

Batch RL

= fake RL

↓ no interaction!

Control

Data! → Policies

~~MDP~~

Who collected →

How was the collected?

Do we know this?

Yes No

Experimental Design: How to collect data

How to use the data? O.P.O.

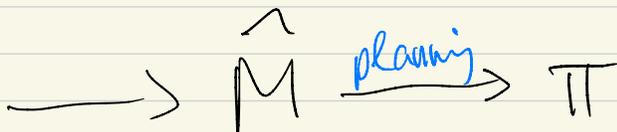
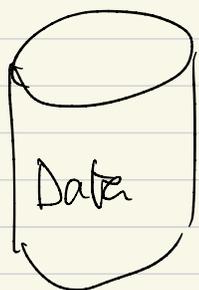
↓ O.P.E

Certification of performance

conservative explanation

Limits / how to approach?

Model-based P.O.



Model-sensitivity

$$0 \leq \sigma < 1$$

$$\hat{M} = (S, A, \hat{P}, \hat{r}, \sigma)$$

\uparrow
 $\pi(a|s)$

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

$$(*) \quad \hat{v}^\pi \geq \hat{v}^* - \varepsilon \mathbb{1}$$

$$\Rightarrow v^\pi \geq v^* - \delta \mathbb{1}$$

$$\delta \leq f(\varepsilon, \frac{1}{1-\sigma}, M, \hat{M})$$

$$f = ?$$

Polig error bound:

$$\textcircled{1} \quad \underline{M_\pi q^*} \geq v^* - \underline{\varepsilon \mathbb{1}} \Rightarrow v^\pi \geq v^* - \frac{\varepsilon}{1-\sigma} \mathbb{1}$$

$$\textcircled{2} \quad \underline{M_\pi q} = Mq \Rightarrow v^\pi \geq v^* - \frac{2\|q - q^*\|_\infty}{1-\sigma} \mathbb{1}$$

$$q = \hat{q}^\pi \quad \hat{q}^\pi \leq \hat{q}^* \Rightarrow M\hat{q}^\pi \leq M\hat{q}^* = v^*$$

$$\underline{M\hat{q}^\pi} \geq \underline{M_\pi \hat{q}^\pi} = \underline{\hat{v}^\pi} \stackrel{(*)}{\geq} \underline{\hat{v}^* - \varepsilon \mathbb{1}} = \underline{M\hat{q}^* - \varepsilon \mathbb{1}} \geq \underline{M\hat{q}^\pi - \varepsilon \mathbb{1}}$$

$$M_\pi \hat{q}^\pi = M\hat{q}^\pi + z, \quad \|z\|_\infty \leq \varepsilon$$

$$M_\pi q = Mq + z$$

$$M_\pi q^* = M_\pi q + M_\pi(q^* - q) \leq Mq + M_\pi(q^* - q) + z$$

$$= Mq^* + \underbrace{Mq - Mq^*}_{\leq -\|q - q^*\|_\infty} + M_\pi(q^* - q) + z$$

$$\|Mq - Mq^*\|_\infty \leq \|q - q^*\|_\infty$$

$$Mq - Mq^* \geq -\|q - q^*\|_\infty \mathbf{1}$$

$$\|M_\pi(q^* - q)\|_\infty \leq \|q - q^*\|_\infty \Rightarrow M_\pi(q^* - q) \geq -\|q - q^*\|_\infty \mathbf{1}$$

$$z \geq -\|z\|_\infty \mathbf{1}$$

$$M_\pi q \geq Mq - \underbrace{(2\|q - q^*\|_\infty + \|z\|_\infty)}_{\leq \epsilon} \mathbf{1}$$

$$\Rightarrow \boxed{v^\pi \geq v^* - \frac{2\|q - q^*\|_\infty + \|z\|_\infty}{1 - \sigma} \mathbf{1}}$$

$$q := \hat{q}^\pi$$

$$v^\pi \geq v^* -$$

$$\boxed{\frac{2\|\hat{q}^\pi - q^*\|_\infty + \epsilon}{1 - \sigma} \mathbf{1}}$$

$$\|\hat{q}^\pi - q^*\|_\infty \leq \underbrace{\|\hat{q}^\pi - \hat{q}^*\|_\infty}_{\leq \epsilon} + \underbrace{\|\hat{q}^* - q^*\|_\infty}_{\text{?}}$$

← general.
 T σ -contraction, $Tx = x$

$$\forall y \quad \|x - y\| \leq \frac{\|Ty - y\|}{1 - \sigma}$$

$$\|x - y\| \leq \|x - Ty\| + \|Ty - y\| \leq \sigma \|x - y\| + \|Ty - y\|$$

$$\| \hat{q}^* - q^* \|_\infty \leq \frac{\| \hat{T} q^* - q^* \|_\infty}{1 - \sigma}$$

$x = q^* \quad T := \hat{T}$

$$\begin{aligned} \hat{T}q &= \hat{r} + \sigma \hat{P}Mq \\ Tq &= r + \sigma P Mq \end{aligned}$$

$$\| \hat{T} q^* - q^* \|_\infty = \| \hat{T} q^* - T q^* \|_\infty$$

$$\leq \| \hat{r} + \sigma \hat{P} v^* - (r + \sigma P v^*) \|_\infty$$

$$\leq \| r - \hat{r} \|_\infty + \sigma \| (\hat{P} - P) v^* \|_\infty$$

$$\| \hat{q}^* - q^* \| \leq \frac{\| T \hat{q}^* - \hat{q}^* \|_\infty}{1 - \sigma}$$

$$\leq \| \hat{P} - P \|_\infty \frac{\| v^* \|_\infty}{1 - \sigma}$$

$$\| T \hat{q}^* - \hat{q}^* \|_\infty \leq \| r - \hat{r} \|_\infty + \sigma \| (\hat{P} - P) v^* \|_\infty$$

Thm : $\forall \pi \quad \hat{v}^\pi \geq v^* - \epsilon \mathbb{1}$

$\Rightarrow N^\pi \geq v^* - \delta \mathbb{1}$

$\begin{cases} \min(a, b) \\ = a \wedge b \end{cases}$

$\delta \leq \frac{\epsilon \sqrt{2}}{1-\gamma} + \frac{2}{(1-\gamma)^2} \left[\|r - \hat{r}\|_\infty + \gamma \|(\hat{P} - P)v^*\|_\infty \right]$
 $\wedge \left[\|r - \hat{r}\|_\infty + \gamma \|(\hat{P} - P)v^*\|_\infty \right]$

$\leq \frac{\epsilon \sqrt{2}}{1-\gamma} + \frac{2}{(1-\gamma)^2} \left[\|r - \hat{r}\|_\infty + \frac{\gamma \|P - \hat{P}\|_\infty}{1-\gamma} \right]$

repl

Batch RL

① Tabular

② Featurized

① trajectories; π_b
 ↑
 behavior policy / logging policy

② offline access to sim $N(s, a)$

$t = 1, \dots, n$

$\{S_{t-1}, A_{t-1}, S_t, R_t\}$

$S_t' \sim P_{A_t}(S_{t-1}, \cdot)$ (states)
 $R_t \sim \tilde{P}_{A_t}(S_{t-1}, \cdot)$ (reals)

$\hat{r}_a(s) = \begin{cases} \frac{1}{N(s, a)} \sum_{t=1}^n \mathbb{I}(S_t = s, A_t = a) R_t \\ 0 \end{cases}$
 $N(s, a) = 0$

$$N(s, a) = \sum_{t=1}^n \mathbb{I}(S_t = s, A_t = a) \quad N(s, a) \neq 0$$

$$\hat{P}_\alpha(s, s') = \begin{cases} \frac{1}{N(s, a)} \sum_{t=1}^n \mathbb{I}(S_t = s, A_t = a) \mathbb{I}(S_t = s') \\ \frac{1}{|S|}, \quad N(s, a) = 0 \end{cases}$$

$$N(s, a) = n(s, a) \quad \text{if } > 0 \quad R_t \in [0, 1]$$

$$|\hat{r}_\alpha(s) - r_\alpha(s)| \leq \sqrt{\frac{\log\left(\frac{SA}{\delta}\right)}{2n(s, a)}} \quad \left. \begin{array}{l} \text{wp } 1-\delta \\ \forall s, a, s' \end{array} \right\}$$

$$|\hat{P}_\alpha(s, s') - P_\alpha(s, s')| \leq \sqrt{\frac{\log\left(\frac{S^2 A}{\delta}\right)}{2n(s, a)}}$$

$$\|\hat{P} - P\|_\infty = \max_{s, a} \|\hat{P}_\alpha(s) - P_\alpha(s)\|$$

$$\leq S \sqrt{\frac{\log\left(\frac{S^2 A}{\delta}\right)}{2n(s, a)}}$$

$$\|r - \hat{r}\|_\infty \leq \sqrt{\frac{\log\left(\frac{S^2 A}{\delta}\right)}{2n(s, a)}} \leq \sqrt{\frac{\log\left(\frac{SA}{\delta}\right)}{n_{\min}}}$$

$n_{\min} = \min_{(s, a)} n(s, a)$

$$\delta \leq \frac{2}{1-\delta^2} \left[\sqrt{\frac{\log\left(\frac{S^2 A}{\delta}\right)}{2n_{\min}}} + \frac{S}{1-\delta} \sqrt{\frac{\log\left(\frac{S^2 A}{\delta}\right)}{2n_{\min}}} \right] = \text{UB}(\delta)$$

$$n(s, a) = ?$$

$$\sum_{s, a} n(s, a) = n$$

max n_{\min}

$$n(s, a) = \frac{1}{SA} = n_{\min}$$

$$\delta \leq \frac{2}{1-\delta^2} \left[\sqrt{\frac{SA \epsilon}{2n}} \left(1 + \frac{S}{1-\delta} \right) \right]$$

$$\epsilon = H^3 \left[\sqrt{\frac{SA}{n}} S \right]$$

$$\frac{\epsilon}{SH^3} = \sqrt{\frac{SA}{n}}$$

$$n \approx SA \left(\frac{SH^3}{\epsilon} \right)^2$$

$$\left(\frac{SH^3}{\epsilon} \right)^2 SA$$

$$= \frac{SA}{\epsilon^2} H^6$$

traje donec

→ φ