

MDP (tabular analysis) $M = (S, A, P, R, \gamma)$

Protocol: for each (s, a) pair.

n iid samples of $s' \sim P(\cdot | s, a)$.
data $D_{s,a}$ and $r \sim R(\cdot | s, a)$.

$$\frac{\Delta}{n} \left\| V_M^* - \underset{\Delta}{V_M} \right\|_{\infty} \leq \varepsilon.$$

Alg: $\hat{R}(s, a) = \frac{1}{n} \sum_{r_i \in D_{s,a}} r_i$

$$\hat{P}(s' | s, a) = \frac{1}{n} \sum_{s'_i \in D_{s,a}} \mathbb{I}[s'_i = s']$$

const. \uparrow
i.i.d. \uparrow

$$\hat{P}(\cdot | s, a) = \frac{1}{n} \mathbf{e}_{s'}$$

$$\hat{M} = (S, A, \hat{P}, \hat{R}, \gamma)$$



Concentration: $\hat{M} \approx M$

$$\underbrace{(w.p. 1 - \delta)}_{\forall s, a} \left\{ \left| \hat{R}(s, a) - R(s, a) \right| \leq \varepsilon_R \right.$$

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

② "error propagation analysis".

$$\| V_M^* - \underline{V}_M^{\pi_M^*} \|_\infty \leq f_n(\epsilon_r, \epsilon_p)$$

$$\forall (s, a) \in S \times A, \quad \text{w.p.} \geq 1 - \delta, \quad | \hat{R}(s, a) - R(s, a) | \leq R_{\max} \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}.$$

$$\forall (s, a, s'), \quad \text{w.p.} \geq 1 - \delta, \quad | \hat{P}(s' | s, a) - P(s' | s, a) | \leq \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}.$$

Given total failure budget δ .

$\left\{ \begin{array}{l} \text{all reward estim: } \frac{\delta}{2} \rightarrow \text{evenly over } S \times A. \\ \text{all transition: } \frac{\delta}{2} \rightarrow \text{evenly over } S \times A \times S. \end{array} \right.$

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

$$\forall (s, a).$$

$$\left| \hat{R}(s, a) - R(s, a) \right| \leq R_{\max} \sqrt{\frac{1}{2n} \ln \frac{4|S \times A|}{\delta}}$$

$$\| \hat{P}(\cdot | s, a) - P(\cdot | s, a) \|_1 \leq |S| \sqrt{\frac{1}{2n} \ln \frac{4|S \times A \times S|}{\delta}}$$

"Simulation Lemma" $\forall \pi: S \rightarrow A$. $R_{\max}/(1-\gamma)$

$$\Delta. \quad \underbrace{\| V_M^\pi - \hat{V}_M^\pi \|_\infty} \leq \frac{\epsilon_R}{1-\gamma} + \frac{\gamma \epsilon_P \overbrace{V_{\max}}^{\| \cdot \|}}{1-\gamma}$$

$$\| V_M^{\pi^*} - \hat{V}_M^{\pi^*} \|_\infty \leq \| V_M^{\pi^*} - \hat{V}_M^{\pi^*} + \hat{V}_M^{\pi^*} - V_M^{\pi^*} \|_\infty$$

$$\leq 2 \cdot \max_{\pi: S \rightarrow A} \| V_M^\pi - \hat{V}_M^\pi \|_\infty$$

Proof of sim. Lemma: $\forall s$.

$$V^\pi(s) = Q^\pi(s, \pi)$$

$$\left| V_M^\pi(s) - \hat{V}_M^\pi(s) \right|$$

$$= \left| R(s, \pi) + \gamma \langle P(\cdot | s, \pi), V_M^\pi(\cdot) \rangle - \left[R(s, \pi) + \gamma \langle \hat{P}(\cdot | s, \pi), \hat{V}_M^\pi(\cdot) \rangle \right] \right|$$

$$\begin{aligned}
& - \hat{R}(s, \pi) - \gamma \langle \hat{P}(\cdot | s, \pi), V_{\hat{M}}^{\pi}(\cdot) \rangle \\
& \leq \epsilon_R + \gamma \left| \langle P(\cdot | s, \pi), V_M^{\pi} \rangle - \langle \hat{P}(\cdot | s, \pi), V_M^{\pi} \rangle \right. \\
& \quad \left. - \langle P(\cdot | s, \pi), V_{\hat{M}}^{\pi} \rangle + \langle P(\cdot | s, \pi), V_{\hat{M}}^{\pi} \rangle \right| \\
& \leq \epsilon_R + \gamma \left| \langle \hat{P}(\cdot | s, \pi) - P(\cdot | s, \pi), V_{\hat{M}}^{\pi} \rangle \right| \leq O(\epsilon_P).
\end{aligned}$$

$$\begin{aligned}
& + \gamma \left| \langle P(\cdot | s, \pi), V_M^{\pi} - V_{\hat{M}}^{\pi} \rangle \right| \\
& \leq \gamma \cdot \|V_M^{\pi} - V_{\hat{M}}^{\pi}\|_{\infty}
\end{aligned}$$

Hölder's ineq: for $\|\cdot\|$ and $\|\cdot\|_*$.

$$|\langle u, v \rangle| \leq \|u\| \cdot \|v\|_*$$

Example: $\|\cdot\|_p$ & $\|\cdot\|_q$ for $\frac{1}{p} + \frac{1}{q} = 1$.

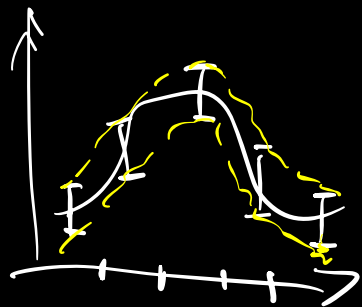
e.g. $p = q = 2$.

$p = 1, q = \infty$.

$$| \langle \hat{P}(\cdot | s, \pi) - P(\cdot | s, \pi), V_M^\pi \rangle |$$

$$\leq \| \hat{P}(\cdot | s, \pi) - P(\cdot | s, \pi) \|_1 \cdot \| V_M^\pi \|_\infty$$

$$\leq \epsilon_p \cdot \frac{R_{\max}}{1-\gamma} \rightarrow V_{\max}$$



$\hat{P} (= P(\cdot | s, \pi))$
 s

$\hat{P} (= \hat{P}(\cdot | s, \pi))$
 s

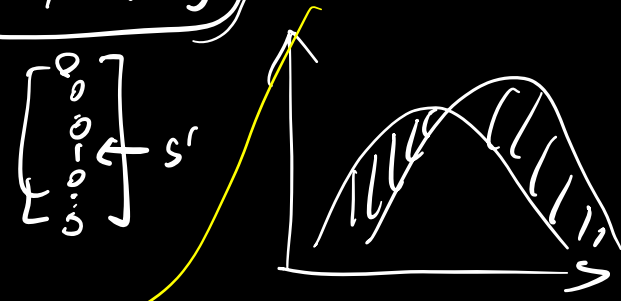
$$\| \phi - \hat{\phi} \|_1 = \max_{u \in \{-1, 1\}^S} u^T (\phi - \hat{\phi})$$

$$u^T \phi - u^T \hat{\phi}$$

$$= u^T \phi - u^T \frac{1}{n} \sum_{s' \in D_{s,n}} e_{s'}$$

$$= \mathbb{E}[u(s)] - \frac{1}{n} \sum_{s' \in D_{s,n}} u(s')$$

$$\stackrel{w.p. 1-\delta}{\leq} 2 \cdot \sqrt{\frac{1}{n} \ln \frac{2}{\delta}}$$



$$\sum_s | \phi(s) - \hat{\phi}(s) |$$

$$= \sum_s \text{sign}[\phi(s) > \hat{\phi}(s)] (\phi(s) - \hat{\phi}(s))$$

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

$$\text{w.p. } \geq 1-\delta, \|\hat{\beta} - \beta\|_1 \leq 2 \sqrt{\frac{1}{2n} \ln \frac{2 \cdot 2^s}{\delta}} \approx \sqrt{\frac{|S|}{n} \ln \frac{1}{\delta}}$$

$$2 \left(\frac{\sum_R \downarrow}{1-\delta} + \delta \frac{\sum_P \uparrow \sqrt{\max}}{1-\delta} \right).$$

Alt. bound: $\forall f \in \mathbb{R}^{S \times A}$

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

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

$$Q_M^*(s,a) - Q_M^*(s,a)$$

$$= Q_M^*(s,a) - R(s,a) - \gamma \langle \hat{P}(s,a), \underline{V_M^*} \rangle$$

$$= Q_M^*(s,a) - \hat{R}(s,a) - \gamma \langle \hat{P}(\cdot|s,a), \underline{V_M^*} \rangle$$

$$+ \gamma \langle \hat{P}(\cdot|s,a), \underline{V_M^*} \rangle - \gamma \langle \hat{P}(\cdot|s,a), \underline{V_M^*} \rangle$$

$$Q_M^*(s,a) - \frac{1}{n} \sum_{r_i \in D_{s,n}} r_i - \frac{1}{n} \langle \sum_{s_i \in D_{s,n}} e_{s_i}, \underline{V_M^*} \rangle$$

$$= Q_M^*(s,a) - \frac{1}{n} \sum_{(r_i, s_i) \in D_{s,n}} (r_i + \gamma V_M^*(s_i))$$

w.p.z. $(-\gamma)$ [0, V_max]

\leq

$$V_{\max} \sqrt{\frac{1}{2\alpha} \ln \frac{2}{\delta}}$$

$$\leq \frac{2 \max_{\pi} \|V_{\gamma}^{\pi} - V_{\gamma'}^{\pi}\|_{\infty}}{\Delta}$$

$$Q_{\gamma}^* - Q_{\gamma'}^*$$

$$\leq \frac{Q_{\gamma}^{\pi_{\gamma}^*} - Q_{\gamma'}^{\pi_{\gamma}^*}}{\Delta}$$

Lemma: $\forall f \in \mathbb{R}^S$ fix do.

$$\mathbb{E}_{s \sim d_0} [f(s_0)] - \mathbb{E}_{s \sim d_0} [V_M^\pi(s)] \quad \leftarrow \quad \text{↯}$$

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

$$\begin{aligned} \gamma &\sim R(\cdot | s, a) \\ s' &\sim P(\cdot | s, a) \end{aligned}$$

$$\begin{aligned} &\parallel \\ &T^\pi f. \\ &= \end{aligned}$$

Let $f = V_M^\pi$.

$$\mathbb{E}_{d_0} [V_M^\pi - V_M^\pi] = \frac{1}{1-\gamma} \mathbb{E}_{s \sim d^\pi, a \sim \pi} \left[V_M^\pi(s) - (T_M^\pi V_M^\pi)(s) \right]$$

$$\left(T_M^\pi V_M^\pi \right) (s)$$

$$\begin{aligned} &= R(s, \pi) + \gamma \langle P(\cdot | s, \pi), V_M^\pi \rangle \\ &- \hat{R}(s, \pi) - \gamma \langle \hat{P}(\cdot | s, \pi), V_M^\pi \rangle. \end{aligned}$$

$$(1-\gamma) \sum_{t=1}^{\infty} \gamma^{t-1} d_t^\pi$$

Proof: $\frac{1}{1-\gamma} \mathbb{E}_{s \sim d_1^\pi, a \sim \pi} [f(s) - r - \gamma f(s')]$

$= \sum_{t=1}^{\infty} \gamma^{t-1} \mathbb{E}_{s \sim d_t^\pi, a \sim \pi} [f(s) - r - \gamma f(s')]$

$\mathbb{E}_{s \sim d_1^\pi, a \sim \pi} [f(s) - r - \gamma f(s')]$
 $+ \gamma \mathbb{E}_{s \sim d_2^\pi, a \sim \pi} [f(s) - r - \gamma f(s')]$
 $+ \gamma^2 \mathbb{E}_{s \sim d_3^\pi, a \sim \pi} [f(s) - r - \gamma f(s')]$
 \vdots

$J_n(\pi) = \mathbb{E}_{s \sim d_0} [V_n^\pi]$

$\Phi = \mathcal{F} \Leftrightarrow \forall f \in \mathcal{F}, \mathbb{E}_\Phi[f] = \mathbb{E}_\mathcal{F}[f]$