

CS292F StatRL Lecture 14

From OPE to Uniform OPE

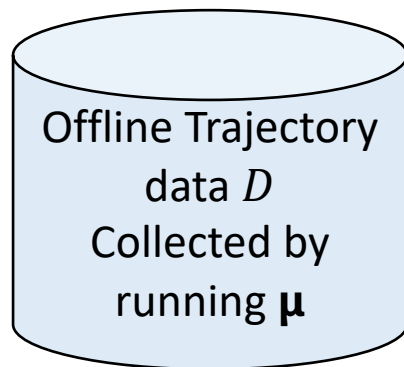
Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

Recap: Offline Reinforcement Learning, aka. Batch RL

- Task 1: Offline Policy Evaluation. (OPE)

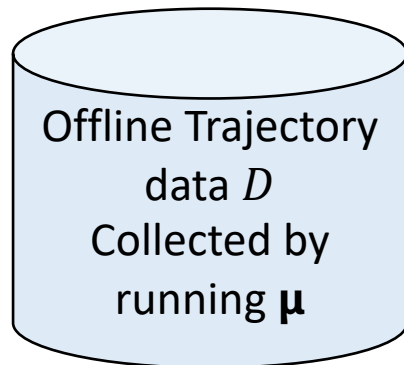


Task: design OPE methods

Evaluate fixed Target Policy π

Via
Uniform
OPE

- Task 2: Offline Policy Learning. (OPL)



Task: design OPO methods

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

Recap: Lecture 13

- Offline Policy Evaluation for RL

- Focused on Finite-Horizon Episodic MDPs
- Offline data collected by running logging / behavioral policy
- Evaluate a target policy $\pi = \{\pi_1, \dots, \pi_H\}$

$$\mathcal{M} = \{\mathcal{M}_1, \mathcal{M}_2, \dots, \mathcal{M}_H\}$$
$$\mathcal{M}_t = \mathcal{S}_t \rightarrow \Delta(\mathcal{A}_t)$$

- OPE RL estimators:

- Trajectory-IS, Per-Step IS, DR
- Curse of Horizon
- Overcoming the curse of horizon with MIS / TMIS
- Connections to DM --- a model-based approach.

$$\mathcal{J}(\pi(e^H)) \rightarrow \text{poly}(H)$$

Recap: Notations

$$P_1, \dots, P_{T+1}: S \times A \rightarrow \Delta(S)$$

$$Q_t^{\bar{i}}, Q_t^*$$

$$V_t^{\bar{i}}, V_t^*$$

- Finite Horizon Episodic MDP

$$(S, A, P, r, H, d)$$

- Offline Data by running π

$$(s_1^{(i)}, a_1^{(i)}, r_1^{(i)}) (s_2^{(i)}, a_2^{(i)}, r_2^{(i)}) \dots (s_H^{(i)}, a_H^{(i)}, r_H^{(i)}) \quad i = 1, 2, \dots, n$$

- Policies / importance weights

Target π :

$$\frac{\pi_t(a|s)}{\mu_t(a|s)} = p_t(s, a)$$

$$\frac{\pi_t(a_t^{(i)}|s_t^{(i)})}{\mu_t(a_t^{(i)}|s_t^{(i)})} = p_t^{(i)}$$

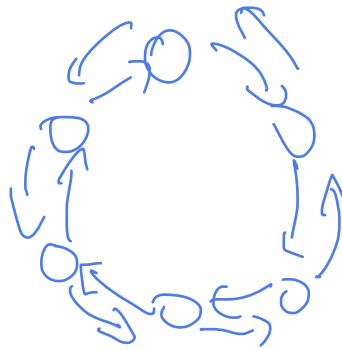
Recap: Main challenge of OPE in RL: The curse of Horizon

$$\hat{v}_{IS}^{\pi} = \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \left[\prod_{t=1}^h \frac{\pi(a_t^{(i)} | s_t^{(i)})}{\mu(a_t^{(i)} | s_t^{(i)})} \right] \underline{r}_h^{(i)}.$$

$\left(\max_{t, s, a} \frac{P(a|s)}{\mu(a|s)} \right)^h$

The curse of horizon. (Liu et al, 2018 NeurIPS)

- The variance is exponential in H!



Recap: Marginalized Importance Sampling

$$\hat{v}_{IS}^{\pi} = \frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H \left[\prod_{t=1}^h \frac{\pi(a_t^{(i)} | s_t^{(i)})}{\mu(a_t^{(i)} | s_t^{(i)})} \right] r_h^{(i)}.$$

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

IS at time t
reverse
 $P_t(s_t^{(i)} | a_t^{(i)})$

Xie, W., and Ma. (2019): Towards Optimal OPE for RL using Marginalized Importance Sampling. NeurIPS 2019.

Recap: Recursive estimation of State-visitation of target policy

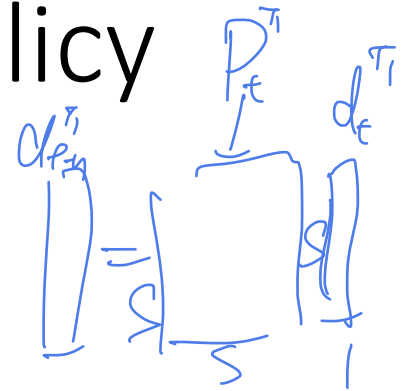
$$d_t^\pi(s_t) = \sum_{s_{t-1}} P_t^\pi(s_t | s_{t-1}) d_{t-1}^\pi(s_{t-1}),$$

$$\hat{d}_t^\pi = \hat{P}_t^\pi \hat{d}_{t-1}^\pi,$$

IS for estimating P_t^π

$$\text{where } \hat{P}_t^\pi(s_t | s_{t-1}) = \frac{1}{n_{s_{t-1}}} \sum_{i=1}^n \frac{\pi(a_{t-1}^{(i)} | s_{t-1})}{\mu(a_{t-1}^{(i)} | s_{t-1})} \mathbf{1}((s_{t-1}^{(i)}, s_t^{(i)}) = (s_{t-1}, s_t));$$

$$\hat{r}_t^\pi(s_t) = \frac{1}{n_{s_t}} \sum_{i=1}^n \frac{\pi(a_t^{(i)} | s_t)}{\mu(a_t^{(i)} | s_t)} r_t^{(i)} \mathbf{1}(s_t^{(i)} = s_t),$$



Xie, W., and Ma. (2019): Towards Optimal OPE for RL using Marginalized Importance Sampling. NeurIPS 2019.

Recap: Tabular MIS (finite action space)

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

- With a minor change to the following recursive estimation

$$\hat{P}_t^{\pi}(s_t | s_{t-1}) = \frac{1}{n_{s_{t-1}}} \sum_{i=1}^n \pi(a_{t-1}^{(i)} | s_{t-1}) \cdot \mathbf{1}((s_{t-1}^{(i)}, s_t^{(i)}, a_t^{(i)}) = (s_{t-1}, s_t, a_t));$$

$$\hat{r}_t^{\pi}(s_t) = \frac{1}{n_{s_t}} \sum_{i=1}^n \pi(a_t^{(i)} | s_t) \cdot \mathbf{1}(s_t^{(i)} = s_t).$$

Changed to Empirical Estimate.

Recap: OPE error bound of MIS and Tabular MIS

- The MSE of MIS estimator obeys:

$$\frac{1}{n} \sum_{t=1}^H \mathbb{E}_{\mu} \left[\frac{d_t^{\pi}(s_t)^2}{d_t^{\mu}(s_t)^2} \text{Var}_{\mu} \left[\frac{\pi_t(a_t|s_t)}{\mu_t(a_t|s_t)} (V_{t+1}^{\pi}(s_{t+1}) + r_t) \middle| s_t \right] \right] + \tilde{O}(n^{-1.5})$$

Xie, W., and Ma. (2019): Towards Optimal OPE for RL using Marginalized Importance Sampling. NeurIPS 2019.

- TMIS obeys:

Yin & W. (2020). Asymptotically efficient off-policy evaluation for tabular reinforcement learning. In AISTATS-2020

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

$$\sum_a \left(\frac{\pi(a)}{\mu(a)} \right)^2$$

Recap: Asymptotic rates for MIS and TMIS

- Assumptions: $\sum_{t=1}^H \tau_t$ is large $\implies \exists \forall t, s, a, \frac{\bar{\tau}_t(a|s)}{\mu_t(a|s)} \leq \bar{\tau}_a$

- MIS:
$$\frac{1}{n} \sum_{t=1}^H \mathbb{E}_\mu \left[\frac{d_t^\pi(s_t)^2}{d_t^\mu(s_t)^2} \text{Var}_\mu \left[\frac{\pi_t(a_t|s_t)}{\mu_t(a_t|s_t)} (V_{t+1}^\pi(s_{t+1}) + r_t) \middle| s_t \right] \right] + \tilde{O}(n^{-1.5})$$

$$\leq \frac{1}{n} \sum_{t=1}^H \tau_t \mathbb{E}_\mu \left[\frac{\text{Var}_\mu [V_{t+1}^\pi(s_{t+1}) + r_t]}{\mu_t(a_t|s_t)} \right] \leq \frac{1}{n} \tau_s \tau_a \leq O\left(\frac{1}{n} \tau_s^3 \tau_a\right) \leq O\left(\frac{1}{n}\right)$$

- TMIS
$$\frac{1}{n} \sum_{h=0}^H \sum_{s_h, a_h} \frac{d_h^\pi(s_h)^2}{d_h^\mu(s_h)^2} \frac{\pi(a_h|s_h)^2}{\mu(a_h|s_h)} \text{Var} \left[(V_{h+1}^\pi(s_{h+1}^{(1)}) + r_h^{(1)}) \middle| s_h^{(1)} = s_h, a_h^{(1)} = a_h \right] + O(n^{-1.5})$$

$$\leq \frac{1}{n} \tau_s \tau_a \left[\sum_{h=1}^H \text{Var} [V_{h+1}^\pi(s_{h+1}^{(1)}) + r_h^{(1)} / s_h, a_h] \right] \leq O\left(\frac{1}{n}\right)$$

Recap: TMIS vs on-policy evaluation

Lemma 3.4. For any policy π and any MDP.

$$H^2 \geq \text{Var}_\pi \left[\sum_{t=1}^H r_t^{(1)} \right] = \sum_{t=1}^H \left(\mathbb{E}_\pi \left[\text{Var} \left[r_t^{(1)} + V_{t+1}^\pi(s_{t+1}^{(1)}) \mid s_t^{(1)}, a_t^{(1)} \right] \right] + \mathbb{E}_\pi \left[\text{Var} \left[\mathbb{E} \left[r_t^{(1)} + V_{t+1}^\pi(s_{t+1}^{(1)}) \mid s_t^{(1)}, a_t^{(1)} \right] \mid s_t^{(1)} \right] \right] \right).$$

≥ 0

Combined with the previous observation:

1. TMIS has an error that is linear in H.

Absolute error $\approx H \sqrt{\frac{1}{n}}$

2. TMIS is better than MC even when we are doing on-policy evaluation

How do MIS / TMIS overcome the curse of horizon?

IS DR

- Leverage the model assumption: MDP
- We visit the same state many times.
- As a matter of fact, TMIS is equivalent to a model-based approach

Recap: TMIS is equivalent to DM -
 -- a model-based approach

$$\hat{m} = (S, A, \hat{p}, \hat{\gamma}, \hat{H}, \hat{d}_1)$$

$$\hat{V}_{MIS} = \frac{1}{n} \sum_{i=1}^n \frac{H}{t+1} \frac{d_t^{\tau}(s_t^{(i)})}{d_t^{\tau}(s_t^{(i)})} \gamma_{\tau}(s_t^{(i)}) = \sum_{S \in \mathcal{S}} \frac{H}{t+1} d_t^{\tau}(s) \gamma_{\tau}(s) = \hat{V}_{DM}$$

$$d_t^{\tau}(s) = \sum_{s'} \mathbb{P}_{\tau}(s'|s) d_{t-1}^{\tau}(s')$$

$$\mathbb{I} = \square \square$$

This lecture

- Fitted Q Iterations
- Uniform Convergence in OPE in RL

Fitted Q Iterations (Munos, 2003) (Munos & Szepesvari, 2008)

- Recall Bellman Optimality equation and the Bellman operator

$$\mathcal{T}f(s, a) := r(s, a) + \mathbb{E}_{s' \sim P(\cdot | s, a)} \max_{a' \in \mathcal{A}} f(s', a').$$



$$Q^*(s, a) = r(s, a) + \sum_{s'} P(s' | s, a) \max_{a'} Q^*(s', a')$$

- Given offline transition data and a function class

FQI: $f_t \in \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \left(f(s_i^{\circ}, a_i) - r_i - \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s_i^{\circ}, a_i) \right)^2$.

Iteratively from some initialization.

Handwritten notes: $f(s_i, a_i) = \theta^T \phi(s_i, a_i)$ and t (with a scribble).

- For the finite horizon episodic case:

$f_{t+1} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n (f(s_{H,i}^{(i)}, a_{H,i}^{(i)}) - r_{H,i}^{(i)})^2$

$f_{t+1} = \operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^n \left(f(s_{H,i}^{(i)}, a_{H,i}^{(i)}) - r_{H,i}^{(i)} - \max_{a' \in \mathcal{A}} f_{t-1}(s_{H,i}^{(i)}, a') \right)^2$

Handwritten notes: $\min_{f \in \mathcal{F}} \sum_{i=1}^n (f(s_i, a_i) - r_i - \gamma \max_{a' \in \mathcal{A}} f_{t-1}(s_i, a'))^2$

Fitted Q iterations for OPE (Duan and Wang, 2020)

- Recall Bellman equation for a fixed policy

$$Q_{h-1}^\pi(s, a) = r(s, a) + \mathbb{E}[V_h^\pi(s') \mid s, a]$$

$\mathbb{E}[\mathbb{E}_{\pi(a'|s)}[Q_h^\pi(s', a') \mid s'] \mid s, a]$

- Given offline transition data and a function class

$$\hat{Q}_{H+1}^\pi := 0 \text{ and for } h = H, H - 1, \dots, 0,$$

$$\hat{Q}_h^\pi = \arg \min_{f_h \in \mathcal{F}} \sum_{i=1}^n \left(f_h(s_h^{(i)}, a_h^{(i)}) - r_h^{(i)} - \sum_{a' \in \mathcal{A}} \pi(a' | s_{h+1}^{(i)}) f_{h+1}(s_{h+1}^{(i)}, a') \right)^2$$

$\forall h = H, \dots, 0$
 for $h < H$, take Q_{h+1}^π

$\theta \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$

FQI in the tabular case

$$\hat{Q}_h^\pi = \arg \min_{f_h \in \mathcal{F}} \sum_{i=1}^n \left(f_h(s_h^{(i)}, a_h^{(i)}) - r_h^{(i)} - \sum_{a' \in \mathcal{A}} \pi(a' | s_{h+1}^{(i)}) f_{h+1}(s_{h+1}^{(i)}, a') \right)^2$$

- Let's work out the optimal solution!

$f_h(s_h^{(i)}, a_h^{(i)}) = \theta \cdot \mathbb{1}(s_h^{(i)}, a_h^{(i)}) = \theta_{s_h^{(i)}, a_h^{(i)}}$

$\frac{\partial}{\partial \theta_{s,a}} \sum_{i=1}^n \left(\theta_{s_h^{(i)}, a_h^{(i)}} - r_h^{(i)} - \sum_{a' \in \mathcal{A}} \pi(a' | s_{h+1}^{(i)}) \theta_{s_{h+1}^{(i)}, a'} \right) \cdot \mathbb{1}(s_h^{(i)}, a_h^{(i)} = s, a)$

$\sum_{i=1}^n \mathbb{1}(s_h^{(i)}, a_h^{(i)} = s, a)$
 # of times s, a

$\theta_{s,a} = \frac{\sum_{i=1}^n (r_h^{(i)} + \sum_{a'} \pi(a' | s_{h+1}^{(i)}) \theta_{s_{h+1}^{(i)}, a'}) \cdot \mathbb{1}(s_h^{(i)}, a_h^{(i)} = s, a)}{\sum_{i=1}^n \mathbb{1}(s_h^{(i)}, a_h^{(i)} = s, a)}$

In conclusion, in the tabular MDP case, ~~they are all equivalent.~~ **they are all equivalent.**

$h=1, \dots, H$

- TMIS

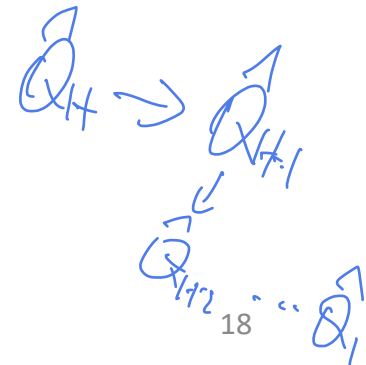
$$\hat{v}_{\text{MIS}}^{\pi} = \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^H \frac{\hat{d}_t^{\pi}(s_t^{(i)})}{\hat{d}_t^{\mu}(s_t^{(i)})} \hat{r}_t^{\pi}(s^{(i)}).$$

- Model-based Plugin

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

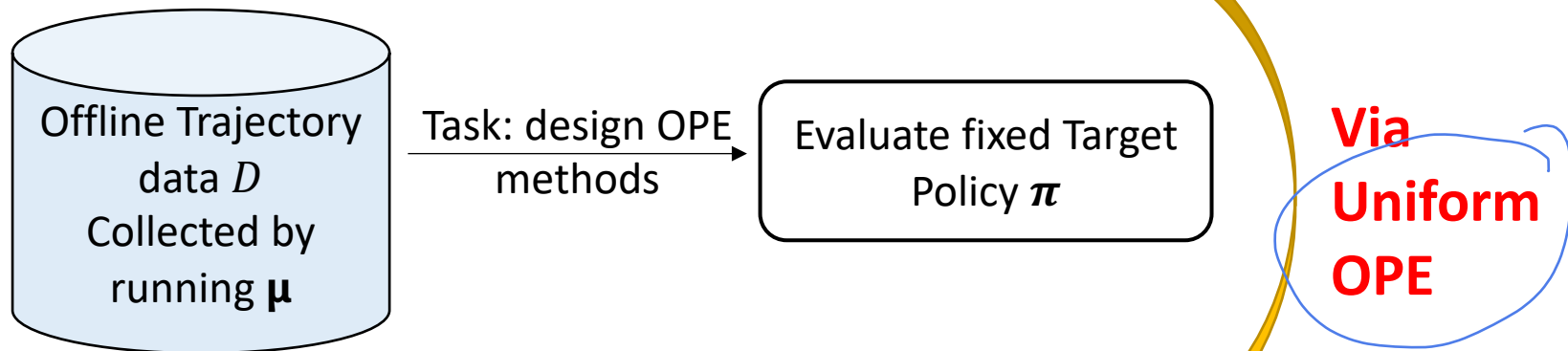
- Fitted Q Iteration

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

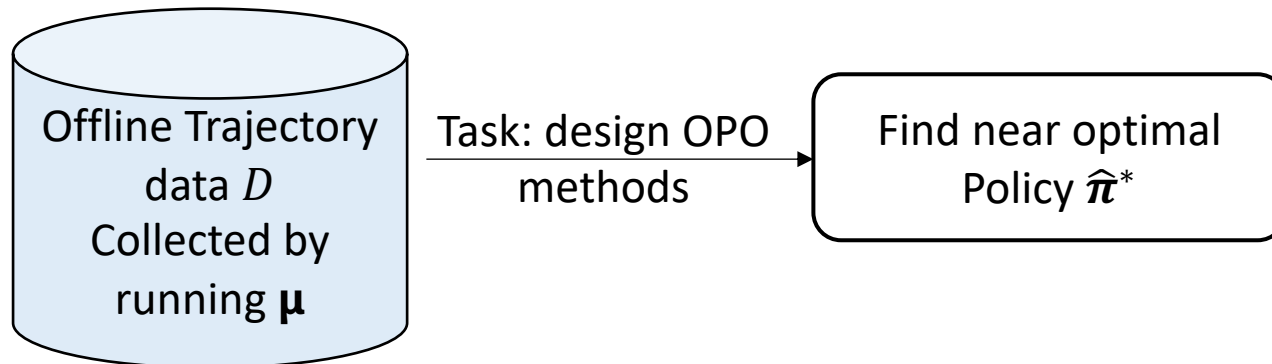


Recap: Offline Reinforcement Learning, aka. Batch RL

- Task 1: Offline Policy Evaluation. (OPE)



- Task 2: Offline Policy Learning. (OPL)



Observation 1: OPE is in its essence a statistical estimation problem.

- But is slightly non-trivial because we are estimating a single number, when the number of parameters describing the distribution are numerous.

$$O(H_{SA}^2 + H_{SA})$$

- Find functions of the data --- estimators, such that

$$|\hat{v}^\pi - v^\pi| \leq \epsilon \quad \text{with high probability}$$

$$\mathbb{E} [|\hat{v}^\pi - v^\pi|^2] \leq \epsilon^2$$

Observation 2: Offline Learning is a statistical learning problem

- But with a structured hypothesis class (the policy class), and structured observations (trajectories).

- Lessons from statistical learning theory:

- ERM suffices and almost necessary.

- In RL context this is: $\hat{\pi} = \arg \max_{\pi \in \Pi} \hat{v}^{\pi}$

(For some estimator \hat{v}^{π})

- Combine with OPE:

$$|\hat{v}^{\pi} - v^{\pi}| \leq \epsilon \text{ w.h.p.}$$

$$\mathbb{E}[|\hat{v}^{\pi} - v^{\pi}|^2] \leq \epsilon^2$$



$$v^{\pi^*} - v^{\hat{\pi}} \leq 2\epsilon \text{ w.h.p.}$$

$$v^{\pi^*} - \mathbb{E}[v^{\hat{\pi}}] \leq 2\epsilon$$

Not quite this easy, the learned policy $\hat{\pi}$ depends on the data

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



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

In standard statistical learning: $\epsilon \asymp \sqrt{d/n}$

Where d is VC-dimension / metric entropy = $\log(\text{Covering \#})$
 $\log|\Pi|$, or implied by Rademacher complexity, etc.

(Much older Empirical process theory , Glivenko-Cantelli style)



Vapnik (1995)

What is a natural complexity measure for the policy class in RL?

We will *not* deal with exploration in offline RL, because we can't

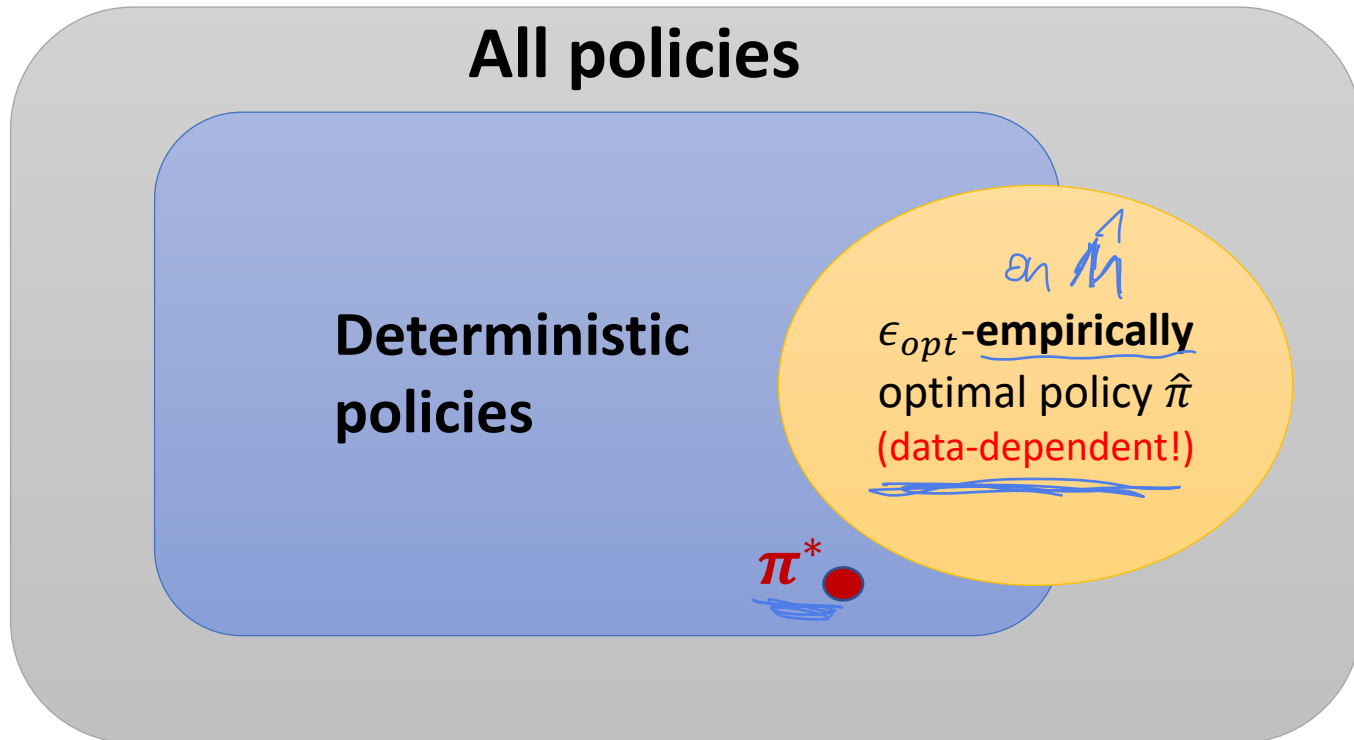
- The logging policy μ is out of our control
- Need to make assumptions about it

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

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

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

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>

Start with the family of all deterministic policies

- The optimal policy is deterministic
- There are a finite number of them
- We have strong pointwise convergence bound from TMIS (last week)

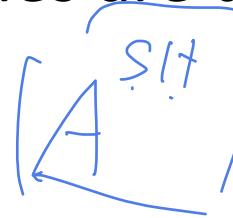
OPE with fixed π

Exercise: counting the number of deterministic policies

- Setting: Tabular MDP with S states, A actions and H steps.

- How many deterministic ^{memoryless} policies are there?

$$\pi = (\pi_1, \dots, \pi_H)$$



$$\log A^{SH} = \underline{\underline{SH \log A}}$$

We need a **high probability version** of the bound we get

Theorem (MSE bound)

(Yin & W., 2020)

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

What we need is something stronger:

$$P\left(\left|\widehat{v}_{\text{TMSIS}}^\pi - v^\pi\right| > \varepsilon\right) \leq e^{-\frac{n\varepsilon^2}{2}}$$

$$\varepsilon = \frac{1}{\sqrt{n\alpha}} \sqrt{\log \frac{A^{SH}}{S}}$$

Fictitious estimator technique

- Fictitious estimator

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

- Define

$$\tilde{r}_t(s_t, a_t) = \hat{r}_t(s_t, a_t)\mathbf{1}(E_t) + r_t(s_t, a_t)\mathbf{1}(E_t^c)$$

$$\tilde{P}_{t+1}(\cdot|s_t, a_t) = \hat{P}_{t+1}(\cdot|s_t, a_t)\mathbf{1}(E_t) + P_{t+1}(\cdot|s_t, a_t)\mathbf{1}(E_t^c).$$

Idea: *hypothetically* plug in the ground truth occasionally

$$\tilde{P}_t^\pi(s_t|s_{t-1}) = \sum_{a_{t-1}} \tilde{P}_t(s_t|s_{t-1}, a_{t-1})\pi(a_{t-1}|s_{t-1}).$$

$$\tilde{v}^\pi := \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t^\pi \rangle, \text{ with } \tilde{d}_t^\pi = \tilde{P}_t^\pi \tilde{d}_{t-1}^\pi$$

multiplication
 change 1/2

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

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

Handwritten notes:

$$\sup_a (v^a - v^{\bar{a}})$$

$$\leq \sup_a (v_{\bar{a}}^a - v_{\bar{a}}^{\bar{a}})$$

$$\leq \sup_a (v_{\bar{a}}^a - v_{\bar{a}}^{\bar{a}})$$

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

Handwritten note: $d_m \leq d_{t,CS,a} \forall S,a,t$

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**.

Dealing with the reward noise

Lemma C.4. *We have with probability $1 - \delta$:*

$$\sup_{\pi \in \Pi} \left| \sum_{t=1}^H \langle \tilde{d}_t^\pi, \tilde{r}_t - r_t \rangle \right| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{n \cdot d_m}}\right)$$

$\langle \tilde{d}_t, \tilde{r}_t - r_t \rangle \leq \|\tilde{d}_t\|_1 \|\tilde{r}_t - r_t\|_\infty$
 w.h.p. $\frac{n}{tSA} > n \cdot d_m \cdot \frac{1}{\epsilon}$
 $\|\tilde{r}_t - r_t\|_\infty \leq \frac{\log(HSA/\delta)}{n \cdot d_m}$
 (Dudley's) $\forall s_{t-1} < s_t \quad \tilde{r}_t(s_t) - \tilde{r}_t(s_{t-1}) \leq \frac{\log(HSA/\delta)}{n \cdot d_m}$
 w.p. $1 - \delta$

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):

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

Let's derive the Martingale Decomposition

$$I = \square I$$

$$\tilde{d}_t^\pi = \pi_t \tilde{T}_t \tilde{d}_{t-1}^\pi$$

$$T_t \in \mathbb{R}^{S \times (S \cdot A)}$$

$$d_t^\pi = \pi_t T_t d_{t-1}^\pi$$

$$(T_t)_{s_t, (s_{t-1}, a_{t-1})} = P_t(s_t | s_{t-1}, a_{t-1})$$

Take the difference of the two

$$\tilde{d}_t^\pi - d_t^\pi = \pi_t (\tilde{T}_t - T_t) \tilde{d}_{t-1}^\pi + \pi_t T_t (\tilde{d}_{t-1}^\pi - d_{t-1}^\pi)$$

Recursively apply the above

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

Dual representation (Value function style):

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

$$\begin{aligned}
 X &= \sum_{t=1}^H \left(\sum_{h=2}^t \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle + \langle r_t, \Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi) \rangle \right) \\
 &= \sum_{t=1}^H \left(\sum_{h=2}^t \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \right) + \sum_{h=1}^H \langle r_t, \Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi) \rangle \\
 &= \sum_{t=2}^H \left(\sum_{h=2}^t \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \right) + \sum_{h=1}^H \langle r_t, \Gamma_{1:t} (\tilde{d}_1^\pi - d_1^\pi) \rangle \\
 &= \sum_{h=2}^H \left(\sum_{t=h}^H \langle r_t, \Gamma_{h+1:t} \pi_h (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \right) + \sum_{h=1}^H \langle (\pi_1^T \Gamma_{1:t}^T r_t)(s), (\tilde{d}_1^\pi - d_1^\pi)(s) \rangle \\
 &= \sum_{h=2}^H \left(\underbrace{\left\langle \sum_{t=h}^H \pi_h^T \Gamma_{h+1:t}^T r_t, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \right\rangle}_{V_h^\pi(s)} \right) + \underbrace{\left\langle \left(\sum_{h=1}^H \pi_1^T \Gamma_{1:t}^T r_t \right)(s), (\tilde{d}_1^\pi - d_1^\pi)(s) \right\rangle}_{V_1^\pi(s)}
 \end{aligned}$$

Let's check that this is a Martingale

$$X_t := \mathbb{E}[X | \mathcal{D}_t] = \sum_{h=2} \langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle + \langle V_1^\pi, \tilde{d}_1^\pi - d_1^\pi \rangle.$$

$$\mathbb{E}[X | \mathcal{D}_t] = \sum_{h=t+1}^H \mathbb{E} \left[\langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_t \right] + \sum_{h=2}^t \langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle + \langle V_1^\pi, (\tilde{d}_1^\pi - d_1^\pi) \rangle.$$

Note for $h \geq t + 1$, $\mathcal{D}_t \subset \mathcal{D}_{h-1}$, so by total law of expectation (tower property) we have

$$\begin{aligned} & \mathbb{E} \left[\langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_t \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\langle V_h^\pi, (\tilde{T}_h - T_h) \tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_{h-1} \right] \middle| \mathcal{D}_t \right] \\ &= \mathbb{E} \left[\langle V_h^\pi, \mathbb{E} \left[(\tilde{T}_h - T_h) \middle| \mathcal{D}_{h-1} \right] \tilde{d}_{h-1}^\pi \rangle \middle| \mathcal{D}_t \right] = 0 \end{aligned}$$

Final results after applying a complex martingale concentration

- And a union bound.

Theorem E.6. *With probability $1 - \delta$, we have*

$$\left| \sum_{t=1}^H \langle \tilde{d}_t^\pi - d_t^\pi, r_t \rangle \right| \leq O\left(\sqrt{\frac{H^2 \log(HSA/\delta)}{nd_m}} + \sqrt{\frac{H^4 SA \cdot \log(H^2 S^2 A^2 / \delta) \log(HSA/\delta)}{n^2 d_m^2}} \right)$$

where $O(\cdot)$ absorbs only the absolute constants.

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.

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.