

# CS292F StatRL Lecture 13

## OPE in Reinforcement Learning

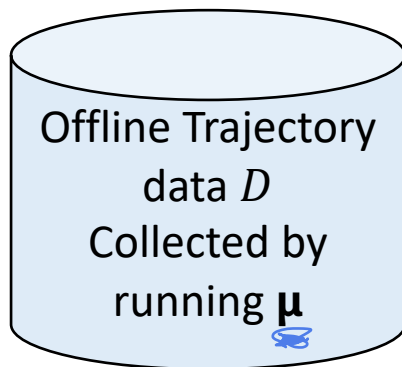
Instructor: Yu-Xiang Wang

Spring 2021

UC Santa Barbara

# Recap: Offline Reinforcement Learning, aka. Batch RL

- Task 1: Offline Policy Evaluation. (OPE)

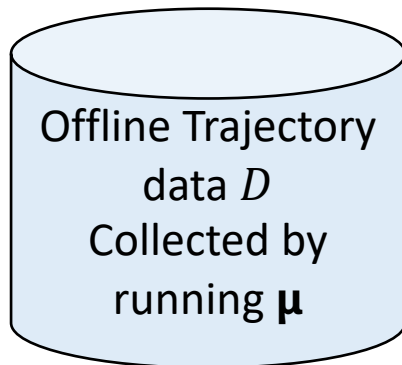


Task: design OPE methods

Evaluate fixed Target Policy  $\pi$

Via  
Uniform  
OPE

- Task 2: Offline Policy Learning. (OPL)



Task: design OPO methods

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

# Recap: Lecture 12

- OPE algorithms in (Contextual) Bandits
  - DM, IS, WIS, DR, SWITCH
- Comparing DM and IS in Multi-armed Bandits:
  - DM is asymptotically more efficient
  - IS is better.
- More generally:
  - DM is asymptotically more efficient if we assume realizability
  - IS cannot be improved when we don't

$$f(s) \rightarrow \mathbb{R}$$

$$\frac{\pi(a|s)}{u(a|s)}$$

$$\frac{\pi}{u}$$

# Recap: Two standard approaches

- Direct method / regression estimator

$$\hat{v}_{\text{DM}}^{\pi} = \frac{1}{n} \sum_{i=1}^n \sum_{a \in \mathcal{A}} \hat{r}(x_i, a) \pi(a|x_i)$$

- Importance sampling / Inverse Propensity Score /

$$\hat{v}_{\text{IPS}}^{\pi} = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i|x_i)}{\mu(a_i|x_i)} r_i$$



# Recap: Combining DM and IS

- Doubly Robust Estimation
  - Remains unbiased, but limited benefits to the variance
- SWITCH
  - Introduce bias, but drastically reduce variance

# This lecture

- Generalizing the bandits OPE idea to RL

Finite ~~Horizon~~ Horizon MDPs

finite state

finite Action

- Curse of Horizon

- Marginalized Importance Sampling

MDP:  $(S, A, P, d, r)$

$n$  trajectories i.i.d  
 $i$ -th:  $(S_1^{(i)}, a_1^{(i)}, r_1^{(i)}), (S_2^{(i)}, a_2^{(i)}, r_2^{(i)}) \dots$   
 $S^{(i)} \sim d(s)$   
 $a_h^{(i)} \sim \pi(a_h | s_h^{(i)})$   
 $r_h^{(i)} \sim p(r_h | s_h, a_h^{(i)})$

# OPE in Reinforcement Learning

- Importance sampling on the entire trajectory

$$V_{\text{OPE}}^{\pi} = \frac{1}{n} \sum_{i=1}^n \left( \sum_{h=1}^H r_h^{(i)} \right) \cdot \frac{\pi_{1:H}(S_1^{(i)}, a_1^{(i)}, \dots, a_H^{(i)})}{\mu_{1:H}(S_1^{(i)}, a_1^{(i)}, \dots, a_H^{(i)})}$$

- (Per-Step) Importance Sampling

$$V_{\text{IS}}^{\pi} = \sum_{h=1}^H \frac{1}{n} \sum_{i=1}^n \left( \frac{1}{\pi} p_h^{(i)} \right) \cdot r_h^{(i)}$$

$$\frac{1}{n} \sum_{i=1}^n \sum_{h=1}^H p_h^{(i)} \sum r_h^{(i)}$$

when  $p_h^{(i)} = \frac{\pi(a_h^{(i)} | s_h^{(i)})}{\mu(a_h^{(i)} | s_h^{(i)})}$

- Exercise:

- Infinite horizon discounted version?
- Weighted Importance Sampling Extension?

# Doubly Robust OPE in Reinforcement Learning

- An alternative form for the Per-Step IS

$$V_{\text{step-IS}}^0 := 0, \text{ and for } t = 1, \dots, H,$$

$$V_{\text{step-IS}}^{H+1-t} := \rho_t \left( r_t + \gamma V_{\text{step-IS}}^{H-t} \right).$$

$\gamma=1$   
 $H=2$

$V_{\text{IS}} = \cancel{P_2 r_2 + P_1 r_1} = P_1 r_1 + P_2 (P_2 r_2)$

- Given a value function approximator

$$V_{\text{DR}}^0 := 0, \text{ and for } t = 1, \dots, H,$$

$$V_{\text{DR}}^{H+1-t} := \hat{V}(s_t) + \rho_t \left( r_t + \gamma V_{\text{DR}}^{H-t} - \hat{Q}(s_t, a_t) \right).$$

$\hat{V}(s_t) = \sum \pi(a_t|s_t) \hat{Q}(s_t, a_t)$

# Mean and Variance of Doubly Robust OPE in RL

- Doubly Robust OPE in RL is unbiased
- Variance

**Theorem 1.**  $V_{DR}$  is an unbiased estimator of  $v^{\pi_1, H}$ , whose variance is given recursively as follows:  $\forall t = 1, \dots, H$ ,

$$\begin{aligned} \mathbb{V}_t [V_{DR}^{H+1-t}] &= \mathbb{V}_t [V(s_t)] + \mathbb{E}_t \left[ \mathbb{V}_t [\rho_t \Delta(s_t, a_t) \mid s_t] \right] \\ &+ \mathbb{E}_t \left[ \rho_t^2 \mathbb{V}_{t+1} [r_t] \right] + \mathbb{E}_t \left[ \gamma^2 \rho_t^2 \mathbb{V}_{t+1} [V_{DR}^{H-t}] \right], \quad (11) \end{aligned}$$

where  $\Delta(s_t, a_t) := \hat{Q}(s_t, a_t) - Q(s_t, a_t)$  and  $\mathbb{V}_{H+1} [V_{DR}^0 \mid s_H, a_H] = 0$ .

$$\mathbb{V}(V_{DR}) = \sum_{h=1}^H \sum_{\pi} (\pi P_t)^2 \mathbb{V}_{\pi} (V(s') + r \mid s, a)$$

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

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

*iid trajectory*

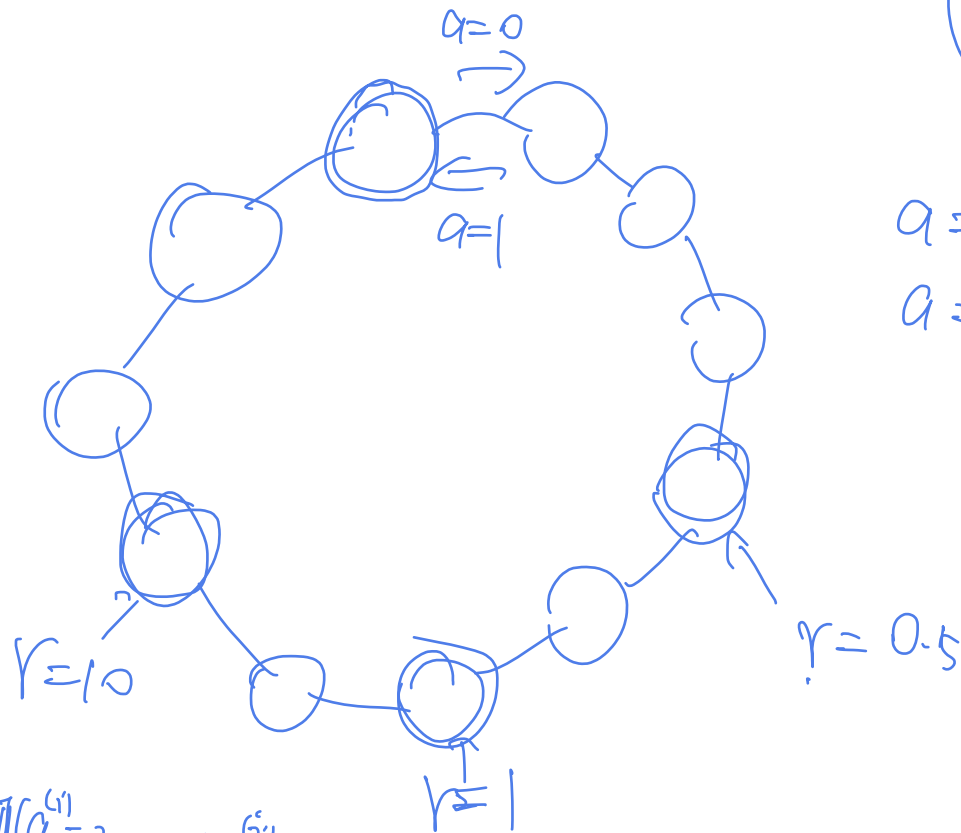
The curse of horizon. (Liu et al, 2018 NeurIPS)

- The variance is exponential in H!

$$\text{Var}(\hat{V}_S^{\pi}) = \frac{1}{n} \text{Var}(V)$$

# Example on Curse of Horizon

(Liu et al. 2018)



$a=0$   
 $a=1$

~~$u(s)$~~   
 $u(s) = \begin{cases} 0 & \text{u.p.p} \\ 1 & \text{u.p.l-p} \end{cases}$

~~$\tau(s)$~~   
 $\tau(s) = \begin{cases} 0 & \text{u.p.p} \\ 1 & \text{u.p.l-p} \end{cases}$

$H = 200$

$$\Delta \pi$$

$$V_{15} = \frac{1}{n} \sum_{i=1}^n \sum_{t=1}^H \left(\frac{1}{p}\right)^t \cdot \mathbb{1}(a_t=0 \text{ for } t \leq h) \cdot \gamma^{(i)}$$

$$\frac{h}{n} \frac{\tau(a(c))}{u(s)}$$

$$\text{Var}(V_{15}^{(i)}) = O\left(\frac{1}{p}\right)^x$$

# From Importance Sampling to Marginalized Importance Sampling

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

$\hat{d}_t^{\bar{u}}(s) \rightarrow \underline{d}_t^{\bar{u}}(s)$       $\hat{r}_t^{\bar{u}}(s) \rightarrow \underline{r}_t^{\bar{u}}(s)$   
 $\hat{d}_t^u(s) \rightarrow \underline{d}_t^u(s)$       $\hat{r}_t^u(s) \rightarrow \underline{r}_t^u(s)$

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

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



# What are some ideas for estimating the marginalized importance weight?

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

Handwritten notes:

$$\hat{r}_t^{\pi}(s) = \frac{1}{n_{s,t}} \sum_{i=1}^n \mathbb{1}(s_t^{(i)}=s) \cdot r_t^{(i)} \cdot \frac{\pi(a_t^{(i)})}{\mu(a_t^{(i)})}$$

- Idea 1: averaging over multiple visits to the same state.

Handwritten note (crossed out):

$$\hat{d}_t^{\mu}(s) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}(s_t^{(i)}=s)$$

Handwritten note:

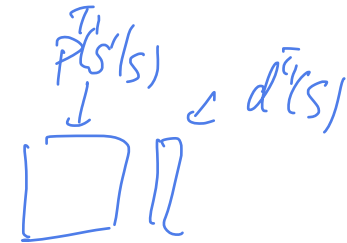
$$\frac{\hat{d}_t^{\pi}(s)}{\hat{d}_t^{\mu}(s)} = \frac{1}{n_{s,t}} \sum_{i=1}^n \frac{r_t^{(i)}}{\gamma} \underbrace{P_h^{(i)}(s_{n+1}^{(i)}, a_t^{(i)})}_{\frac{\mu(s_{n+1})}{\mu(s_t)}} \cdot \mathbb{1}(s_t^{(i)}=s)$$

Handwritten note:  $\frac{1}{n_{s,t}}$

# Idea 2: Recursive estimation

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

add  
dt



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

$$d_t^\pi = \frac{1}{n} \sum \mathbb{1}(s_t^{(i)} = s)$$

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

use IS

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

use IS

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

# Results: OPE error bound of MIS

- The MSE of MIS estimator obeys:

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

Handwritten annotations:  $\leq \tau_s^2$  above the first fraction,  $\tau_a^2$  above the second fraction,  $\leq o(H)$  below the variance term, and  $\sum_{s_t \in S} d_t^{\mu}(s_t) \cdot d_t^{\pi}(s_t)$  below the sum.

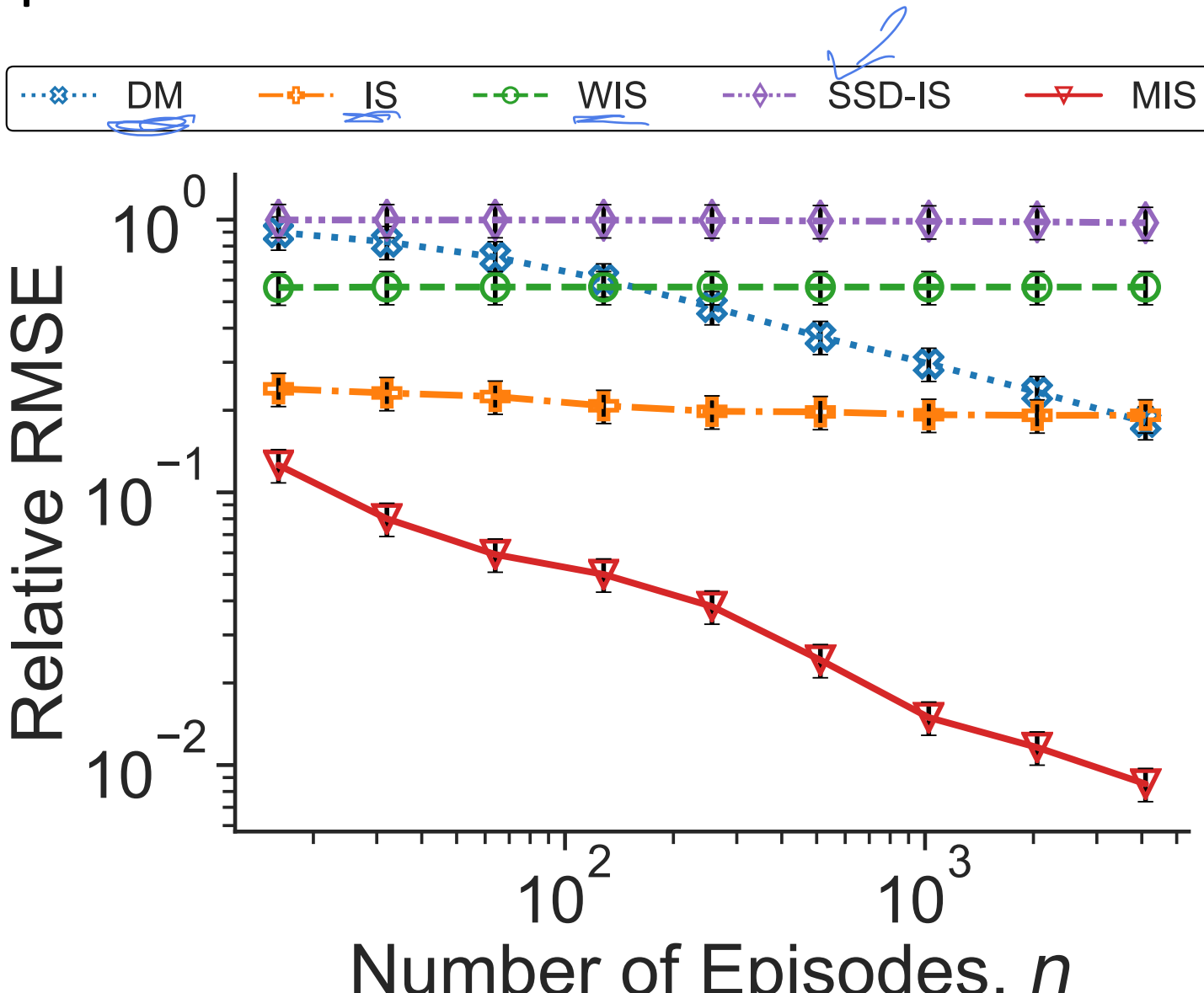
$$\frac{1}{n} H \tau_s^2 \tau_a^2 H^2 = \frac{H^3}{n} \tau_s^2 \tau_a^2$$

$$\begin{aligned} \max_{s,t} \left| \frac{d_t^{\pi}(s)}{d_t^{\mu}(s)} \right| &\leq \tau_s \\ \max_{s,a,t} \frac{\pi_t(a|s)}{\mu_t(a|s)} &\leq \tau_a \end{aligned}$$

Handwritten notes:  $0 \leq \pi \leq 1, \text{Var} r \leq 1$

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

# Experiment on mountain car



# Challenges of the analysis

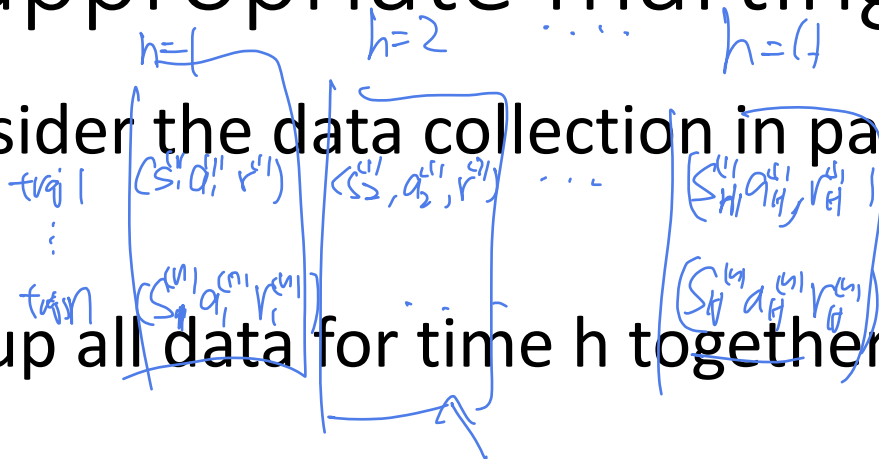
- 1. Dependent data:** The data within each trajectory are not independent
- 2. An annoying bias:** there is a non-zero probability that some states are not visited at all. And it affects all future estimates
- 3. Error propagation** from recursive estimation

$$\underline{d_t} = \underline{P_t} \underline{d_{t+1}}$$

$\uparrow$  Error       $\uparrow$  Error

# Addressing Challenge 1: Define an appropriate martingale

- Consider the data collection in parallel



- Group all data for time  $h$  together



- Conditioning on the number of times states are visited

# Addressing Challenge 2: Fictitious estimator technique

$$\underline{\tilde{v}^\pi} := \sum_{t=1}^H \sum_{s_t} \tilde{d}_t^\pi(s_t) \tilde{r}_t^\pi(s_t) \quad \text{where} \quad \tilde{d}_t^\pi = \tilde{\mathbb{P}}_{t,t-1}^\pi \tilde{d}_{t-1}^\pi$$

$$\tilde{r}_t^\pi(s_t) = \begin{cases} \hat{r}_t^\pi(s_t) & \text{if } n_{s_t} \geq n d_t^\mu(s_t)(1 - \delta) \\ \underline{r_t^\pi(s_t)} & \text{otherwise;} \end{cases}$$

true value

$$\tilde{\mathbb{P}}_{t,t-1}^\pi(\cdot | s_{t-1}) = \begin{cases} \hat{\mathbb{P}}_{t,t-1}^\pi & \text{if } n_{s_{t-1}} \geq n d_t^\mu(s_{t-1})(1 - \delta) \\ \mathbb{P}_{t,t-1}^\pi & \text{otherwise.} \end{cases}$$

$$E(\tilde{v}^\pi) = v^\pi$$

# Multiplicative Chernoff Bound

**Lemma A.1** (Multiplicative Chernoff bound [Chernoff et al., 1952]).

Let  $X$  be a Binomial random variable with parameter  $p, n$ .

For any  $\delta > 0$ , we have that

$$\mathbb{P}[X < (1 - \delta)pn] < e^{-\frac{\delta^2 pn}{2}}$$

- Apply to our problem

$n_{S,t}$  by running  $\mathcal{M}$  for  $n$  trials.

$\sim \text{Bin}(d_t^{\text{true}}(s), n)$

$$P(n_{S,t} < \underline{nd_t^{\text{true}}(s)}(1-\delta)) \leq e^{-\frac{\delta^2 n \cdot d_t^{\text{true}}(s)}{2}}$$

Union bound for all  $t=1, \dots, T$ , all  $S \in \mathcal{S}$

$\frac{1}{H \cdot S}$

~~the tip~~  
 $P(\text{Fictitious Estimator} \neq M(S))$   
 $\leq H \cdot S \cdot e^{-\frac{\delta^2 n \cdot \min_{s \in \mathcal{S}} d_t^{\text{true}}(s)}{2}}$   
 $\frac{\delta = \tilde{O}\left(\frac{1}{\sqrt{n}}\right)}{h \cdot \frac{1}{\sqrt{n}}}$



# Address Challenge 3: Empirical / Offline version of Bellman equation of variance

$$\text{Var}[\tilde{v}^\pi] = \sum_{h=0}^H \sum_{s_h} \mathbb{E} \left[ \frac{\tilde{d}_h^\pi(s_h)^2}{n_{s_h}} \mathbf{1} \left( n_{s_h} \geq \frac{n d_h^\mu(s_h)}{(1-\delta)^{-1}} \right) \right] \text{Var}_\mu \left[ \frac{\pi(a_h^{(1)} | s_h)}{\mu(a_h^{(1)} | s_h)} (V_{h+1}^\pi(s_{h+1}^{(1)}) + r_h^{(1)}) \mid s_h^{(1)} = s_h \right].$$

Handwritten annotations:
 

- Blue arrow pointing from the title to the expectation term in the equation.
- Blue arrow pointing from the expectation term to the variance term.
- Handwritten expression:  $(1-\delta) \cdot \frac{d^\pi(s_h)^2}{d^\mu(s_h)}$
- Handwritten expression:  $n_{s_h} = n \cdot d_h^\mu(s_h)$
- Handwritten expression:  $\mathbb{E} \left[ \frac{1}{d_h^\mu(s_h)} \right]$
- Handwritten expression:  $\left( \frac{d_h^\pi}{d_h^\mu} \right) \geq \dots$

# Bounding error propagation

$$\mathbb{E} \left[ \frac{\tilde{d}_h^\pi(s_h)^2}{n_{s_h}} \mathbf{1} \left( n_{s_h} \geq \frac{nd_h^\mu(s_h)}{(1-\delta)^{-1}} \right) \right] \leq \frac{(1-\delta)^{-1}}{n} \left( \frac{d_h^\pi(s_h)^2}{d_h^\mu(s_h)} + \text{Var} \left[ \tilde{d}_h^\pi(s_h) \right] \right),$$

$$\mathbb{E} x^2 = \text{Var}(x) + (\mathbb{E} x)^2$$

- Bounding the variance of is somewhat tedious
  - Requires use to bound the covariance

$$\text{Var}[\tilde{d}_h^\pi(s_h)] \leq \frac{2(1-\delta)^{-1} h d_h^\pi(s_h)}{n}.$$

$d_h^\pi$  is unbiased

$$\text{Cov} \left( \frac{d_h^\pi(\cdot)}{\mathbb{E} R^s} \right) = \mathbb{E}[\text{Cov}] + \text{Cov}[\mathbb{E}]$$

# Is MIS optimal for OPE in RL?

- It depends on the settings.
- For finite-state / infinite action space, we conjecture that it is.  
(Still an open problem now.)
- For fully tabular setting, it is not optimal, at least asymptotically.

# Tabular MIS

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

- With a minor change to the following recursive estimation

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

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

↕ plugin model based approach

# A short detour: How shall we do DM in RL?

- How would you do DM in this case?

1. Estimate MDP

$$\hat{M} = (S, A, \hat{P}, \hat{r}, \hat{d}_0)$$

2. Plug-in the target policy

Value iteration

$$\begin{aligned} \hat{V}^{\pi}(d_0) &= E^{\pi} [r_1 + \gamma V_1^{\pi}(s_1) + \dots + \gamma^{H-1} r_H] \\ &\parallel \\ &= \sum_{s \sim d_0} d_0(s) \hat{V}_1^{\pi}(s) \\ &\parallel \\ &= \sum_{h=0}^{H-1} \sum_{s \sim d_h^{\pi}(s)} \gamma^h \hat{r}^{\pi}(s) \end{aligned}$$

# TMIS is equivalent to DM --- a model-based approach

$$\begin{aligned}
 \hat{V}_{MIS} &= \frac{1}{n} \sum_{i=1}^n \frac{H}{T+1} \frac{\hat{d}_T^{\pi}(s_t^{(i)})}{\hat{d}_T^{\mu}(s_t^{(i)})} \quad \hat{V}_{DM}(s_t^{(i)}) = \sum_{S \in \mathcal{S}} \frac{H}{T+1} \frac{\hat{d}_T^{\pi}(s)}{\hat{d}_T^{\mu}(s)} \gamma(s) \cdot \mathbb{1}(S_t = s) \\
 &= \sum_{S \in \mathcal{S}} \frac{H}{T+1} \frac{\hat{d}_T^{\pi}(s) \gamma(s)}{\hat{d}_T^{\mu}(s)} \frac{1}{n} \sum_{i=1}^n \mathbb{1}(S_t^{(i)} = s) \\
 &= \hat{V}_{DM}
 \end{aligned}$$

$\hat{d}_1(s) = \text{plug in}$

$$\hat{d}_2(s) = \sum_S \hat{P}_1^{\pi}(s'/s) \hat{d}_1(s')$$

an exact transition dynamics in  $\hat{M}$

# MSE of the TMIS / model-based OPE estimator

- Theorem 3.1 (Yin and W., 2020)

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

Handwritten annotations:

- $\tau(a|s) \cdot \tau_g$  (circled)
- $\frac{d_h^\pi(s_h)^2}{d_h^\mu(s_h)}$  (underlined)
- $\frac{\pi(a_h|s_h)^2}{\mu(a_h|s_h)}$  (circled)
- $\mathbb{E}_a \frac{\pi(a|s)^2}{\mu(a|s)}$  (circled)
- $\text{Var} \left( \frac{\tau}{\mu} \otimes \text{Var}_g \right) = \mathbb{E} \text{Var} + \text{Var} \left( \frac{\tau}{\mu} \right)$  (circled)
- $\frac{\tau}{\mu} \otimes \text{Var}_g$  (crossed out)

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

# TMIS vs on-policy evaluation

**Lemma 3.4.** For any policy  $\pi$  and any MDP.

$\sum_{t=1}^H O(H^2)$

on policy evaluation  $\rightarrow$

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

$\frac{1}{n} \text{Var}(\sum r_t^{(1)})$

$O(H^2)$

Combined with the previous observation:

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

$$\frac{[H^3] \tau \tau_s}{n} \Rightarrow \frac{H^2 \tau \tau_s}{n}$$

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