## Recap

$$M = (S, A, P, r, \gamma)$$

<u>Generative model (simulator):</u> $\forall (s,a) \in S \times A,$ draw $N$ samples $s' \sim P(\cdot | s, a)$

<u>Model-based estimator</u> ("Plug-in"): $\hat{P}(s' | s, a) = \dfrac{\text{count}(s, a, s')}{N}$

$$\hat{M} = (S, A, \hat{P}, r, \gamma) \xrightarrow{\text{planning}} \pi^*_{\hat{M}} =: \hat{\pi}$$

Two metrics:

- (Value estimate)  $\| Q^* - \hat{Q}^* \|_\infty$
- (policy estimate)  $\| Q^* - Q^{\hat{\pi}} \|_\infty$

$$H = \frac{1}{1-\gamma} \; : \; \text{horizon}$$

## Result recaps & plan :

- lower bound ( value & policy) : #samples $= N \cdot |S| \cdot |A| = \tilde{\Omega}\left( \frac{H^3 \; |S| \cdot |A|}{\varepsilon^2} \right)$

     (Azar et al. '18)

- upper bounds:

• (coarse analysis wf uniform convergene) : #samples$_{\text{policy/value}} = \tilde{O}\left( \frac{H^4 \; |S|^2 \cdot |A|}{\varepsilon^2} \right)$

• Today:

(*) #samples$_{\text{value}} = \tilde{O}\left( H^4 \cdot |S| \cdot |A| / \varepsilon^2 \right)$

(**) #samples$_{\text{value}} = \tilde{O}\left( H^3 \cdot \frac{|S| \cdot |A|}{\varepsilon^2} \right)$     (Azar et al. '13)

     #samples$_{\text{policy}} = \tilde{O}\left( H^3 \frac{|S| \cdot |A|}{\varepsilon^2} \right)$ for $\varepsilon \in (0, \sqrt{\frac{H}{\gamma |S|}})$   (Azar)

(***) #sample$_{\text{policy}} = \tilde{O}\left( H^3 \frac{|S| \cdot |A|}{\varepsilon^2} \right)$ $\forall \varepsilon \in (0, 1)$    (Agarwal '20)

(*)

Theorem

(value estimation)

$$\text{wp } 1-\delta$$
$$\| Q^* - \hat{Q}^* \|_\infty \leq \gamma H^2 \sqrt{\frac{2 \log\left(|S| \cdot |A| / \delta\right)}{N}}$$

Claim:

$$\| Q^* - \hat{Q}^* \|_\infty \leq \frac{\gamma}{1-\gamma} \| (P - \hat{P}) V^* \|_\infty$$

Proof:

$$\| Q^* - \hat{Q}^* \|_\infty = \| \gamma P^{\pi^*} Q^* - \gamma \hat{P}^{\hat{\pi}} \hat{Q}^* \|_\infty$$

$$= \gamma \| P^{\pi^*} Q^* - \hat{P}^{\pi^*} Q^* + \hat{P}^{\pi^*} Q^* - \hat{P}^{\hat{\pi}} \hat{Q}^* \|_\infty$$

$$\leq \gamma \| P^{\pi^*} Q^* - \hat{P}^{\pi^*} Q^* \|_\infty + \gamma \| \hat{P}^{\pi^*} Q^* - \hat{P}^{\hat{\pi}} \hat{Q}^* \|_\infty$$

$$= \gamma \| P V^* - \hat{P} V^* \|_\infty + \gamma \| \hat{P} V^* - \hat{P} \hat{V}^* \|_\infty$$

$$\leq \gamma \| (P - \hat{P}) V^* \|_\infty + \gamma \underbrace{\| V^* - \hat{V}^* \|_\infty}_{\leq \| Q^* - \hat{Q}^* \|_\infty}$$

$$\| (P - \hat{P}) V^{*} \|_{\infty} = \max_{(s,a)} \left| \left( P(\cdot | s,a) - \hat{P}(\cdot | s,a) \right)^{\top} V^{*} \right|$$

$$\leq \frac{1}{1-\gamma} \sqrt{\frac{\log (|S| \cdot |A| / \delta)}{N}} \quad \text{(Hoeffding's)}$$

Theorem :    (**)

- (value bound)    $\|Q^* - \hat{Q}^*\|_\infty \leq \gamma \sqrt{H^3 \dfrac{\log(|S| \cdot |A|/\delta)}{N}} + \dfrac{\gamma}{(1-\gamma)^3} \dfrac{\log(|S| \cdot |A|/\delta)}{N}$

## Lemma (point-wise bounds)

$$Q^* - \hat{Q}^* \leq \gamma \left(I - \gamma \hat{P}^{\pi^*}\right)^{-1} (P - \hat{P}) V^*$$

$$Q^* - \hat{Q}^* \geq \gamma \left(I - \gamma \hat{P}^{\hat{\pi}}\right)^{-1} (P - \hat{P}) V^*$$

**Proof:**

• $Q^* - \hat{Q}^* \leq Q^* - \hat{Q}^{\pi^*} = \gamma \left(I - \gamma \hat{P}^{\pi^*}\right)^{-1} (P - \hat{P}) V^*$

• Recall: $Q^* = r + \gamma P Q^*$, $\hat{Q}^* = r + \gamma \hat{P}^{\hat{\pi}} \hat{Q}^*$

**Thus,** $Q^* - \hat{Q}^* = Q^* - \left(I - \gamma \hat{P}^{\hat{\pi}}\right)^{-1} r$

$$= Q^* - \left(I - \gamma \hat{P}^{\hat{\pi}}\right)^{-1} \left(I - \gamma P^{\pi^*}\right) Q^*$$

$$= \left(I - \gamma \hat{P}^{\hat{\pi}}\right)^{-1} \left[ \left(I - \gamma \hat{P}^{\hat{\pi}}\right) - \left(I - \gamma P^{\pi^*}\right) \right] Q^*$$

$$= \gamma \left(I - \gamma \hat{P}^{\hat{\pi}}\right)^{-1} \left(P^{\pi^*} - \hat{P}^{\hat{\pi}}\right) Q^*$$

**Note:** $P^{\pi^*} Q^* = P V^*$; $\hat{P}^{\hat{\pi}} Q^* \leq \hat{P}^{\pi^*} Q^* = \hat{P} V^*$

( point-wise bounds)

Claim

①

$$Q^* - \hat{Q}^* \leq \gamma \left( I - \gamma \hat{P}^{\pi^*} \right)^{-1} (P - \hat{P}) V^*$$

$$Q^* - \hat{Q}^* \geq \gamma \left( I - \gamma \hat{P}^{\hat{\pi}} \right)^{-1} (P - \hat{P}) V^*$$

Bernstein's inequality:

②

$$\left| (P - \hat{P}) V^* \right| \leq \sqrt{\text{Var}_P(V^*) \frac{\log(|S| \cdot |A| / \delta)}{N}} + H \frac{\log(|S| \cdot |A| / \delta)}{N} \mathbb{1}$$

here $\left[ \text{Var}_P(V^*) \right](s,a) = \underset{s' \sim P(\cdot | s, a)}{\text{Var}} \left[ V^*(s') \right]$

Plug in ② into ①:

$$\bullet \quad Q^* - \hat{Q}^* \leq \gamma \left( I - \gamma \hat{P}^{\pi^*} \right)^{-1} \sqrt{\widetilde{\text{Var}_P(V^*)}} \cdot \tilde{O}\left( \frac{1}{\sqrt{N}} \right)$$

$$+ \quad \tilde{O}\left( \left( \frac{1}{1-\gamma} \right)^2 \frac{1}{N} \right)$$

$$\bullet \quad Q^* - \hat{Q}^* \geq -\gamma \left( I - \gamma \hat{P}^{\hat{\pi}} \right)^{-1} \sqrt{\widetilde{\text{Var}_P(V^*)}} \cdot \tilde{O}\left( \frac{\hat{1}}{\sqrt{N}} \right) +$$

It remains to upper bound:

$$\left\| \left( I - \gamma \hat{P}^{\pi} \right)^{-1} \sqrt{\text{Var}_P(V^{\pi})} \right\|_{\infty}$$

where $\pi \in \{ \pi^*, \hat{\pi} \}$

- A trivial upper bound: $\left( \frac{1}{1-\gamma} \right)^2$

- A more intricate analysis gives: $\left( \frac{1}{1-\gamma} \right)^{3/2} + \tilde{O} \left( \frac{1}{1-\gamma} \cdot \frac{1}{N^{1/4}} + \left( \frac{1}{1-\gamma} \right)^2 \frac{1}{\sqrt{N}} \right)$

⑤

$$\boxed{\text{Claim} \qquad \left\| \left(I - \gamma P^\pi\right)^{-1} \sqrt{\text{Var}_P(V^\pi)} \right\|_\infty \leq \sqrt{H^3}}$$

- Bellman equation for variance

$$\boxed{\Sigma_M^\pi = \gamma^2 \left(I - \gamma^2 P^\pi\right)^{-1} \text{Var}_P(V_M^\pi)}$$

where:

$$\Sigma_N^\pi = \mathbb{E}\left[ \left( \sum_{t=0}^\infty \gamma^t r_t - V_M^\pi(s,a) \right)^2 \middle| (s_0,a_0)=(s,a) \right]$$

- $\left\| \Sigma_M^\pi \right\|_\infty \leq H^2$

- $\left\| \left(I - \gamma P^\pi\right)^{-1} \sqrt{V} \right\|_\infty = \dfrac{1}{1-\gamma} \left\| (1-\gamma)\left(1-\gamma P^\pi\right)^{-1} \sqrt{V} \right\|_\infty$

$$\because \left\| \left(1-\gamma P^\pi\right)^{-1} V \right\|_\infty \leq \left\| 2 \left(1 - \gamma^2 P^\pi\right)^{-1} V \right\|_\infty$$

$$\leq \frac{1}{1-\gamma} \left\| \sqrt{(1-\gamma)\left(1-\gamma P^\pi\right)^{-1} V} \right\|_\infty$$

$$\leq \frac{1}{1-\gamma} \left\| \sqrt{2(1-\gamma)\left(1-\gamma^2 P^\pi\right)^{-1} V} \right\|_\infty$$

$$= \frac{2}{\sqrt{1-\gamma}} \left\| \sqrt{\left(1-\gamma^2 P^\pi\right)^{-1} V} \right\|_\infty$$

$$\left\{ \begin{array}{l} \text{Var}_P \left( V^{\#} \right) \longrightarrow \text{Var}_{\hat{P}} \left( \hat{V}^{k} \right) \ ? \\[20pt] \text{Var}_P \left( V^{*} \right) \longrightarrow \text{Var}_{\hat{P}} \left( \hat{V}^{\hat{\pi}} \right) \ ? \end{array} \right.$$

$\longleftarrow$ Crude bounds
are fine !
(e.g. Hoeffding's ineq)

$$\bullet \quad \mathrm{Var}_P(V^*) = \mathrm{Var}_P(V^*) - \mathrm{Var}_{\hat{P}}(V^*) + \mathrm{Var}_{\hat{P}}(V^*)$$

$$= \left[ P(V^*)^2 - (P V^*)^2 \right] - \left[ \hat{P}(V^*)^2 - (\hat{P} V^*)^2 \right] + \mathrm{Var}_{\hat{P}}(V^*)$$

$$= (P - \hat{P})(V^*)^2 + \underbrace{(\hat{P} - P) V^* \left[ (\hat{P} + P) V^* \right]}_{} + \mathrm{Var}_{\hat{P}}(V^*)$$

$$= \breve{O}\left( \left( \frac{1}{1-\gamma} \right)^2 \frac{1}{\sqrt{N}} \right) + \breve{O}\left( \left( \frac{1}{1-\gamma} \right)^2 \frac{1}{\sqrt{N}} \right) + \mathrm{Var}_{\hat{P}}(V^*)$$

$$\bullet \quad \mathrm{Var}_{\hat{P}}(V^*) = \mathrm{Var}_{\hat{P}}\left( \hat{V}^* + V^* - \hat{V}^* \right)$$

$$\leq \left( \sqrt{\mathrm{Var}_{\hat{P}}(\hat{V}^*)} + \sqrt{\mathrm{Var}_{\hat{P}}(V^* - \hat{V}^*)} \right)^2$$

$$\leq 2 \left( \mathrm{Var}_{\hat{P}}(\hat{V}^*) + \mathrm{Var}_{\hat{P}}(V^* - \hat{V}^*) \right)$$

$$\leq \| V^* - \hat{V}^* \|_\infty^2$$

$$= \breve{O}\left( \left( \frac{1}{1-\gamma} \right)^4 \frac{1}{N} \right)$$

## Theorem (policy bound)  (✳✳✳)

- (policy bound)

$$\|Q^* - Q^{\hat{\pi}}\|_\infty \leq \gamma \sqrt{H^3 \frac{\log(|s| \cdot |A| / \delta)}{N}} \; , \; \text{if } N \geq H^2$$

more subtle

Ref:  Agarwal & Kakade & Yang : "Model-based RL w/ a generative model is minimax optimal", 2020

Error decomposition:

- $\underbrace{Q^* - Q^{\hat{\pi}}}_{\geq 0} = Q^* - \hat{Q}^{\pi^*} + \hat{Q}^{\pi^*} - Q^{\hat{\pi}}$

$$\leq \underbrace{Q^* - \hat{Q}^{\pi^*} + \hat{Q}^{\hat{\pi}} - Q^{\hat{\pi}}}$$
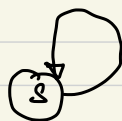
.

$$\Rightarrow \| Q^* - Q^{\hat{\pi}} \|_\infty \leq \| Q^* - \hat{Q}^{\pi^*} \|_\infty + \| \hat{Q}^{\hat{\pi}} - Q^{\hat{\pi}} \|_\infty$$

Simulation lemma:

$$Q^* - \hat{Q}^{\pi^*} = \gamma (I - \gamma P^{\pi^*})(P - \hat{P}) \hat{V}^{\pi^*}$$
$$\hat{Q}^{\hat{\pi}} - Q^{\hat{\pi}} = \gamma (I - \gamma P^{\hat{\pi}})(P - \hat{P}) \hat{V}^{\hat{\pi}}$$

<u>idea:</u>     s-absorbing MDP:

- For any $s \in S$, $u \in \mathbb{R}$, $M_{s,u}$ is identical to $M$ except that:



$$\begin{cases} P_{M_{s,u}}(s \mid s,a) = 1 \quad \forall a \\ r_{M_{s,u}}(s,a) = u \quad \forall a \end{cases}$$

- $$(P - \hat{P})_{s,a} \hat{V}^* = (P - \hat{P})_{s,a} V^*_{\widehat{M}_{s,u}} + (P - \hat{P})_{s,a} \underbrace{\left( \hat{V}^* - V^*_{\widehat{M}_{s,u}} \right)}_{\leq \; \| \hat{V}^* - V^*_{\widehat{M}_{s,u}} \|_\infty =: \Delta}$$

Note: $\hat{P}_{s,a} \perp V^*_{\widehat{M}_{s,u}}$

<u>Bernstein:</u>   Fix $(s,a)$,   $\forall u \in \mathcal{U}$:

$$(P_{s,a} - \hat{P}_{s,a}) \hat{V}^* \leq \sqrt{\frac{\mathrm{Var}_{P_{s,a}}(V^*_{\widehat{M}_{s,u}}) \log(|\mathcal{U}|)}{N}} + \frac{\log(|\mathcal{U}|)}{(1-\delta) N} + \Delta$$

$$\leq \sqrt{\frac{\text{Var}_{P_{s,a}}\left(\hat{V}^{*} + V^{*}_{\hat{M}_{s,a}} - \hat{V}^{*}\right) \cdot \log u}{N}} + \frac{\log u}{(1-\sigma)N} + \triangle$$

$$= \sqrt{\frac{\text{Var}_{P_{s,a}}\left(\hat{V}^{*}\right)}{N} \log(|u|)} + \inf_{u \in U} \| \hat{V}^{*} - V^{*}_{\tilde{M}_{s,u}} \| \left(1 + \frac{\log u}{\sqrt{N}}\right)$$
$$+ \frac{\log u}{(1-\sigma)N}$$

Controlling: $\left\| V^{*}_{\hat{M}_{s,4}} - V^{*}_{\hat{M}} \right\|_{\infty}$   (Goal: $\tilde{O}\left(\frac{1}{\sqrt{N}}\right)$ )

$$= \left| u - \hat{V}^{*}(s) \right|$$

choose $\mathcal{U}$: even interval of $\left[ V^{*}(s) \pm \tilde{O}\left( \left(\frac{1}{1-\gamma}\right)^{2} \frac{1}{\sqrt{N}} \right) \right]$

$$\min_{u \in \mathcal{U}} \left| u - \hat{V}^{*}(s) \right| \leq \frac{\tilde{O}\left( \frac{1}{(1-\gamma)^{2}} \frac{1}{\sqrt{N}} \right)}{|\mathcal{U}| - 1} = \tilde{O}\left( \frac{1}{\sqrt{N}} \right)$$

$\Rightarrow$ choos $|\mathcal{U}| - 1 = \frac{1}{(1-\gamma)^{2}}$