

Tabular Analysis of model-based RL

Setting: for each $(s,a) \in S \times A$.

$D_{s,a}$ = a bag of n iid (r,s') , where $\left. \begin{array}{l} r \sim R(s,a) \\ s' \sim P(\cdot | s,a) \end{array} \right\} \begin{array}{l} r \in [0, R_{\max}] \\ \end{array}$

Algorithm: Build empirical MDP: $\forall s,a$.

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

$$\hat{P}(\cdot | s,a) = \frac{1}{n} \sum_{s' \in D_{s,a}} e_{s'} \quad \begin{array}{l} \text{empirical freq.} \\ \text{[0, 0, \dots, 0, 1, 0, \dots, 0]} \\ \text{\small s'-th coordinate} \end{array}$$

let $\hat{\pi}$ be the opt. policy in $\hat{M} = (S, A, \hat{P}, \hat{R}, \gamma)$

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

want to bound $\|V_M^* - V_M^{\hat{\pi}}\|_{\infty}$. certainty-equivalence

Analysis: $\left\{ \begin{array}{l} \text{Step 1: establish } \hat{M} \approx M \\ \text{Step 2: bound } V_M^* - V_M^{\hat{\pi}} \text{ using the error} \end{array} \right.$

How to measure the error of \hat{M} (w.r.t. M)?

$$\left[\begin{array}{l} \epsilon_R = \max_{s,a} |R(s,a) - \hat{R}(s,a)| \\ \epsilon_P = \max_{s,a} \|P(\cdot | s,a) - \hat{P}(\cdot | s,a)\|_1 \end{array} \right.$$

$\|u\|_1 = \sum_i |u_i|$

Analysis | Start w/ step 2:

$$\begin{aligned} \forall s, V_M^*(s) - V_M^{\frac{1}{\pi}}(s) &= V_M^{\pi_M^*}(s) - V_M^{\pi_{\hat{M}}^*}(s) \\ &\leq V_M^{\pi_M^*}(s) - V_M^{\pi_{\hat{M}}^*}(s) + V_M^{\pi_{\hat{M}}^*}(s) - V_M^{\pi_M^*}(s) \\ &\geq 0 \text{ b/c } \pi_{\hat{M}}^* \text{ is opt. in } \hat{M} \end{aligned}$$

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

Simulation Lemma $\forall \pi: S \rightarrow A$. $\leq \frac{R_{\max}}{1-\gamma}$.

$$\|V_M^{\pi} - V_{\hat{M}}^{\pi}\|_{\infty} \leq \frac{\epsilon_R + \gamma \epsilon_P V_{\max}/2}{1-\gamma}$$

Proof: $\forall s, |V_M^{\pi}(s) - V_{\hat{M}}^{\pi}(s)|$

$$= |R(s, \pi) + \gamma \langle P(\cdot | s, \pi), V_M^{\pi}(\cdot) \rangle$$

$$- \hat{R}(s, \pi) - \gamma \langle \hat{P}(\cdot | s, \pi), V_{\hat{M}}^{\pi}(\cdot) \rangle|$$

$$\leq \epsilon_R + \gamma | \langle P(s, \pi), V_M^{\pi} \rangle - \langle \hat{P}(s, \pi), V_M^{\pi} \rangle + \langle \hat{P}(s, \pi), V_M^{\pi} \rangle - \langle \hat{P}(s, \pi), V_{\hat{M}}^{\pi} \rangle |$$

$$\leq \epsilon_R + \gamma \left| \langle P(s, \pi) - \hat{P}(s, \pi), V_M^\pi \rangle \right| \quad (A)$$

$$+ \gamma \left| \langle \hat{P}(s, \pi), V_M^\pi - V_{\hat{M}}^\pi \rangle \right| \quad (B)$$

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

Slightly tighter: $\langle P(s, \pi), \vec{1} \rangle = 1$.

$$(A) = \left| \langle P(s, \pi) - \hat{P}(s, \pi), V_M^\pi - \frac{V_{\max}}{2} \vec{1} \rangle \right|$$

$$\leq \|P(s, \pi) - \hat{P}(s, \pi)\|_1 \cdot \left\| V_M^\pi - \frac{V_{\max}}{2} \vec{1} \right\|_\infty$$

$$\leq \epsilon_P \cdot \frac{V_{\max}}{2} \quad \left[-\frac{V_{\max}}{2}, \frac{V_{\max}}{2} \right]$$

Hölder's inequality:

$$\forall p, q, \text{ s.t. } \frac{1}{p} + \frac{1}{q} = 1$$

$$|\langle u, v \rangle| \leq \|u\|_p \cdot \|v\|_q$$

Special cases:

(1) $p = q = 2$.

(2) $p = 1, q = \infty$.

$$(B) \leq \|V_M^\pi - V_{\hat{M}}^\pi\|_\infty$$

Dual norm: $\|\cdot\|, \|\cdot\|_*$

$$\|x\|_* = \sup_{\|y\| \leq 1} y^T x$$

Put together: $\|V_M^\pi - V_{\hat{M}}^\pi\|_\infty \leq \epsilon_R + \gamma \cdot \epsilon_P \cdot \frac{V_{\max}}{2}$

$$+ \gamma \cdot \|V_M^\pi - V_{\hat{M}}^\pi\|_\infty$$

$$\Rightarrow \|V_M^\pi - V_{\hat{M}}^\pi\|_\infty \leq \frac{\epsilon_R + \gamma \epsilon_P V_{\max} / 2}{1 - \gamma}$$

$$\|V_M^* - V_M^\pi\|_\infty \leq 2 \max_{\pi} \|V_M^\pi - V_{\hat{M}}^\pi\| \leq 2 \cdot (\cdot)$$

Step 1: $M \approx \hat{M}$ (i.e. upper bound ϵ_R & ϵ_p).

Fix any s, a . w.p. $\geq 1 - \delta$,

assume stoch. rewards bounded in $[0, R_{\max}]$.

$$|\hat{R}(s, a) - R(s, a)| \leq R_{\max} \cdot \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}$$

Fix any s, a ,

Naive analysis:

$$\begin{aligned} \|\hat{P}(\cdot | s, a) - P(\cdot | s, a)\|_1 &= \sum_{s'} |\hat{P}(s' | s, a) - P(s' | s, a)| \\ &\leq |S| \cdot \max_{s'} |\hat{P}(s' | s, a) - P(s' | s, a)| \end{aligned}$$

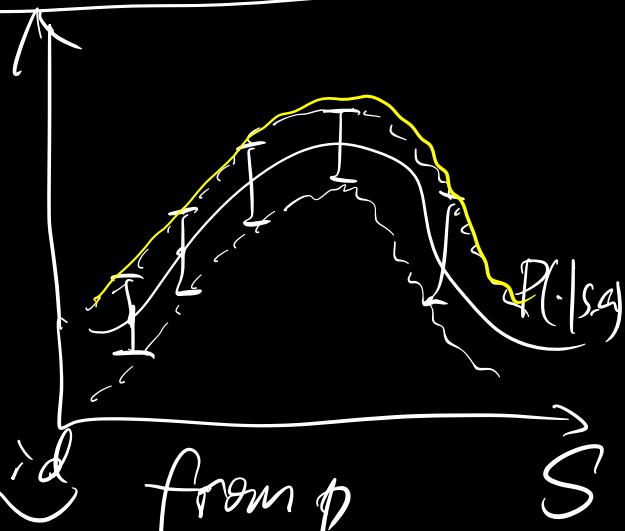
$$\frac{1}{n} \sum_{\tilde{s} \in \mathcal{D}_{s, a}} \mathbb{I}[\tilde{s} = s'] \leq |S| \cdot O\left(\frac{1}{\sqrt{n}}\right)$$

Lemma: Let p be a multinomial dist. over S .

Let \hat{p} be emp. dist.

estimated from n samples i.i.d from p .

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



Proof: For any vector $v \in \mathbb{R}^s$.

$$\|v\|_1 = \max_{u \in \{-1, 1\}^s} u^T v.$$

$$\leq \underbrace{\|u\|_\infty}_{=1} \cdot \|v\|_1 = \|v\|_1$$

Proof of LHS=RHS:

$$\|v\|_1 \leq \text{RHS.}$$

$$\text{RHS} \leq \|v\|_1 \text{ (Hölder).}$$

Therefore: $\|p - \hat{p}\|_1$

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

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

Since $\hat{p} = \frac{1}{n} \sum_{s \in D} e_s$, where D is the iid samples drawn $\sim \varphi$.

$$u^T p - u^T \hat{p} = u^T p - \frac{1}{n} \sum_{s \in D} u^T e_s$$

Now, $u^T e_s$ is iid (scalar valued) r.v. with mean $u^T p$.

$$\text{Also, } |u^T e_s| \leq \|u\|_\infty \cdot \|e_s\|_1 = 1.$$

$\hookrightarrow e_s \in \{-1, 1\}$.

\Rightarrow Fixing u , Hoeffding's: v.p. $\geq 1 - \delta$.

$$|u^T \phi - \frac{1}{n} \sum_{s \in \mathcal{D}} u^T e_s| \leq 2 \cdot \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}$$

By union bound: v.p. $\geq 1 - \delta$,

$$\|\phi - \hat{\phi}\|_1 = \max_{u \in \{-1, 1\}^S} |u^T \phi - u^T \hat{\phi}| \leq 2 \cdot \sqrt{\frac{1}{2n} \ln \frac{2 \cdot 2^{|S|}}{\delta}} \approx \sqrt{\frac{|S|}{n}}$$

Finally, union bound over $|S \times A|$ reward estimations
& $|S \times A|$ transition estimations.

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

$$\epsilon_\phi \leq 2 \cdot \sqrt{\frac{1}{2n} \cdot \ln \frac{2 \cdot 2 \cdot |S \times A| \cdot 2^{|S|}}{\delta}}$$

Eventually: v.p. $\geq 1 - \delta$.

$$\|V_M^* - V_M^{\hat{\pi}}\|_\infty \leq \boxed{\approx 0} \left(\frac{\sqrt{S} \cdot V_{\max}}{\sqrt{n} (1 - \gamma)} \right)$$

ignoring logarithms

Alternative Analysis.

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

$$\|V_M^* - V_M^{\hat{\pi}}\|_{\infty} \leq \tilde{O}\left(\frac{V_{\max}}{\sqrt{n}(1-\gamma)^2}\right).$$

In the previous analysis: $\forall \pi, V_M^{\pi} \approx V_M^{\hat{\pi}}$.

We will use a different way to establish the near-opt of $\hat{\pi}$.

Recall: $\forall f \in \mathbb{R}^{S \times A}, \|V_M^* - V_M^{\pi_f}\|_{\infty} \leq \frac{2 \|f - Q_M^*\|_{\infty}}{1-\gamma}$.

If we want to let $\hat{\pi} = \pi_f$, what is f ?

Answer: $f = Q_M^*$. ($\hat{\pi} = \pi_{Q_M^*}$).

$$\Rightarrow \|V_M^* - V_M^{\hat{\pi}}\|_{\infty} \leq \frac{2}{1-\gamma} \cdot \| \underbrace{Q_M^* - Q_M^*}_{\Delta} \|_{\infty}.$$

be careful: $Q_M^* - Q_M^* \neq Q_M^{\pi} - Q_M^{\pi}$ for the same π .

\downarrow \downarrow

π_M^* π_M^*

$$\|Q_M^* - Q_M^*\|_{\infty} = \| \underbrace{Q_M^* - T_M^{\hat{\pi}} Q_M^*}_{\Delta} + T_M^{\hat{\pi}} Q_M^* - Q_M^* \|_{\infty}$$

$$\leq \underbrace{\|T_{\hat{M}} Q_{\hat{M}}^* - T_{\hat{M}} Q_M^*\|_{\infty}} + \|T_{\hat{M}} Q_M^* - Q_M^*\|_{\infty}$$

$$\leq \underbrace{\gamma \cdot \|Q_{\hat{M}}^* - Q_M^*\|_{\infty}} + \boxed{\|T_{\hat{M}} Q_M^* - Q_M^*\|_{\infty}}$$

$\forall s, a$

$$|(T_{\hat{M}} Q_M^*)(s, a) - Q_M^*(s, a)| \quad \uparrow V_M^*(s')$$

$$= \left| \hat{R}(s, a) + \gamma \mathbb{E}_{s' \sim \hat{P}(\cdot | s, a)} \left[\max_{a'} Q_M^*(s', a') \right] - Q_M^*(s, a) \right|$$

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

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

$[0, 0, \dots, 0, 1, 0, 0, \dots, 0]$

\uparrow
s'-th coordinate

$$= \left| \frac{1}{n} \sum_{(r, s') \in D_{s, a}} \underbrace{(r + \gamma V_M^*(s'))}_{\text{empirical Bellman update}} - Q_M^*(s, a) \right| \quad \boxed{(*)}$$

empirical Bellman update

Notice that since (r, s') in $D_{s,a}$ are iid.

$\Rightarrow r + \gamma V_M^*(s')$ are also iid for $(r, s') \in D_{s,a}$.

scalar. $\in [0, \frac{R_{\max}}{1-\gamma}]$.

$$r \leq R_{\max}. \quad V_M^* \leq \frac{R_{\max}}{1-\gamma}.$$

$$\begin{aligned} \text{Fix } s, a. \quad \mathbb{E}_{r, s'} [r + \gamma V_M^*(s')] & \quad \max_{a'} Q_M^*(s', a') \\ & \quad \parallel \\ & = R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [V_M^*(s')]. \\ & = \left(\mathcal{T}_M Q_M^* \right) (s, a) = Q_M^*(s, a). \end{aligned}$$

By Hoeffding's, w.p. $\geq 1 - \delta$.

$$(*) \leq \frac{R_{\max}}{1-\gamma} \cdot \sqrt{\frac{1}{2n} \ln \frac{2}{\delta}}.$$

By union bound, w.p. $\geq 1 - \delta$.

$$\| \mathcal{T}_M \hat{Q}_M^* - Q_M^* \|_{\infty} \leq V_{\max} \sqrt{\frac{1}{2n} \ln \frac{2|S \times A|}{\delta}}$$

$$\Rightarrow \| V_M^* - \hat{V}_M^* \|_{\infty} \leq \mathcal{O} \left(\frac{V_{\max}}{\sqrt{n} (1-\gamma)^2} \right)$$

Remark 1. $\forall f \in \mathbb{R}^{S \times A}$. $\mathbb{E}[\dots | s, a] = (Tf)(s, a)$

$$(T_{\hat{M}} f)(s, a) = \frac{1}{n} \sum_{(r, s') \in D_{s, a}} (r + \gamma \max_{a'} f(s', a'))$$

Previous analysis

emp. Bellman update
of f .

$T_{\hat{M}} \approx T_M$ in the sense that

$$\forall f \text{ (w/ bounded range)} \quad T_{\hat{M}} f \approx T_M f$$

This analysis: $T_{\hat{M}} f \approx T_M f$ only holds
for $f = Q_M^*$.

Remark 2: this analysis uses

$$\|Q_M^* - Q_{\hat{M}}^*\|_{\infty} \leq \frac{\|T_M Q_M^* - T_{\hat{M}} Q_M^*\|_{\infty}}{1 - \gamma}$$

But the following also holds:

$$\|Q_M^* - Q_{\hat{M}}^*\|_{\infty} \leq \frac{\|T_{\hat{M}} Q_{\hat{M}}^* - T_M Q_{\hat{M}}^*\|_{\infty}}{1 - \gamma}$$



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

Hoeffding's
doesn't apply

$$\left. - (R(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [V_M^*(s')]) \right|$$

Random

For fixed $f \in \mathbb{R}^S$.

$$\left| \frac{1}{n} \sum_{(r, s') \in D_{s, a}} (r + \gamma f(s')) \right.$$

Hoeffding's applies!

$$\left. - (R(s, a) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, a)} [f(s')]) \right|$$