{I}[s' = \tilde{s}].$$

$$\sqrt{|\mathcal{S}|} \sqrt{\frac{1}{2a} \ln \frac{2|\mathcal{S}|}{\delta}}$$

$$\|\hat{P}(\cdot | \Delta(s)) - P(\cdot | \Delta(s))\|_1.$$

$$= \max_{u \in \{-1, 1\}^{\mathcal{S}}} u^T (\hat{P} - P).$$



$$\|x\|_x = \max_{\|y\| \leq 1} y^T x$$

$$u^T \hat{P} = u^T \frac{1}{n} \sum_{s' \in \mathcal{S}} e_{s'}$$

$$= \frac{1}{n} \sum_{s' \in \mathcal{D}_{s,n}} u^T e_{s'}$$

$$= \frac{1}{n} \sum_{s'} u(s')$$

$$\mathbb{E}_{s' \sim p} [u(s')] = \langle P, u \rangle$$

$$u^T \hat{P} - u^T P \stackrel{[-1,1]}{\leq}$$

Fix any  $u$ ,  $|u^T \hat{P} - u^T P| \leq 2 \cdot \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}$

$$\|\hat{P} - P\|_1 = \max_{u \in \{-1,1\}^S} |u^T \hat{P} - u^T P|$$

$$\leq 2 \cdot \sqrt{\frac{1}{2n} \ln \frac{2 \cdot 2^S}{\delta}}$$

$$\|V_M^* - V_M^{\pi_M^*}\|_\infty$$

$$\pi_M^* = \pi_{Q_M^*}$$

$\forall f: \mathbb{R}^{S \times A}$

$$\|V_M^* - V_M^{\pi_f}\|_\infty$$

$$\leq \frac{2 \|f - Q_M^*\|_\infty}{1 - \gamma}$$

$$\leq \frac{2 \cdot \|Q_M^* - Q_M\|_\infty}{1 - \gamma}$$

$$1 - \gamma$$



$$\begin{aligned}
& \| Q_{\hat{M}}^* - Q_M^* \|_{\infty} \leq \frac{1}{1-\gamma} V_{\max} \sqrt{\frac{1}{n} \ln \frac{2|S_A|}{\delta}} \\
& = \| Q_{\hat{M}}^* - T_{\hat{M}} Q_M^* + T_{\hat{M}} Q_M^* - Q_M^* \|_{\infty} \\
& \leq \| T_{\hat{M}} Q_{\hat{M}}^* - T_{\hat{M}} Q_M^* \|_{\infty} + \| T_{\hat{M}} Q_M^* - T_M Q_M^* \|_{\infty} \\
& \leq \gamma \cdot \| Q_{\hat{M}}^* - Q_M^* \|_{\infty} \quad (\text{II})
\end{aligned}$$

$$\boxed{R(s,a) + \gamma \langle P(\cdot|s,a), V_M^* \rangle} \quad \left[ \text{max}_{a'} \langle P(\cdot|s,a'), V_M^* \rangle \right]$$

$$(\text{II})(s,a) = \boxed{\hat{R}(s,a) + \gamma \langle \hat{P}(\cdot|s,a), V_M^* \rangle}$$

$$- R(s,a) - \gamma \langle P(\cdot|s,a), V_M^* \rangle$$

$$\frac{1}{n} \sum_{(r,s') \in D_{s,a}} r + \gamma \langle \frac{1}{n} \sum_{\Delta} e_{s'}, V_M^* \rangle$$

$$= \frac{1}{n} \sum_{(r,s') \in D_{s,a}} \left( r + \gamma \langle e_{s'}, V_M^* \rangle \right)$$

$$= \frac{1}{n} \sum_{(r,s') \in D_{s,a}} \left( r + \gamma V_M^*(s') \right)$$

$$\left| \frac{1}{n} \sum_{\Delta} (r + \gamma V_M^*(s')) - Q_M^*(s, a) \right|$$

$\mathbb{E}_{r, s'}[\cdot] = Q_M^*(s, a)$

$$R_{\max} + \gamma \cdot \frac{R_{\max}}{1-\gamma} = \frac{R_{\max}}{1-\gamma} = V_{\max}$$

$$\leq V_{\max} \cdot \sqrt{\frac{1}{2n} \ln \frac{2|S \times A|}{\delta}}$$