

Review of Off-Policy Evaluation (OPE)

Chengchun Shi

Outline

1. Off-Policy Evaluation (OPE) Introduction
2. OPE in Contextual Bandits
3. OPE in Reinforcement Learning

Outline

1. Off-Policy Evaluation (OPE) Introduction

2. OPE in Contextual Bandits

3. OPE in Reinforcement Learning

What is OPE and Why OPE

- **Objective:** Evaluate the impact of a **target policy** offline using historical data generated from a different **behavior policy**
- **Motivation:** In many applications, it can be **dangerous** to evaluate a **target policy** by directly running this policy



(a) Health Care

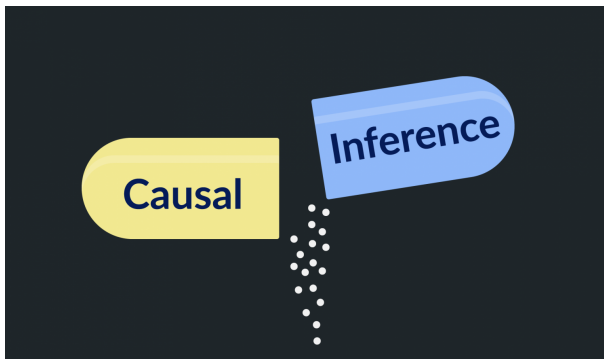


(b) Ridesharing

- **Healthcare:** which **medical treatment** to suggest for a patient
- **Ridesharing:** which **driver** to assign for a call order

Causal Inference

Off-policy evaluation is closely related to **causal inference**, whose objective is to learn the difference between a new treatment and a standard treatment



Outline

1. Off-Policy Evaluation (OPE) Introduction

2. OPE in Contextual Bandits

3. OPE in Reinforcement Learning

Contextual Bandits

- Extension of MAB with **contextual** information.
- A **widely-used** model in medicine and technological industries.
- At time t , the agent
 - Observe a context S_t ;
 - Select an action A_t ;
 - Receives a reward R_t (depends on both S_t and A_t).
- **Objective:** Given an i.i.d. offline dataset $\{(S_t, A_t, R_t) : 0 \leq t < T\}$ generated by a behavior policy b , i.e.,

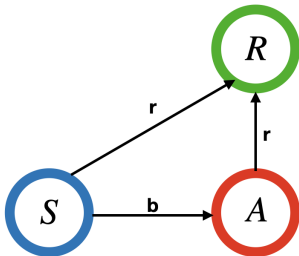
$$\Pr(A_t = a | S_t = s) = b(a|s),$$

we aim to evaluate the mean outcome under a target policy π , i.e.,

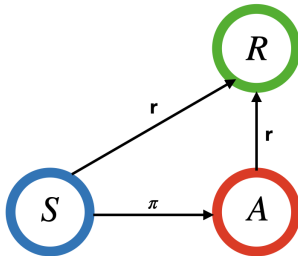
$$\Pr(A_t = a | S_t = s) = \pi(a|s).$$

Challenge

- **Confounding:** State serves as confounding variables that confound the action-reward pair
- **Distributional shift:** The target policy generally differs from the behavior policy



historical data



what we want to evaluate

Challenge (Cont'd)

- Suppose π is a nondynamic policy, i.e., there exists some \mathbf{a} such that $\pi(\mathbf{a}|\mathbf{s}) = 1$ for any \mathbf{s} . We aim to evaluate the value under a given action \mathbf{a} . A naive estimator is

$$\frac{\sum_{t=0}^{T-1} \mathbf{R}_t \mathbb{I}(\mathbf{A}_t = \mathbf{a})}{\sum_{t=0}^{T-1} \mathbb{I}(\mathbf{A}_t = \mathbf{a})} \xrightarrow{P} \mathbb{E}(\mathbf{R}|\mathbf{A} = \mathbf{a})$$

- This estimator is valid only when no confounding variables exist
- According to the causal diagram, the target policy's value equals

$$\mathbb{E}[\mathbb{E}(\mathbf{R}|\mathbf{A} = \mathbf{a}, \mathbf{S})] \neq \mathbb{E}(\mathbf{R}|\mathbf{A} = \mathbf{a})$$

OPE Estimators

- With a general target policy π , the target policy's value equals

$$\sum_{\mathbf{a}} \mathbb{E}[\pi(\mathbf{a}|\mathbf{S})\mathbb{E}(\mathbf{R}|\mathbf{A} = \mathbf{a}, \mathbf{S})] = \sum_{\mathbf{a}} \mathbb{E}[\pi(\mathbf{a}|\mathbf{S})\mathbf{r}(\mathbf{S}, \mathbf{a})],$$

where $\mathbf{r}(\mathbf{s}, \mathbf{a}) = \mathbb{E}(\mathbf{R}|\mathbf{A} = \mathbf{a}, \mathbf{S} = \mathbf{s})$

- Direct estimator
- Importance sampling estimator
- Doubly robust estimator

Direct Estimator

- Given that the target policy's value is given by

$$\sum_{\mathbf{a}} \mathbb{E}[\pi(\mathbf{a}|\mathbf{S})r(\mathbf{S}, \mathbf{a})]$$

- The expectation can be approximated by the sample average, i.e.,

$$\frac{1}{T} \sum_{\mathbf{a}} \sum_{t=0}^{T-1} [\pi(\mathbf{a}|\mathbf{S}_t)r(\mathbf{S}_t, \mathbf{a})]$$

- The reward function can be replaced with some estimator \hat{r} . This yields the direct estimator

$$\frac{1}{T} \sum_{\mathbf{a}} \sum_{t=0}^{T-1} [\pi(\mathbf{a}|\mathbf{S}_t)\hat{r}(\mathbf{S}_t, \mathbf{a})]$$

Importance Sampling Estimator

- Given that the target policy's value is given by

$$\sum_{\mathbf{a}} \mathbb{E}[\pi(\mathbf{a}|\mathbf{S})r(\mathbf{S}, \mathbf{a})]$$

- By the change of measure theory, it equals

$$\sum_{\mathbf{a}} \mathbb{E} \left[\mathbf{b}(\mathbf{a}|\mathbf{S}) \frac{\pi(\mathbf{a}|\mathbf{S})}{\mathbf{b}(\mathbf{a}|\mathbf{S})} r(\mathbf{S}, \mathbf{a}) \right] = \mathbb{E} \left[\frac{\pi(\mathbf{A}|\mathbf{S})}{\mathbf{b}(\mathbf{A}|\mathbf{S})} r(\mathbf{S}, \mathbf{A}) \right] = \mathbb{E} \left[\frac{\pi(\mathbf{A}|\mathbf{S})}{\mathbf{b}(\mathbf{A}|\mathbf{S})} R \right]$$

- This yields the following importance sampling (IS) estimator [Zhang et al., 2012]

$$\frac{1}{T} \sum_{t=0}^{T-1} \frac{\pi(\mathbf{A}_t|\mathbf{S}_t)}{\widehat{\mathbf{b}}(\mathbf{A}_t|\mathbf{S}_t)} R_t,$$

for a given estimator $\widehat{\mathbf{b}}$

Direct Estimator v.s. IS Estimator

- Bias/Variance Trade-Off
- The direct estimator has **some bias**, since r needs to be estimated from data
- The IS estimator has **zero bias** when b is known as in randomized studies
- The IS estimator might have a **large variance** when π differs significantly from b
- Suppose $R = r(S, A) + \varepsilon$ for some ε independent of (S, A) ,

$$\begin{aligned}\text{Var} \left[\frac{\pi(A|S)}{b(A|S)} R \right] &= \mathbb{E} \left[\frac{\pi(A|S)}{b(A|S)} \{R - r(S, A)\} \right]^2 + \text{some term} \\ &= \sigma^2 \mathbb{E} \left[\frac{\pi^2(A|S)}{b^2(A|S)} \right] + \text{some term},\end{aligned}$$

where $\sigma^2 = \text{Var}(\varepsilon)$

Extensions

- When π differs from b significantly, IS estimator suffers from **large variance** and becomes **unstable**
- Solutions sought by using **self-normalized** and/or **truncated** IS
- **Self-normalized** IS

$$\left[\frac{1}{T} \sum_{t=0}^{T-1} \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{b(\mathbf{A}_t | \mathbf{S}_t)} \right]^{-1} \frac{1}{T} \sum_{t=0}^{T-1} \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{b(\mathbf{A}_t | \mathbf{S}_t)} R_t$$

- **Truncated** IS

$$\frac{1}{T} \sum_{t=0}^{T-1} \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{\max(\hat{b}(\mathbf{A}_t | \mathbf{S}_t), \epsilon)} R_t,$$

for some $\epsilon > 0$

Doubly Robust Estimator

- Direct estimator

$$\frac{1}{T} \sum_{\mathbf{a}} \sum_{t=0}^{T-1} [\pi(\mathbf{a} | \mathbf{S}_t) \hat{r}(\mathbf{S}_t, \mathbf{a})]$$

requires \hat{r} to be consistent

- IS estimator

$$\frac{1}{T} \sum_{t=0}^{T-1} \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{\hat{b}(\mathbf{A}_t | \mathbf{S}_t)} R_t,$$

requires \hat{b} to be consistent

- Doubly robust (DR) estimator combines both, and requires **either \hat{r} or \hat{b}** to be consistent (“**doubly-robustness**” property)

Doubly Robust Estimator (Cont'd)

- Consider the estimating function

$$\phi(\mathbf{S}, \mathbf{A}, \mathbf{R}) = \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{S}) \hat{r}(\mathbf{S}, \mathbf{a}) + \frac{\pi(\mathbf{A}|\mathbf{S})}{\hat{b}(\mathbf{A}|\mathbf{S})} [\mathbf{R} - \hat{r}(\mathbf{S}, \mathbf{A})]$$

- First term on the RHS is the estimating function of the direct estimator
- Second term corresponds to the **augmentation term**
 - Zero mean when $\hat{r} = r$
 - Debias the bias of the direct estimator
 - Offering additional robustness against model misspecification of \hat{r}
- DR estimator given by $T^{-1} \sum_{t=0}^{T-1} \phi(\mathbf{S}_t, \mathbf{A}_t, \mathbf{R}_t)$

Fact 1: Double Robustness

- The estimating function

$$\phi(\mathbf{S}, \mathbf{A}, \mathbf{R}) = \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{S}) \hat{r}(\mathbf{S}, \mathbf{a}) + \frac{\pi(\mathbf{A}|\mathbf{S})}{\hat{b}(\mathbf{A}|\mathbf{S})} [\mathbf{R} - \hat{r}(\mathbf{S}, \mathbf{A})]$$

- In large sample size, DR estimator converges to $\mathbb{E}\phi(\mathbf{S}, \mathbf{A}, \mathbf{R})$
- When $\hat{r} = r$, the augmentation term has zero mean. It follows that

$$\mathbb{E}\phi(\mathbf{S}, \mathbf{A}, \mathbf{R}) = \sum_{\mathbf{a}} \mathbb{E}[\pi(\mathbf{a}|\mathbf{S}) r(\mathbf{S}, \mathbf{a})] = \text{target policy's value}$$

- When $\hat{b} = b$, it has the same mean as the IS estimator

$$\begin{aligned} \mathbb{E}\phi(\mathbf{S}, \mathbf{A}, \mathbf{R}) &= \mathbb{E} \left[\frac{\pi(\mathbf{A}|\mathbf{S})}{b(\mathbf{A}|\mathbf{S})} \mathbf{R} \right] + \mathbb{E} \left[\sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{S}) \hat{r}(\mathbf{S}, \mathbf{a}) - \frac{\pi(\mathbf{A}|\mathbf{S})}{b(\mathbf{A}|\mathbf{S})} \hat{r}(\mathbf{S}, \mathbf{A}) \right] \\ &= \mathbb{E} \left[\frac{\pi(\mathbf{A}|\mathbf{S})}{b(\mathbf{A}|\mathbf{S})} \mathbf{R} \right] = \text{target policy's value} \end{aligned}$$

Fact 2: Efficiency

- When $\hat{\mathbf{b}} = \mathbf{b}$, the estimating function

$$\phi(\mathbf{S}, \mathbf{A}, \mathbf{R}) = \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{S}) \hat{r}(\mathbf{S}, \mathbf{a}) + \frac{\pi(\mathbf{A}|\mathbf{S})}{\mathbf{b}(\mathbf{A}|\mathbf{S})} [\mathbf{R} - \hat{r}(\mathbf{S}, \mathbf{A})]$$

- The MSE of DR estimator is proportional to the variance of $\phi(\mathbf{S}, \mathbf{A}, \mathbf{R})$

$$\text{Var}(\phi(\mathbf{S}, \mathbf{A}, \mathbf{R})) = \mathbb{E}[\text{Var}(\phi(\mathbf{S}, \mathbf{A}, \mathbf{R})|\mathbf{S}, \mathbf{A})] + \text{Var}[\mathbb{E}(\phi(\mathbf{S}, \mathbf{A}, \mathbf{R})|\mathbf{S}, \mathbf{A})]$$

- The first term on the RHS is independent of $\hat{\mathbf{r}}$
- The second term is minimized when $\hat{\mathbf{r}} = \mathbf{r}$
- A good working model for \mathbf{r} improves the estimator's efficiency
- When $\hat{\mathbf{r}} = \mathbf{r}$, the estimator achieves the **efficiency bound** [e.g., smallest MSE among a class of regular estimators; see Tsiatis, 2007]

Fact 3: Efficiency

- When $\hat{\mathbf{b}}$ is estimated from data and the model is **correctly specified**, the estimator's MSE would be **generally smaller than** the one that uses the oracle behavior policy \mathbf{b} [Tsiatis, 2007]
- Estimating $\hat{\mathbf{b}}$ yields a more efficient estimator, even if we know the oracle \mathbf{b}
- **Multi-armed bandit** example without context information
 - **Objective:** evaluate $\mathbb{E}(\mathbf{R} | \mathbf{A} = \mathbf{a})$ for a given \mathbf{a}
 - IS estimator with **known** $\Pr(\mathbf{A} = \mathbf{a})$

$$\frac{\sum_{t=0}^{T-1} \mathbb{I}(\mathbf{A}_t = \mathbf{a}) \mathbf{R}_t}{T \Pr(\mathbf{A}_t = \mathbf{a})}$$

- IS estimator with **estimated** $\Pr(\mathbf{A} = \mathbf{a})$ has a **smaller** asymptotic variance

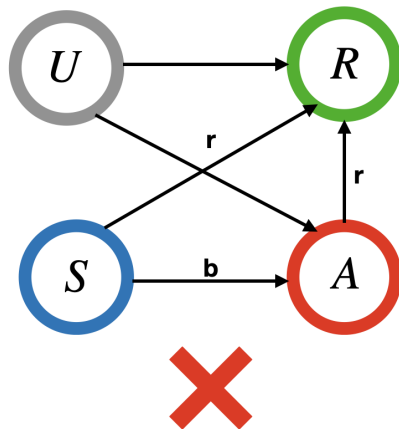
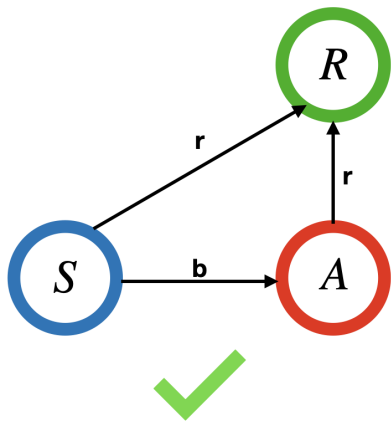
$$\frac{\sum_{t=0}^{T-1} \mathbb{I}(\mathbf{A}_t = \mathbf{a}) \mathbf{R}_t}{\sum_{t=0}^{T-1} \mathbb{I}(\mathbf{A}_t = \mathbf{a})}$$

- Same holds true in **RL** [Hanna et al., 2019, 2021, Zhou et al., 2025]

Fact 4: Asymptotic Normality

- The DR estimator converges at a parametric rate and is asymptotically normal even when both \hat{r} and \hat{b} converge **slower** than the parameter rate (i.e., root- n rate)
- This observation allows us to apply machine learning methods to estimate both nuisance functions, leading to the **double machine learning** estimator [Chernozhukov et al., 2017]
- Indeed, it only requires \hat{r} and \hat{b} to converge at a rate of $o_p(n^{-1/4})$, due to the double robustness property

Assumption: No Unmeasured Confounders



Lecture Outline

1. Off-Policy Evaluation (OPE) Introduction
2. OPE in Contextual Bandits
- 3. OPE in Reinforcement Learning**

General OPE Problem

- **Objective:** Given an offline dataset $\{(\mathbf{S}_{i,t}, \mathbf{A}_{i,t}, \mathbf{R}_{i,t}) : 1 \leq i \leq N, 0 \leq t \leq T\}$ generated by a behavior policy \mathbf{b} , where i indexes the i th episode and t indexes the t th time point, we aim to evaluate the mean return under a target policy π

$$\mathbb{E}^{\pi} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{R}_t \right] = \mathbb{E} \mathbf{V}^{\pi}(\mathbf{S}_0)$$

When $\gamma = 1$, the task is assumed to be episodic

- We focus on the case where both π and \mathbf{b} are **stationary** policies
- Challenge: **Distributional shift**
 - In the offline dataset, actions are generated according to \mathbf{b}
 - The target policy π we wish to evaluate is different from \mathbf{b}

Direct Estimator

- The target policy's value is given by $\mathbb{E} \mathbf{V}^\pi(\mathbf{S}_0)$, or equivalently,

$$\mathbb{E}[\sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{S}_0) Q^\pi(\mathbf{S}_0, \mathbf{a})]$$

- The expectation can be approximated via the **empirical initial state distribution**
- Q-learning is an **off-policy** algorithm. Can be applied to learn Q^π offline
- This yields the direct estimator

$$\frac{1}{N} \sum_{i=1}^N \sum_{\mathbf{a}} \pi(\mathbf{a}|\mathbf{S}_{i,0}) \hat{Q}(\mathbf{S}_{i,0}, \mathbf{a})$$

- It remains to compute \hat{Q}

Fitted Q-Evaluation [Le et al., 2019]

- Bellman equation

$$\mathbb{E}[R_t + \gamma \pi(a|S_{t+1}) Q^\pi(S_{t+1}, a) | S_t, A_t] = Q^\pi(S_t, A_t)$$

- Both LHS and RHS involves Q^π
- Repeat the following procedure
 1. Compute \hat{Q} as the argmin of

$$\arg \min_{\tilde{Q}} \sum_t \left[R_{i,t} + \gamma \sum_a \pi(a|S_{i,t+1}) \tilde{Q}(S_{i,t+1}, a) - Q(S_{i,t}, A_{i,t}) \right]^2$$

2. Set $\tilde{Q} = \hat{Q}$
- Designed for learning Q^π
 - Do **not** require actions to follow the target policy

Other Direct Estimators

- Sieve-based estimator [Shi et al., 2022c]
 - Use linear sieves to parametrize Q^π
 - Estimate regression coefficients by solving the Bellman equation
- Kernel-based estimator [Liao et al., 2021]
 - Use RHKSs to parametrize Q^π
 - Estimate parameters by solving a coupled optimization [Farahmand et al., 2016]
- Limiting distributions of value estimators are derived in the two papers

Stepwise IS Estimator [Zhang et al., 2013]

- Consider episodic task where T is the termination time
- Importance sampling ratio needs to be employed

$$\begin{aligned}\mathbb{E}^{\pi} R_0 &= \mathbb{E}^b \left[\frac{\pi(\mathbf{A}_0 | \mathbf{S}_0)}{b(\mathbf{A}_0 | \mathbf{S}_0)} R_0 \right] \\ \mathbb{E}^{\pi} R_1 &= \mathbb{E}^b \left[\frac{\pi(\mathbf{A}_0 | \mathbf{S}_0)}{b(\mathbf{A}_0 | \mathbf{S}_0)} \frac{\pi(\mathbf{A}_1 | \mathbf{S}_1)}{b(\mathbf{A}_1 | \mathbf{S}_1)} R_1 \right] \\ &\vdots \\ \mathbb{E}^{\pi} R_t &= \mathbb{E}^b \left[\frac{\pi(\mathbf{A}_0 | \mathbf{S}_0)}{b(\mathbf{A}_0 | \mathbf{S}_0)} \cdots \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{b(\mathbf{A}_t | \mathbf{S}_t)} R_t \right]\end{aligned}$$

Stepwise IS Estimator (Cont'd)

- According to this logic, the target policy's value can be represented by

$$\mathbb{E} \left[\sum_{t=0}^T \gamma^t \left\{ \prod_{j=0}^t \frac{\pi(\mathbf{A}_j | \mathbf{S}_j)}{\mathbf{b}(\mathbf{A}_j | \mathbf{S}_j)} \right\} R_t \right]$$

- This yields the stepwise IS estimator

$$\frac{1}{N} \sum_{i=1}^N \left[\sum_{t=0}^T \gamma^t \left\{ \prod_{j=0}^t \frac{\pi(\mathbf{A}_{i,j} | \mathbf{S}_{i,j})}{\widehat{\mathbf{b}}(\mathbf{A}_{i,j} | \mathbf{S}_{i,j})} \right\} R_{i,t} \right]$$

for a given estimator $\widehat{\mathbf{b}}$ computed using supervised learning algorithms

Limitation

- Stepwise IS suffers from a **large variance**
- In particular, the IS ratio at time t is the product of individual ratios from the **initial** time to time t

$$\prod_{j=0}^t \frac{\pi(\mathbf{A}_j | \mathbf{S}_j)}{b(\mathbf{A}_j | \mathbf{S}_j)}$$

- Variance of the ratio grows **exponentially** with respect to t , referred to as the **curse of horizon** [Liu et al., 2018]
- Extension: **Doubly-robust** estimator by [Jiang and Li, 2016]

Pros & Cons of Direct v.s. Stepwise IS

- Bias/Variance Trade-Off
- When \mathbf{b} is known, stepwise IS is an **unbiased** estimator since

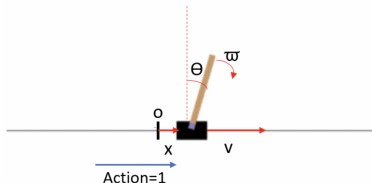
$$\mathbb{E}^{\pi} R_t = \mathbb{E}^{\mathbf{b}} \left[\frac{\pi(\mathbf{A}_0 | \mathbf{S}_0)}{\mathbf{b}(\mathbf{A}_0 | \mathbf{S}_0)} \cdots \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{\mathbf{b}(\mathbf{A}_t | \mathbf{S}_t)} R_t \right]$$

- Direct estimator has **some bias**, since Q^{π} needs to be estimated from data
- Stepwise IS suffers from **curse of horizon** and a **large variance**
- Direct estimator has a much lower variance

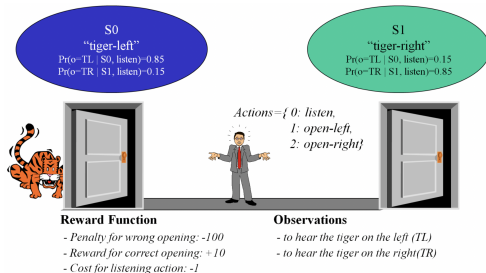
Pros & Cons of Direct v.s. Stepwise IS (Cont'd)

- Direct estimator exploits **Markov** & **stationary** properties
- Relies on the **Bellman equation**
- More **efficient** in MDP environments

frame: 53, Obs: (0.018, 0.669, 0.286, 0.618)
Action: 1.0, Cumulative Reward: 47.0, Done: 1



- SIS does **not** exploit these properties
- More **flexible** in non-MDP environments (e.g., POMDP)



Marginalized IS Estimator

- As we have discussed, stepwise IS suffers from **curse of horizon**
- Curse of horizon is **unavoidable** in general **Non-Markov decision processes** (e.g., POMDP)
- Under some additional model assumptions (e.g., Markovianity & time-homogeneity), it is possible to break the curse of horizon using **marginalized IS** estimator
- Stepwise IS does **not** exploit these properties

Marginalized IS Estimator (Cont'd)

- Stepwise IS uses the **cumulative** IS ratio

$$\mathbb{E}^{\pi} R_t = \mathbb{E}^b \left[\frac{\pi(\mathbf{A}_0 | \mathbf{S}_0)}{b(\mathbf{A}_0 | \mathbf{S}_0)} \cdots \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{b(\mathbf{A}_t | \mathbf{S}_t)} R_t \right]$$

- Under Markovianity (TMDP), marginalized IS uses the **marginalized** IS ratio

$$\mathbb{E}^{\pi} R_t = \mathbb{E}^b \left[\frac{p_t^{\pi}(\mathbf{S}_t, \mathbf{A}_t)}{p_t^b(\mathbf{S}_t, \mathbf{A}_t)} R_t \right] \quad (1)$$

where p_t^{π} and p_t^b are the marginal density functions of $(\mathbf{S}_t, \mathbf{A}_t)$ under π and b

- The resulting marginalized IS estimator can be derived from (1)

Marginalized IS Estimator

- Under Markovianity and time-homogeneity (MDP),

$$\mathbb{E} V^{\pi}(\mathbf{s}_0) = \mathbb{E}^{\mathbf{b}} \left[\frac{\sum_{t=0}^{\infty} \gamma^t \mathbf{p}_t^{\pi}(\mathbf{s}, \mathbf{a})}{\mathbf{p}_{\infty}(\mathbf{s}, \mathbf{a})} R \right] \quad (2)$$

where \mathbf{p}_{∞} denotes the limiting state-action distribution under \mathbf{b} and the numerator corresponds to the γ -discounted state-action visitation probability

- The resulting marginalized IS estimator can be derived from (2)
- Marginal IS ratio can be estimated via **minimax learning** [Uehara et al., 2019]
- Closed-form expression is available when using **linear sieves**
- Coupled optimization can also be employed when using **RKHSs** [Liao et al., 2020]
- Alternatively, we can use **RKHSs** to parametrize the discriminator class, use **neural networks** to parametrize the ratio and apply SGD for parameter estimation

Double RL [Kallus and Uehara, 2019]

- Double RL extends DR in **contextual bandits** to the general RL problem
- Similar to DR, the estimator can be represented as

Direct Estimator + Augmentation Term

- **Augmentation** term is to **debias** the bias of direct estimator and offer protection against model misspecification of Q^π ; it relies on the marginalized IS ratio
- Similar to DR, the estimator is **doubly-robust**, e.g., consistent when either Q^π or the marginalized IS ratio is correct
- Similar to DR, the estimator achieves the **efficiency bound** in MDPs

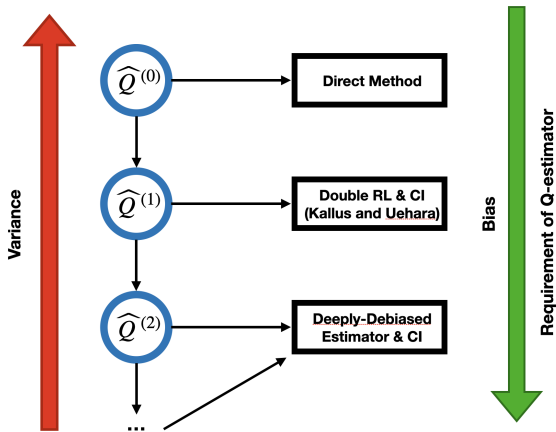
Fact 5: Efficiency

- Direct estimators (based on linear sieves or RKHSs) also achieve the **efficiency bound** in MDPs [Liao et al., 2021, Shi et al., 2022b]
- Marginalized IS estimators (based on linear sieves) also achieve the **efficiency bound** in MDPs
- When using linear sieves,

direct estimator = marginalized IS estimator = double RL estimator

See the proof of Theorem 1-LSTD of Wen et al. [2025] and proof of Theorem 8 of Zhou et al. [2025].

Deeply-Debiased OPE [Shi et al., 2021]



- Constructed based on high-order influence function [Robins et al., 2008, 2017]
- Ensures bias decays much faster than standard deviation
- Allows to provide valid **uncertainty quantification** (e.g., confidence interval)

Other Topics

- **Partially observable** environments [Uehara et al., 2023, Hu and Wager, 2023]
- **Heavy-tailed** rewards [Xu et al., 2022, Liu et al., 2023, Rowland et al., 2023, Zhu et al., 2024, Behnamnia et al., 2025]
- **Confounded OPE**
 - Confounded POMDP [Tennenholtz et al., 2020, Bennett and Kallus, 2021, Nair and Jiang, 2021, Shi et al., 2022a]
 - Confounded MDPs [Zhang and Bareinboim, 2016, Wang et al., 2021, Fu et al., 2022, Shi et al., 2024, Xu et al., 2023, Yu et al., 2024]
- **Experimental designs** [Hanna et al., 2017, Mukherjee et al., 2022, Wan et al., 2022, Li et al., 2023, Sun et al., 2024, Wen et al., 2025]
- Evaluation of the **optimal policy**
 - Inference is challenging in **nonregular** settings where the optimal policy is not unique
 - m -out-of- n bootstrap [Chakraborty et al., 2013]
 - Martingale-based method [Luedtke and Van Der Laan, 2016, Shi et al., 2022c]
 - Subagging-based method [Shi et al., 2020]

References I

- Armin Behnamnia, Gholamali Aminian, Alireza Aghaei, Chengchun Shi, Vincent Y. F. Tan, and Hamid R. Rabiee. Log-sum-exponential estimator for off-policy evaluation and learning. In *International Conference on Machine Learning*. PMLR, 2025.
- Andrew Bennett and Nathan Kallus. Proximal reinforcement learning: Efficient off-policy evaluation in partially observed markov decision processes. *arXiv preprint arXiv:2110.15332*, 2021.
- Bibhas Chakraborty, Eric B Laber, and Yingqi Zhao. Inference for optimal dynamic treatment regimes using an adaptive m-out-of-n bootstrap scheme. *Biometrics*, 69(3): 714–723, 2013.
- Victor Chernozhukov, Denis Chetverikov, Mert Demirer, Esther Duflo, Christian Hansen, and Whitney Newey. Double/debiased/neyman machine learning of treatment effects. *American Economic Review*, 107(5):261–65, 2017.

References II

- Amir-massoud Farahmand, Mohammad Ghavamzadeh, Csaba Szepesvári, and Shie Mannor. Regularized policy iteration with nonparametric function spaces. *The Journal of Machine Learning Research*, 17(1):4809–4874, 2016.
- Zuyue Fu, Zhengling Qi, Zhaoran Wang, Zhuoran Yang, Yanxun Xu, and Michael R Kosorok. Offline reinforcement learning with instrumental variables in confounded markov decision processes. *arXiv preprint arXiv:2209.08666*, 2022.
- Josiah Hanna, Scott Niekum, and Peter Stone. Importance sampling policy evaluation with an estimated behavior policy. In *International Conference on Machine Learning*, pages 2605–2613. PMLR, 2019.
- Josiah P. Hanna, Philip S. Thomas, Peter Stone, and Scott Niekum. Data-efficient policy evaluation through behavior policy search. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1394–1403. PMLR, 06–11 Aug 2017.

References III

- Josiah P Hanna, Scott Niekum, and Peter Stone. Importance sampling in reinforcement learning with an estimated behavior policy. *Machine Learning*, 110(6):1267–1317, 2021.
- Yuchen Hu and Stefan Wager. Off-policy evaluation in partially observed Markov decision processes under sequential ignorability. *The Annals of Statistics*, 51(4):1561 – 1585, 2023. doi: 10.1214/23-AOS2287.
- Nan Jiang and Lihong Li. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pages 652–661. PMLR, 2016.
- Nathan Kallus and Masatoshi Uehara. Efficiently breaking the curse of horizon in off-policy evaluation with double reinforcement learning. *arXiv preprint arXiv:1909.05850*, 2019.
- Hoang M Le, Cameron Voloshin, and Yisong Yue. Batch policy learning under constraints. *arXiv preprint arXiv:1903.08738*, 2019.

References IV

- Ting Li, Chengchun Shi, Jianing Wang, Fan Zhou, et al. Optimal treatment allocation for efficient policy evaluation in sequential decision making. *Advances in Neural Information Processing Systems*, 36:48890–48905, 2023.
- Peng Liao, Zhengling Qi, and Susan Murphy. Batch policy learning in average reward markov decision processes. *arXiv preprint arXiv:2007.11771*, 2020.
- Peng Liao, Predrag Klasnja, and Susan Murphy. Off-policy estimation of long-term average outcomes with applications to mobile health. *Journal of the American Statistical Association*, 116(533):382–391, 2021.
- Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. Breaking the curse of horizon: Infinite-horizon off-policy estimation. In *Advances in Neural Information Processing Systems*, pages 5356–5366, 2018.

References V

- Weidong Liu, Jiyuan Tu, Yichen Zhang, and Xi Chen. Online estimation and inference for robust policy evaluation in reinforcement learning. *arXiv preprint arXiv:2310.02581*, 2023.
- Alexander R Luedtke and Mark J Van Der Laan. Statistical inference for the mean outcome under a possibly non-unique optimal treatment strategy. *Annals of statistics*, 44(2):713, 2016.
- Subhojyoti Mukherjee, Josiah P. Hanna, and Robert D Nowak. Revar: Strengthening policy evaluation via reduced variance sampling. In James Cussens and Kun Zhang, editors, *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 1413–1422. PMLR, 01–05 Aug 2022.
- Yash Nair and Nan Jiang. A spectral approach to off-policy evaluation for pomdps. *arXiv preprint arXiv:2109.10502*, 2021.

References VI

- James Robins, Lingling Li, Eric Tchetgen, Aad van der Vaart, et al. Higher order influence functions and minimax estimation of nonlinear functionals. In *Probability and statistics: essays in honor of David A. Freedman*, pages 335–421. Institute of Mathematical Statistics, 2008.
- James M Robins, Lingling Li, Rajarshi Mukherjee, Eric Tchetgen Tchetgen, Aad van der Vaart, et al. Minimax estimation of a functional on a structured high-dimensional model. *The Annals of Statistics*, 45(5):1951–1987, 2017.
- Mark Rowland, Yunhao Tang, Clare Lyle, Rémi Munos, Marc G Bellemare, and Will Dabney. The statistical benefits of quantile temporal-difference learning for value estimation. In *International Conference on Machine Learning*, pages 29210–29231. PMLR, 2023.

References VII

- Chengchun Shi, Wenbin Lu, and Rui Song. Breaking the curse of nonregularity with subagging: inference of the mean outcome under optimal treatment regimes. *Journal of Machine Learning Research*, 21, 2020.
- Chengchun Shi, Runzhe Wan, Victor Chernozhukov, and Rui Song. Deeply-debiased off-policy interval estimation. In *International Conference on Machine Learning*, pages 9580–9591. PMLR, 2021.
- Chengchun Shi, Masatoshi Uehara, Jiawei Huang, and Nan Jiang. A minimax learning approach to off-policy evaluation in confounded partially observable markov decision processes. In *International Conference on Machine Learning*, pages 20057–20094. PMLR, 2022a.

References VIII

- Chengchun Shi, Xiaoyu Wang, Shikai Luo, Hongtu Zhu, Jieping Ye, and Rui Song. Dynamic causal effects evaluation in a/b testing with a reinforcement learning framework. *Journal of the American Statistical Association*, (just-accepted):1–29, 2022b.
- Chengchun Shi, Sheng Zhang, Wenbin Lu, and Rui Song. Statistical inference of the value function for reinforcement learning in infinite-horizon settings. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 84(3):765–793, 2022c.
- Chengchun Shi, Jin Zhu, Ye Shen, Shikai Luo, Hongtu Zhu, and Rui Song. Off-policy confidence interval estimation with confounded markov decision process. *Journal of the American Statistical Association*, 119(545):273–284, 2024.
- Ke Sun, Linglong Kong, Hongtu Zhu, and Chengchun Shi. Optimal treatment allocation strategies for a/b testing in partially observable time series experiments. *arXiv preprint arXiv:2408.05342*, 2024.

References IX

- Guy Tennenholtz, Uri Shalit, and Shie Mannor. Off-policy evaluation in partially observable environments. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 10276–10283, 2020.
- Anastasios Tsiatis. *Semiparametric theory and missing data*. Springer Science & Business Media, 2007.
- Masatoshi Uehara, Jiawei Huang, and Nan Jiang. Minimax weight and q-function learning for off-policy evaluation. *arXiv preprint arXiv:1910.12809*, 2019.
- Masatoshi Uehara, Haruka Kiyohara, Andrew Bennett, Victor Chernozhukov, Nan Jiang, Nathan Kallus, Chengchun Shi, and Wen Sun. Future-dependent value-based off-policy evaluation in pomdps. In *Advances in Neural Information Processing Systems*, volume 36, pages 15991–16008. Curran Associates, Inc., 2023.

References X

- Runzhe Wan, Branislav Kveton, and Rui Song. Safe exploration for efficient policy evaluation and comparison. In *International Conference on Machine Learning*, pages 22491–22511. PMLR, 2022.
- Lingxiao Wang, Zhuoran Yang, and Zhaoran Wang. Provably efficient causal reinforcement learning with confounded observational data. *Advances in Neural Information Processing Systems*, 34:21164–21175, 2021.
- Qianglin Wen, Chengchun Shi, Ying Yang, Niansheng Tang, and Hongtu Zhu. Unraveling the interplay between carryover effects and reward autocorrelations in switchback experiments. In *Forty-second International Conference on Machine Learning*. PMLR, 2025.
- Yang Xu, Chengchun Shi, Shikai Luo, Lan Wang, and Rui Song. Quantile off-policy evaluation via deep conditional generative learning. *arXiv preprint arXiv:2212.14466*, 2022.

References XI

- Yang Xu, Jin Zhu, Chengchun Shi, Shikai Luo, and Rui Song. An instrumental variable approach to confounded off-policy evaluation. In *International Conference on Machine Learning*, pages 38848–38880. PMLR, 2023.
- Shuguang Yu, Shuxing Fang, Ruixin Peng, Zhengling Qi, Fan Zhou, and Chengchun Shi. Two-way deconfounder for off-policy evaluation in causal reinforcement learning. *Advances in Neural Information Processing Systems*, 37:78169–78200, 2024.
- Baqun Zhang, Anastasios A Tsiatis, Eric B Laber, and Marie Davidian. A robust method for estimating optimal treatment regimes. *Biometrics*, 68(4):1010–1018, 2012.
- Baqun Zhang, Anastasios A Tsiatis, Eric B Laber, and Marie Davidian. Robust estimation of optimal dynamic treatment regimes for sequential treatment decisions. *Biometrika*, 100(3):681–694, 2013.

References XII

- Junzhe Zhang and Elias Bareinboim. Markov decision processes with unobserved confounders: A causal approach. Technical report, Technical report, Technical Report R-23, Purdue AI Lab, 2016.
- Hongyi Zhou, Josiah P Hanna, Jin Zhu, Ying Yang, and Chengchun Shi. Demystifying the paradox of importance sampling with an estimated history-dependent behavior policy in off-policy evaluation. In *Forty-second International Conference on Machine Learning*. PMLR, 2025.
- Jin Zhu, Runzhe Wan, Zhengling Qi, Shikai Luo, and Chengchun Shi. Robust offline reinforcement learning with heavy-tailed rewards. In *International Conference on Artificial Intelligence and Statistics*, pages 541–549. PMLR, 2024.