

# CS292F StatRL Lecture 15

## Uniform OPE and Near-Optimal Offline Learning

Instructor: Yu-Xiang Wang

Spring 2021

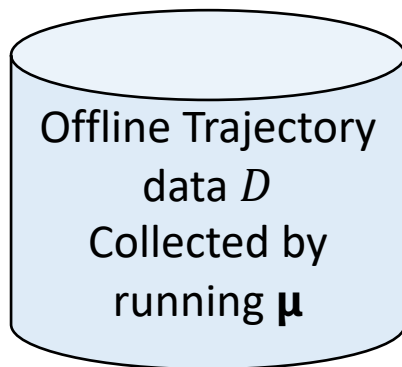
UC Santa Barbara

# Logistics

- Homework 3 is released.
  - Due end of the quarter / June 2
  - 5 questions, but you only need to do either Q4 or Q5.
- Two more lectures on offline RL.
  - including this one.
- I will send schedules for the project presentations this week

# Recap: Offline Reinforcement Learning, aka. Batch RL

- Task 1: Offline Policy Evaluation. (OPE)

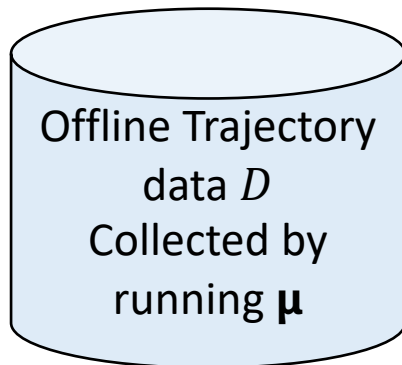


Task: design OPE  
methods

Evaluate fixed Target  
Policy  $\pi$

**Via  
Uniform  
OPE**

- Task 2: Offline Policy Learning. (OPL)



Task: design OPO  
methods

Find near optimal  
Policy  $\hat{\pi}^*$

OPE: the tabular MDP case, they are all equivalent.

- TMIS 
$$\hat{v}_{\text{MIS}}^{\pi} = \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^H \frac{\hat{d}_t^{\pi}(s_t^{(i)})}{\hat{d}_t^{\mu}(s_t^{(i)})} \hat{r}_t^{\pi}(s^{(i)}).$$

- Model-based Plugin

$$\hat{v}_{\text{DM}}^{\pi} = \sum_{h=1}^H \sum_{s \in \mathcal{S}} \hat{d}_h^{\pi}(s) \hat{r}_h^{\pi}(s)$$

- Fitted Q Iteration

$$\hat{v}_{\text{FQI}}^{\pi} = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \hat{d}_1(s) \pi(a|s) \hat{Q}_1(s, a)$$

# They are also information-theoretically optimal

## Offline Policy Evaluation

Simulation lemma (Kearns and Singh, 1998)	IS / DR (Jiang and Li, 2016)	MIS (Xie, Ma, W.,2019)	TMIS (Yin & W. 2020)	Fitted Q-Iteration (Duan and Wang, 2020)
$\sqrt{\frac{H^4 S^2}{n d_m}}$	$\sqrt{\frac{e^H \text{poly}(S, A)}{n}}$	$\sqrt{\frac{H^3}{n d_m}}$	$\sqrt{\frac{H^2}{n d_m}}$	$\sqrt{\frac{H^2}{n d_m}}$

**Per-instance optimal.**

$$\begin{aligned}
 & \mathbb{E}[(\hat{v}_{\text{TMIS}}^\pi - v^\pi)^2] \\
 & \leq \frac{1}{n} \sum_{h=0}^H \sum_{s_h, a_h} \frac{d_h^\pi(s_h)^2 \pi(a_h|s_h)^2}{d_h^\mu(s_h) \mu(a_h|s_h)} \text{Var} \left[ (V_{h+1}^\pi(s_{h+1}^{(1)}) + r_h^{(1)}) \middle| s_h^{(1)} = s_h, a_h^{(1)} = a_h \right] \\
 & \qquad \qquad \qquad + O(n^{-1.5})
 \end{aligned}$$

**Matching Cramer-Rao lower bound up to low-order terms.**

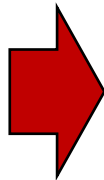
# Recap: From OPE to offline learning

- Empirical Risk Minimization (ERM)?

$$\hat{\pi} = \arg \max_{\pi \in \Pi} \hat{v}^{\pi} \quad (\text{For some OPE estimator } \hat{v}^{\pi})$$

- A uniform convergence argument

$$\sup_{\pi \in \Pi} |\hat{v}^{\pi} - v^{\pi}| \leq \epsilon \quad \text{w.h.p.} \quad v^{\pi^*} - v^{\hat{\pi}} \leq 2\epsilon \quad \text{w.h.p.}$$



$$\mathbb{E} \left[ \sup_{\pi \in \Pi} |\hat{v}^{\pi} - v^{\pi}|^2 \right] \leq \epsilon^2 \quad v^{\pi^*} - \mathbb{E}[v^{\hat{\pi}}] \leq 2\epsilon$$

# Recap: exploration assumption

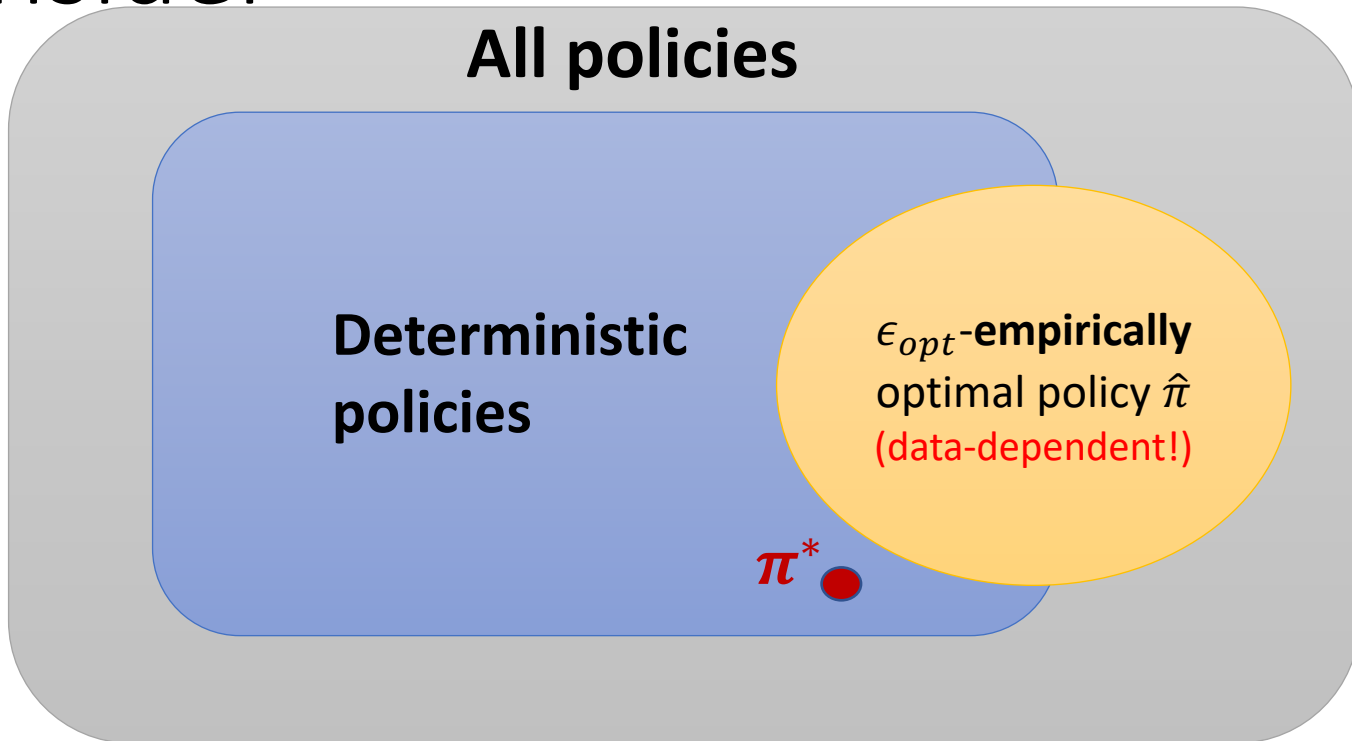
- The logging policy  $\mu$  is out of our control
- Need to make assumptions about it

$$d_m := \min_{t,s,a} d_t^\mu(s,a) > 0 \text{ for all } t, s, a$$

$$\text{s.t. } d_t^\pi(s,a) > 0 \text{ for some } \pi \in \Pi$$

- Assumed to simplify the discussion on optimality
- Sometimes appear only in low-order terms.

# Recap: The policy classes we consider



For ERM, it suffices to consider the smaller policy class.  
But we also want to cover other planning algorithms.

The remainder of the lecture is based on:

Yin, Bai and W. (2020) <https://arxiv.org/pdf/2007.03760.pdf>



# This lecture

- Characterize the uniform OPE on deterministic policy class
- Optimal offline learning via a local uniform OPE

# Recap: counting the number of deterministic policies

- Setting: Tabular MDP with  $S$  states,  $A$  actions and  $H$  steps.
- How many deterministic policies are there?  
(Answer:  $A^{\{SH\}}$ )
- Together with a high-probability OPE bound

# Uniform convergence theorem for all **deterministic** policies

**Theorem 3.5:** with probability  $\geq 1 - \delta$

$$\sup_{\pi \in \Pi_{\text{deterministic}}} |\hat{v}^{\pi} - v^{\pi}| \lesssim \sqrt{\frac{H^3 S}{nd_m} \log\left(\frac{HSA}{\delta}\right)} + O(1/n)$$

- **Optimal in H.**
- **Suboptimal in S?**
- **Proof:** Union bound with a high-probability pointwise OPE bound.

# How do we obtain a high-probability pointwise OPE bound?

- Steps in the analysis

1. Fictitious estimator technique and multiplicative Chernoff bounds.

- Nice event:  $E_t := \{n_{s_t, a_t} \geq n d_t^\mu(s_t, a_t) / 2\}$

2. Error decomposition (reducing to the case with known reward function)

3. Further decomposition of the occupancy measure into a Martingale.

4. Apply Freedman's inequality --- a Bernstein-style Martingale Concentration.

Step 1 Recap: The fictitious estimator is easier to analyze, because:

- Always unbiased.
- Has an *epistemic* Bellman-equation of variance
- Has nice martingale decompositions
- Moreover: Lemma C.3

$$\sup_{\pi \in \Pi} |\tilde{v}^{\pi} - \hat{v}^{\pi}| = 0 \quad \text{w.h.p.}$$

Under mild condition:  $n \gtrsim \frac{1}{d_m} \log \frac{HSA}{\delta}$

Step 2 Recap: The noise in the reward is straightforward to handle.

$$\begin{aligned}
 \sup_{\pi \in \Pi} |\tilde{v}^\pi - v^\pi| &= \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t \rangle - \sum_{t=1}^H \langle d_t^\pi, r_t \rangle \right| \\
 &= \sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t \rangle - \sum_{t=1}^H \langle \tilde{d}_t^\pi, r_t \rangle + \sum_{t=1}^H \langle \tilde{d}_t^\pi, r_t \rangle - \sum_{t=1}^H \langle d_t^\pi, r_t \rangle \right| \\
 &\leq \underbrace{\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right|}_{(*)} + \underbrace{\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle \right|}_{(**)}
 \end{aligned}$$

**Lemma C.4:**  $(**) \lesssim \sqrt{H^2 / (nd_m)}$

Therefore, it suffices to consider the case with **deterministic rewards**.

# Step 3 Recap: Martingale decomposition of the error $\tilde{v}^\pi - v^\pi$

**Primal representation (Marginal distribution style):**

$$\sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle$$

|| (Lemma C.5)

**Dual representation (Value function style):**

$$\langle v_1^\pi(s), (\tilde{d}_1^\pi - d_1^\pi)(s) \rangle + \sum_{h=2}^H \langle v_h^\pi(s), ((\tilde{T}_h - T_h)\tilde{d}_{h-1}^\pi)(s) \rangle$$

(You can prove this by simulation lemma, see HW3 Q5.)

Let's check that this is a Martingale  
(w.r.t. the parallel data sequence)

Recall definition of Martingale:

- a.  $E[X_t | D_{\{t-1\}}] = X_{\{t-1\}}$
- b.  $E[|X_t|]$  is bounded

$$X_t := \sum_{h=2}^t \langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle + \langle V_1^\pi, \tilde{d}_1^\pi - d_1^\pi \rangle.$$



# Step 4: Freedman's inequality

**Lemma A.6** (Freedman's inequality [Tropp et al. \(2011\)](#)). *Let  $X$  be the martingale associated with a filter  $\mathcal{F}$  (i.e.  $X_i = \mathbb{E}[X|\mathcal{F}_i]$ ) satisfying  $|X_i - X_{i-1}| \leq M$  for  $i = 1, \dots, n$ . Denote  $W := \sum_{i=1}^n \text{Var}(X_i|\mathcal{F}_{i-1})$  then we have*

$$\mathbb{P}(|X - \mathbb{E}[X]| \geq \epsilon, W \leq \sigma^2) \leq 2e^{-\frac{\epsilon^2}{2(\sigma^2 + M\epsilon/3)}}.$$

*Or in other words, with probability  $1 - \delta$ ,*

$$|X - \mathbb{E}[X]| \leq \sqrt{8\sigma^2 \cdot \log(1/\delta)} + \frac{2M}{3} \cdot \log(1/\delta), \quad \text{Or } W \geq \sigma^2.$$

- To apply this inequality, we need to
  - bound  $M$
  - Need to work out the variance.

\*In fact we will use a more flexible version of Freedman's inequality due to ([Chung and Liu 2006](#)) that allows the bound  $M$  to hold w.h.p rather than with prob 1.

# Bound M in high probability

- Lemma E.2 With prob at least  $1-\delta$

$$\sup_t |X_t - X_{t-1}| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot d_m}}\right).$$

Proof:

$$\begin{aligned} |X_t - X_{t-1}| &= \langle V_t^\pi, (\tilde{T}_t - T_t) \tilde{d}_{t-1}^\pi \rangle \\ &\leq \|(\tilde{T}_t - T_t)^T V_t^\pi\|_\infty \|\tilde{d}_{t-1}^\pi\|_1 = \|(\tilde{T}_t - T_t)^T V_t^\pi\|_\infty. \end{aligned}$$

# Bound W: Sum of conditional variance

**Lemma E.3.** *We have the following decomposition of conditional variance:*

$$\text{Var}[X_{t+1}|\mathcal{D}_t] = \sum_{s_t, a_t} \frac{\tilde{d}_t^\pi(s_t, a_t)^2 \cdot \mathbf{1}(E_t)}{n_{s_t, a_t}} \cdot \text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)}) | s_t^{(1)} = s_t, a_t^{(1)} = a_t]$$

**Proof:**

# Bound W: Sum of conditional variance

- If we can bound  $\tilde{d}_t^\pi(s_t, a_t)$
- Then, we can apply Lemma 3.4 from the last lecture

$$\sum_{t=1}^H \text{Var}[X_{t+1} | \mathcal{D}_t] \leq O\left(\frac{1}{nd_m} \cdot \sum_{t=1}^H \mathbb{E}[\text{Var}[V_{t+1}^\pi(s_{t+1}^{(1)}) | s_t^{(1)}, a_t^{(1)}]]\right) \leq O\left(\frac{H^2}{nd_m}\right)$$

# Bounding $\tilde{d}_t^\pi(s_t, a_t)$

- Martingale decomposition

$$\tilde{d}_t^\pi(s_t, a_t) - d_t^\pi(s_t, a_t) = \sum_{h=2}^t (\Gamma_{h+1:t} \pi_h(\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi)(s_t, a_t) + (\Gamma_{1:t}(\tilde{d}_1^\pi - d_1^\pi))(s_t, a_t),$$

- Bounded Martingale difference w.h.p

$$\sup_h \|\Gamma'_{h:t}(\tilde{T}_h - T_h)\|_\infty \leq O\left(\sqrt{\frac{1}{n \cdot d_m} \log \frac{H^2 S^2 A^2}{\delta}}\right).$$

- (Chung and Liu, 06)-style Azuma-Hoeffding, and union bound

$$\sup_t \|\tilde{d}_t^\pi - d_t^\pi\|_\infty \leq O\left(\sqrt{\frac{H}{n d_m} \log \frac{H^2 S^2 A^2}{\delta} \log \frac{H S A}{\delta}}\right).$$

# Completing the analysis

- Bounding  $W$

- Bounding  $M$

$$\sup_t |X_t - X_{t-1}| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot d_m}}\right).$$

- Apply (Chung and Liu's) Freedman's inequality

$$|X - \mathbb{E}[X]| \leq \sqrt{8\sigma^2 \cdot \log(1/\delta)} + \frac{2M}{3} \cdot \log(1/\delta),$$

# Uniform convergence theorem for all policies

**Theorem 3.3:** with probability  $\geq 1 - \delta$

$$\sup_{\pi \in \Pi} |\hat{v}^\pi - v^\pi| \lesssim \sqrt{\frac{H^4}{nd_m} \log\left(\frac{HSA}{\delta}\right)} + \sqrt{\frac{H^4 S}{nd_m} \log(SA)}$$

- Optimal in  $S$  if  $\delta < e^{-S}$ , suboptimal in  $H$ .
- Proof idea: Martingale decomposition over  $H$ . Freedman's inequality. Rademacher complexity argument.

# Uniform convergence theorem for near-empirically optimal policies

**Theorem 3.7:** Let  $\Pi_1 := \{\pi : s.t. \|\hat{V}_t^\pi - \hat{V}_t^{\hat{\pi}^*}\|_\infty \leq \epsilon_{opt}, \forall t \in [H]\}$ . Assume  $\epsilon_{opt} \leq \sqrt{H}/S$ , and also let  $n \gtrsim H^2/d_m$ . Then w.p.  $\geq 1 - \delta$ ,

$$\sup_{\pi \in \Pi_1} \left\| \hat{Q}_1^\pi - Q_1^\pi \right\|_\infty \leq c_2 \sqrt{\frac{H^3 \log(HSA/\delta)}{n \cdot d_m}}.$$

- Optimal in all parameters.
- Implies optimal learning bounds for ERM by taking  $\epsilon_{opt} = 0$
- Proof idea: A cute argument that takes the empirical optimal policy as an anchor point.



# Local uniform convergence is sufficient for

Assume generative model

## Offline Policy Learning

Simulation lemma (Kearns and Singh, 1998)	MSBO (Xie and Jiang, 2020)	Variance-Reduction (Sidford et al, 19), (Wainwright, 19)	Model-based (Agarwal, Kakade, Yang, 20)	Model-based via UniformOPE
$\sqrt{\frac{H^4 S^2}{n d_m}}$	$\sqrt{\frac{H^4}{n d_m}}$	$\sqrt{\frac{H^3 S A}{n}}$	$\sqrt{\frac{H^3 S A}{n}} + H \cdot \epsilon_{opt}$	$\sqrt{\frac{H^3}{n d_m}} + \epsilon_{opt}$

Converted from infinite horizon case...

# Proof sketch:

- Apply Simulation Lemma in a different way

$$\begin{aligned}\widehat{Q}_t^\pi - Q_t^\pi &= \sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi (\widehat{P}_h^\pi - P_h^\pi) \widehat{Q}_h^\pi \\ &= \sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi (\widehat{P}_h - P_h) \widehat{V}_h^\pi\end{aligned}$$

- Error decomposition

$$\begin{aligned}\left| \widehat{Q}_t^{\widehat{\pi}} - Q_t^{\widehat{\pi}} \right| &\leq \sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \left| (\widehat{P}_h - P_h) \widehat{V}_h^{\widehat{\pi}} \right| \\ &\leq \underbrace{\sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \left| (\widehat{P}_h - P_h) \widehat{V}_h^{\widehat{\pi}^*} \right|}_{(***)} + \underbrace{\sum_{h=t+1}^H \Gamma_{t+1:h-1}^\pi \left| (\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}}) \right|}_{(****)}\end{aligned}$$

# Bounding (\*\*\*\*)

$$\begin{aligned} \left\| \sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \cdot |(\hat{P}_h - P_h)(\hat{V}_h^{\hat{\pi}^*} - \hat{V}_h^{\hat{\pi}})| \right\|_{\infty} &\leq H \cdot \sup_h \left\| \Gamma_{t+1:h-1}^{\hat{\pi}} \right\|_{\infty} \left\| |(\hat{P}_h - P_h)(\hat{V}_h^{\hat{\pi}^*} - \hat{V}_h^{\hat{\pi}})| \right\|_{\infty} \\ &\leq H \cdot \sup_h \left\| |(\hat{P}_h - P_h)(\hat{V}_h^{\hat{\pi}^*} - \hat{V}_h^{\hat{\pi}})| \right\|_{\infty} \end{aligned}$$

$$\sup_h \left\| |(\hat{P}_h - P_h)(\hat{V}_h^{\hat{\pi}^*} - \hat{V}_h^{\hat{\pi}})| \right\|_{\infty} \leq \epsilon_{\text{opt}} \cdot \sup_h \left\| |\hat{P}_h - P_h| \cdot \mathbf{1} \right\|_{\infty}$$

- Apply the local-uniform assumption
- Apply L1-norm error bound (Recall from Lecture 3)

# Bounding (\*\*\*) $\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\hat{\pi}} \left| (\hat{P}_h - P_h) \hat{V}_h^{\hat{\pi}^*} \right|.$

- Key observation: Conditioning on  $D_{\{h-1\}}$

- $\hat{P}_h$  depends only on the data at step h
- $\hat{V}_h^{\hat{\pi}^*}$  depends only on the data after h

- Results in us saving a factor of S!

$$\left| (\hat{P}_h - P_h) \hat{V}_h^{\hat{\pi}^*} \right|_{(s_{t-1}, a_{t-1})} \leq 4 \sqrt{\frac{\log(1/\delta)}{N}} \sqrt{\text{Var}(\hat{V}_h^{\hat{\pi}^*})_{(s_{t-1}, a_{t-1})}} + \frac{4(H-t)}{3N} \log\left(\frac{1}{\delta}\right)$$

# Putting things together

$$\begin{aligned}
 \left| \widehat{Q}_t^{\widehat{\pi}} - Q_t^{\widehat{\pi}} \right| &\leq \sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\pi} \left| (\widehat{P}_h - P_h) \widehat{V}_h^{\widehat{\pi}} \right| \\
 &\leq \underbrace{\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\pi} \left| (\widehat{P}_h - P_h) \widehat{V}_h^{\widehat{\pi}^*} \right|}_{(***)} + \underbrace{\sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\pi} \left| (\widehat{P}_h - P_h) (\widehat{V}_h^{\widehat{\pi}^*} - \widehat{V}_h^{\widehat{\pi}}) \right|}_{(****)} \\
 &\leq \sum_{h=t+1}^H \Gamma_{t+1:h-1}^{\widehat{\pi}} \left( 4 \sqrt{\frac{\log(HSA/\delta)}{N}} \sqrt{\text{Var}(\widehat{V}_h^{\widehat{\pi}^*})} + \frac{4(H-t)}{3N} \log\left(\frac{HSA}{\delta}\right) \cdot \mathbf{1} \right) \\
 &\quad + c_1 \epsilon_{\text{opt}} \cdot \sqrt{\frac{H^2 S^2 \log(HSA/\delta)}{N}} \cdot \mathbf{1}
 \end{aligned}$$

- Bounding the sum of variance
  - Triangular inequality, then similar to that of Lemma 3.4

Lastly, backward recursion from  $H$  to  $1$ , using the following theorem

**Theorem F.4.** *Conditional on  $N > 0$ , then with probability  $1 - \delta$ , we have for all  $t = 1, \dots, H - 1$*

$$\begin{aligned} \left\| \widehat{Q}_t^{\widehat{\pi}} - Q_t^{\widehat{\pi}} \right\|_{\infty} &\leq 4\sqrt{\frac{H^3 \log(HSA/\delta)}{N}} + 4\sqrt{\frac{\log(HSA/\delta)}{N}} \sum_{h=t+1}^H \left\| \widehat{Q}_h^{\widehat{\pi}} - Q_h^{\widehat{\pi}} \right\|_{\infty} + \frac{4H^2}{3N} \log\left(\frac{HSA}{\delta}\right) \\ &\quad + c_2 \epsilon_{opt} \cdot \sqrt{\frac{H^2 S^2 \log(HSA/\delta)}{N}}. \end{aligned}$$

- **Assume:**  $N \geq 64H^2 \cdot \log(HSA/\delta)$  and  $\epsilon_{opt} \leq \sqrt{H}/S$ .

- We obtain the final result:

$$\left\| \widehat{Q}_1^{\widehat{\pi}} - Q_1^{\widehat{\pi}} \right\|_{\infty} \leq 2(9 + c_2) \sqrt{\frac{H^3 \log(HSA/\delta)}{N}}$$

# Summary

- We finished Offline RL for the tabular setting
- Via local uniform convergence, we showed that the model-based approach / ERM is minimax optimal
- We covered all essential tricks for doing a tight theoretical analysis

# What I did not cover

- Optimal offline RL in
  - Infinite horizon case
  - Finite horizon stationary case
  - (Yin, Bai, W., 2021) <https://arxiv.org/abs/2102.01748>
- Optimal rates in global uniform OPE (for model-based algorithms)
  - (Yin and W., 2021) <https://arxiv.org/abs/2105.06029>
- Offline RL with function approximation
  - Many algorithms that do policy learning without doing any sort of uniform OPE, e.g., FQI