

ST455: Reinforcement Learning

Lecture 11: Off-Policy Evaluation

Chengchun Shi

Lecture Outline

1. Off-Policy Evaluation (OPE) Introduction
2. OPE in Contextual Bandits
3. OPE in Reinforcement Learning
4. State-of-the-Art OPE Algorithms

Lecture Outline

- 1. Off-Policy Evaluation (OPE) Introduction**
2. OPE in Contextual Bandits
3. OPE in Reinforcement Learning
4. State-of-the-Art OPE Algorithms

What is Off-Policy Evaluation

- **Objective:** Evaluate the impact of a **target policy** offline using historical data generated from a different **behavior policy**
- **Setting:** Offline RL with a precollected data



(a) Health Care



(b) Robotics



(c) Ridesharing



(d) Auto-driving

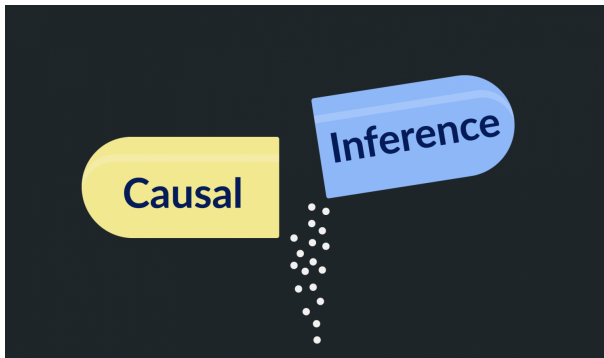
Why Off-Policy Evaluation

In many applications, it can be **dangerous** to evaluate a **target policy** by directly running this policy

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

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



Causal Inference (Cont'd)

[home](#) / [insights](#) / [agenda](#) / [causality and natural experiments](#) the 2021 nobel prize in economic sciences

Causality and natural experiments: the 2021 Nobel Prize in Economic Sciences

26 NOV 2021

Offline Policy Optimisation

- Off-policy evaluation is also related to **offline** policy learning (Lecture 10), whose objective is to learn an optimal policy offline using historical data
- Suppose we are able to evaluate the **value** of any policy, it suffices to pick the policy that maximises the value

Lecture Outline

1. Off-Policy Evaluation (OPE) Introduction
- 2. OPE in Contextual Bandits**
3. OPE in Reinforcement Learning
4. State-of-the-Art OPE Algorithms

Recap: Contextual Bandits

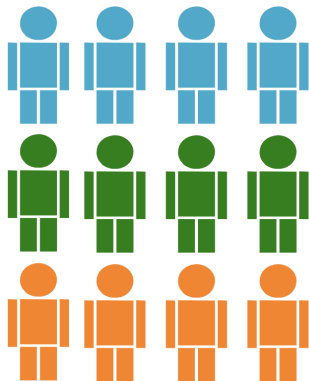
- Extension of MAB with **contextual** information.
- A **widely-used** model in medicine and technological industries.
- At time t , the agent
 - Observe a context S_t ;
 - Select an action A_t ;
 - Receives a reward R_t (depends on both S_t and A_t).
- **Objective:** Given an i.i.d. offline dataset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

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

Application I: Precision Medicine



Patients

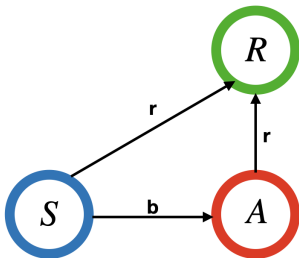


Application II: Personalized Recommendation

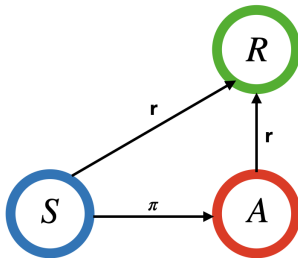


Challenge

- **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

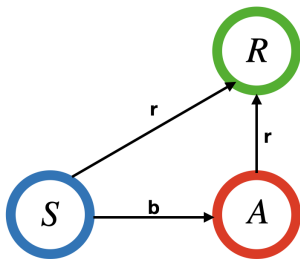
Challenge (Cont'd)

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

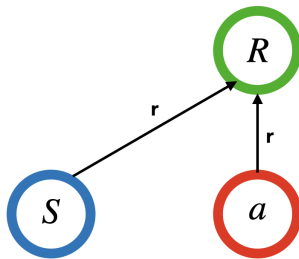
$$\frac{\sum_{t=0}^{T-1} R_t \mathbb{I}(\mathbf{A}_t = \mathbf{a})}{\sum_{t=0}^{T-1} \mathbb{I}(\mathbf{A}_t = \mathbf{a})} \xrightarrow{P} \mathbb{E}(R|\mathbf{A} = \mathbf{a})$$

- This estimator is valid only when no confounding variables exist

Challenge (Cont'd)



historical data



what we want to evaluate

According to the causal diagram, the target policy's value equals

$$\mathbb{E}[\mathbb{E}(\textcolor{green}{R} | \textcolor{red}{A} = \textcolor{red}{a}, \textcolor{blue}{S})] \neq \mathbb{E}(\textcolor{green}{R} | \textcolor{red}{A} = \textcolor{red}{a})$$

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

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

Direct Estimator (Cont'd)

- \hat{r} estimated using supervised learning

$$\begin{aligned} S_0, A_0 &\rightarrow R_0 \\ S_1, A_1 &\rightarrow R_1 \\ &\vdots \\ S_{T-1}, A_{T-1} &\rightarrow R_{T-1} \end{aligned}$$

- Loss function: least square/Huber loss
- Computer parameter that minimizes empirical loss

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 $\widehat{\mathbf{b}}$

Importance Sampling Estimator (Cont'd)

- The ratio $\pi(\mathbf{a}|\mathbf{s})/\mathbf{b}(\mathbf{a}|\mathbf{s})$ is referred to as the **importance sampling ratio**
- It measures the difference between the behavior and target policies
- When $\pi = \mathbf{b}$, the ratio equals **1** for any \mathbf{a} and \mathbf{s}
- In general, the ratio centres at **1**

$$\mathbb{E} \left[\frac{\pi(\mathbf{A}|\mathbf{S})}{\mathbf{b}(\mathbf{A}|\mathbf{S})} \right] = \mathbf{1}$$

Importance Sampling Estimator (Cont'd)

- In **randomized studies**, b is known
- In **observational studies**, b needs to be estimated from data
- \hat{b} estimated using supervised learning

$$\begin{array}{ccc} S_0 & \rightarrow & A_0 \\ S_1 & \rightarrow & A_1 \\ & \vdots & \\ S_{T-1} & \rightarrow & A_{T-1} \end{array}$$

- Loss function: logistic regression loss
- Computer parameter that minimizes empirical loss

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

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

- Equivalent to IS estimator in large samples, by noting that

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

- Stabilize the important sampling ratio in finite samples

Extensions (Cont'd)

- 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 $\epsilon > 0$

- Truncate the behavior policy whose value is smaller than ϵ
- Avoid **extremely large ratio**, thus reducing the variance of the estimator

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 $\mathbf{T}^{-1} \sum_{t=0}^{T-1} \phi(\mathbf{S}_t, \mathbf{A}_t, \mathbf{R}_t)$

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** $\Pr(\mathbf{A} = \mathbf{a})$

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

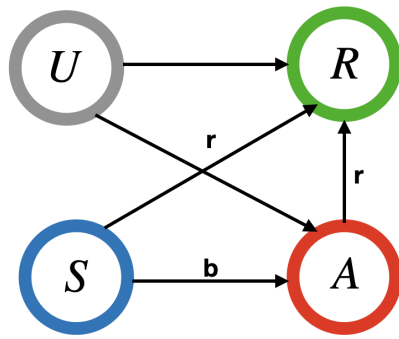
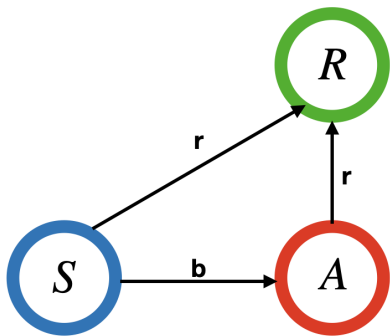
- IS estimator with **estimated** $\Pr(\mathbf{A} = \mathbf{a})$ has a **smaller** asymptotic variance

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

Assumption: No Unmeasured Confounders

- All three estimators (direct estimator, IS, DR) rely on the **no unmeasured confounders** assumption
- They are **biased** when this assumption is violated
- It requires **all** confounders that confound the action-reward relationship are included in the state
- This assumption is **cannot** be verified in practice
- When violated, we may use some **auxiliary variable** (e.g., instrumental variables, mediators) for consistent policy evaluation [Angrist et al., 1996, Pearl, 2009]

Assumption: No Unmeasured Confounders (Cont'd)



Lecture Outline

1. Off-Policy Evaluation (OPE) Introduction
2. OPE in Contextual Bandits
- 3. OPE in Reinforcement Learning**
4. State-of-the-Art OPE Algorithms

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}
- Existing prediction algorithms (e.g., MC, TD) designed in online settings are **not** applicable

Recap: MC Prediction

- **Objective:** learns V^π from experience under π
- MC Policy Evaluation: $V(s) \leftarrow \text{average}[\text{Returns}(s)]$
- Incremental update for every-visit MC prediction:

$$V(s_t) \leftarrow V(s_t) + \alpha_t [G_t - V(s_t)]$$

where α_t is $\frac{1}{\#[\text{Returns}(s_t)]}$ at time t

- The update can be performed after return G_t is observed
- i.e. after the episode is completed

Recap: TD Prediction

- Unlike MC methods, TD methods wait only until **next time step**
- The simplest TD method (so called TD(0)) considers the update

$$V(S_t) \leftarrow V(S_t) + \alpha_t [R_t + \gamma V(S_{t+1}) - V(S_t)]$$

- This update rule has $R_t + \gamma V(S_{t+1})$ as the **target**
- Considered as a **bootstrap** method: update in part based on an existing estimate

Direct Estimator

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

Recap: Fitted Q-Iteration in Offline Setting

- Offline data: $\{(\mathbf{S}_{i,t}, \mathbf{A}_{i,t}, \mathbf{R}_{i,t}) : \mathbf{0} \leq t \leq \mathbf{T}, \mathbf{1} \leq i \leq \mathbf{N}\}$
- Fitted Q-Iteration can be naturally applied by repeating
 1. Compute \hat{Q} as the argmin of

$$\arg \min_Q \sum_t \left[\mathbf{R}_{i,t} + \gamma \max_a \tilde{Q}(\mathbf{S}_{i,t+1}, a) - Q(\mathbf{S}_{i,t}, \mathbf{A}_{i,t}) \right]^2$$

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

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

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

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

Stepwise IS Estimator [Zhang et al., 2013]

- Consider episodic task where T is the termination time
- Standard MC prediction is **not** applicable under **distributional shift**
- 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

- Stepwise IS is similar to an offline version of **MC**
- SIS learns from **complete** sequences
- SIS only works for **episodic** (terminating) environments
- Direct estimator (DE) is similar to an offline version of **TD**
- DE can learn from **incomplete** sequences
- DE works in **continuing** environments

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

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

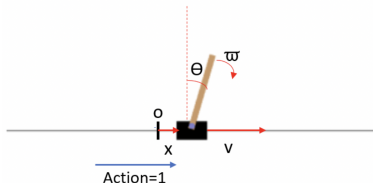
$$\mathbb{E}^{\pi} R_t = \mathbb{E}^{\mathbf{b}} \left[\frac{\pi(\mathbf{A}_0 | \mathbf{S}_0)}{\mathbf{b}(\mathbf{A}_0 | \mathbf{S}_0)} \cdots \frac{\pi(\mathbf{A}_t | \mathbf{S}_t)}{\mathbf{b}(\mathbf{A}_t | \mathbf{S}_t)} R_t \right]$$

- 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

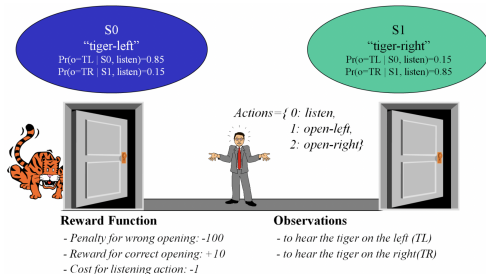
Pros & Cons of Direct v.s. Stepwise IS (Cont'd)

- 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)



Lecture Outline

1. Off-Policy Evaluation (OPE) Introduction
2. OPE in Contextual Bandits
3. OPE in Reinforcement Learning
- 4. State-of-the-Art OPE Algorithms**

Recap: RL Models

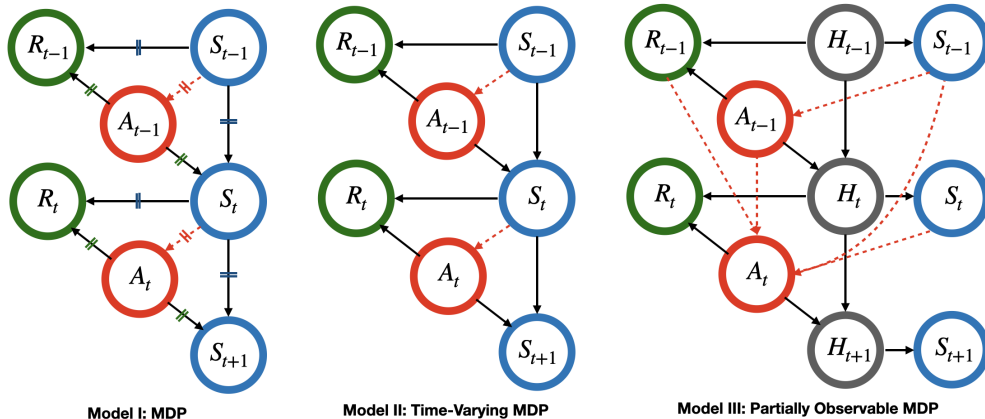


Figure: Causal diagrams for MDPs, TMDPs and POMDPs. Solid lines represent the causal relationships. Dashed lines indicate the information needed to implement the optimal policy. $\{H_t\}_t$ denotes latent variables. The parallel sign \parallel indicates that the conditional probability function given parent nodes is equal.

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)

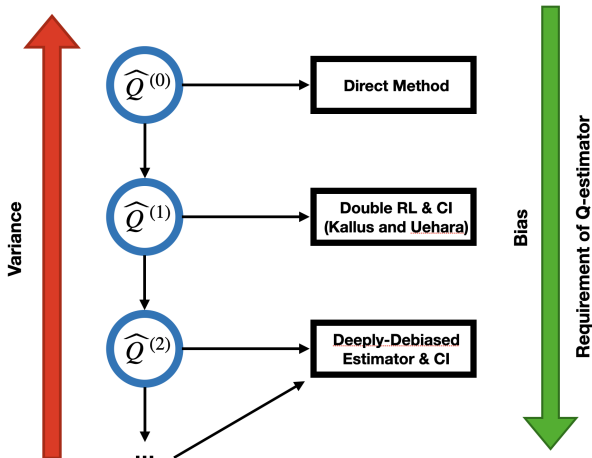
Double RL [Kallus and Uehara, 2019]

- 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

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



- Ensures the bias decays much faster than standard deviation
- Allows to provide valid **uncertainty quantification** (e.g., confidence interval)

Summary

- Off-policy evaluation
- Direct estimator
- Importance sampling estimator
- Doubly robust estimator
- Fitted Q-evaluation
- Stepwise IS/Marginalized IS
- Double reinforcement learning
- Deeply-debiased estimator

References I

- Joshua D Angrist, Guido W Imbens, and Donald B Rubin. Identification of causal effects using instrumental variables. *Journal of the American statistical Association*, 91(434): 444–455, 1996.
- 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 II

- 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.
- Judea Pearl. *Causality*. Cambridge university press, 2009.
- 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.
- Anastasios Tsiatis. *Semiparametric theory and missing data*. Springer Science & Business Media, 2007.
- 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.

References III

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.

Questions