

Intrinsically Efficient, Stable, and Bounded Off-Policy Evaluation for RL

Nathan Kallus , Masatoshi Uehara

Cornell University and Cornell Tech, Harvard University

Summary

We propose new estimators for OPE based on empirical likelihood that are always more efficient than IS (importance sampling), SNIS (self normalized importance sampling), and DR (doubly robust) and satisfy the same stability and boundedness properties as SNIS.

Introduction

Contextual Bandit Setting

- State space \mathcal{X} , Action space \mathcal{A}
- Reward space $[0, R_{\max}]$
- Evaluation Policy; $\pi_e : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$
- Behavior Policy; $\pi_b : \mathcal{X} \times \mathcal{A} \rightarrow [0, 1]$

We want to estimate $\beta_T^{\pi_e} = E_{\pi_e}[r]$ from n observation from the behavior policy $\{x^{(i)}, a^{(i)}, r^{(i)}\}_{i=1}^n$.

Existing approaches

Define

$$\omega = \pi_e(a|x)/\pi_b(a|x), \quad q(a, x) = E[r|a, x]$$

$$v(x) = E_{a \sim \pi_e(x)}[q(a, x)].$$

- **IS**; $\hat{\beta}_{\text{IS}} = E_n[\omega r]$
- **SIS**; $\hat{\beta}_{\text{SIS}} = E_n[\omega r]/E_n[\omega]$
- **DR**; $\hat{\beta}_{\text{dr}} = \hat{\beta}_d(\hat{q}(x, a))$, where $\hat{\beta}_d(m) = E_n[wr - \mathcal{F}(m)]$, $\mathcal{F}(m(x, a)) = wm(x, a) - \{\sum_{a \in \mathcal{A}} m(x, a)\pi_e(a|x)\}$.

Definition

- Local efficiency - When the model $q(x, a; \tau)$ is well-specified, the estimator achieves the efficiency bound.
- α -Boundedness - Bounded by αR_{\max} .
- Stability - If the conditional variance of $r^{(i)}$, given $\mathcal{D}_{x,a}$, is bounded by σ^2 , then the conditional variance of the estimator, given $\mathcal{D}_{x,a}$, is also bounded by σ^2 ($\mathcal{D}_{x,a} = \{(x^{(i)}, a^{(i)}) : i \leq n\}$).
- Intrinsic efficiency — The asymptotic MSE of the estimator is smaller than that of any of $\hat{\beta}_{\text{SIS}}, \hat{\beta}_{\text{IS}}, \hat{\beta}_{\text{SNIS}}, \hat{\beta}_{\text{dr}}$, irrespective of model specification.

REG

We let the REG estimator be $\hat{\beta}_{\text{reg}} = \hat{\beta}_d(\hat{\zeta}_1 + \hat{\zeta}_2 q(x, a; \hat{\tau}))$ where we choose the parameters by minimizing the estimated variance:

$$(\hat{\zeta}, \hat{\tau}) = \arg \min_{\zeta \in \mathbb{R}^2, \tau \in \Theta_\tau} E_n[\{wr - \mathcal{F}(\zeta_1 + \zeta_2 q(x, a; \tau))\}^2].$$

- When $\zeta_1 = 0, \zeta_2 = 1$, more doubly robust estimator (MDR).
- When $\zeta_2 = 1, \zeta_1 = 0$, IS estimator
- When $\zeta_2 = \beta^*, \zeta_1 = 0$, SIS estimator

EMP

We let the EMP estimator be

$$\hat{\beta}_{\text{emp}} = E_n[\hat{c}^{-1} \hat{\kappa}(x, a) \pi_e(a|x) r], \quad \text{where}$$

$$\hat{\kappa}(x, a) = \{\pi_b(a|x)[1 + \mathcal{F}(m(x, a; \hat{\xi}, \hat{\tau}))]\}^{-1},$$

$$\hat{c} = E_n[\{1 + \mathcal{F}(m(x, a; \hat{\xi}, \hat{\tau}))\}^{-1}],$$

$$\hat{\xi}, \hat{\tau} = \arg \max_{\xi \in \mathbb{R}, \tau \in \Theta_\tau} E_n[\log\{1 + \mathcal{F}(m(x, a; \xi, \tau))\}].$$

This is based on empirical likelihood

$$\max_{\kappa} \sum_{i=1}^n \log \kappa^{(i)}, \quad \text{s.t.} \quad \sum_{i=1}^n \kappa^{(i)} \pi_b(a^{(i)}|x^{(i)}) = 1,$$

$$\sum_{i=1}^n \kappa^{(i)} \pi_b(a^{(i)}|x^{(i)}) \mathcal{F}(m(x^{(i)}, a^{(i)}; \xi, \tau)) = 0.$$

The dual formulation derives the estimator.

Practical REG and EMP

We estimate a parameter τ in $q(x, a; \tau)$ as in DM to obtain $\hat{\tau}$, which we assume has a limit, $\hat{\tau} \xrightarrow{P} \tau^\dagger$. Then, we consider solving the following optimization problems for REG and EMP, respectively

$$\hat{\zeta} = \arg \min_{\zeta \in \mathbb{R}^2} E_n[\{wr - \mathcal{F}(m(x, a; \zeta, \hat{\tau}))\}^2],$$

$$\hat{\xi} = \arg \max_{\xi \in \mathbb{R}^2} E_n[\log\{1 + \mathcal{F}(m(x, a; \xi, \hat{\tau}))\}],$$

where $m(x, a; \zeta, \hat{\tau}) = \zeta_1 + \zeta_2 q(x, a; \hat{\tau})$ or $m(x, a; \xi, \hat{\tau}) = \xi_1 + \xi_2 q(x, a; \hat{\tau})$.

- Optimization becomes easier
- Desires properties are still kept

Theoretical results

Theorem; The estimators $\hat{\beta}_{\text{emp}}, \hat{\beta}_{\text{reg}}$ have local and intrinsic efficiency, and $\text{Asmse}[\hat{\beta}_{\text{reg}}] = \text{Asmse}[\hat{\beta}_{\text{emp}}]$

$$n^{-1} \min_{\zeta \in \mathbb{R}^2, \tau \in \mathbb{R}^{d_\tau}} E[\{wr - \mathcal{F}(\zeta + q(x, a; \tau))\}^2 - \beta^{*2}].$$

The estimator $\hat{\beta}_{\text{emp}}$ satisfies 1-boundedness and partial stability.

Theorem; The plug-in- τ versions of $\hat{\beta}_{\text{reg}}$ and $\hat{\beta}_{\text{emp}}$ still satisfy local and intrinsic efficiency, and $\hat{\beta}_{\text{emp}}$ satisfies 1-boundedness and partial stability. Their asymptotic MSEs are

$$n^{-1} \min_{\zeta \in \mathbb{R}^2} E[\{wr - \mathcal{F}(\zeta_1 + \zeta_2 q(x, a; \tau^\dagger))\}^2 - \beta^{*2}]. \quad (1)$$

Extension to RL

In RL, the goal is evaluating $E_{\pi_e}[\sum \gamma^t r_t]$. We can do the same technique. The intrinsic efficiency is defined as a superiority over

$$\hat{\beta}_{\text{snis}} = \frac{E_n[\omega_{0:T-1} \sum_{t=0}^{T-1} \gamma^t r_t]}{E_n[\omega_{0:T-1}]}, \quad \hat{\beta}_{\text{snis}} = E_n\left[\sum_{t=0}^{T-1} \frac{\omega_{0:t}}{E_n[\omega_{0:t}]} \gamma^t r_t\right]$$

$$\hat{\beta}_{\text{dr}} = \hat{\beta}_d(\{q(x, a; \hat{\tau})\}_{t=0}^{T-1}), \quad \text{where } \hat{\beta}_d(\{m_t\}_{t=0}^{T-1}) \text{ is}$$

$$E_n\left[\sum_{t=0}^{T-1} \gamma^t \omega_{0:t} r_t - \gamma^t (\omega_{0:t} m_t(x_t, a_t) - m_t(x_t, \pi_e))\right].$$

To derive REG, consider the class of estimators $\hat{\beta}_d(\{m_t\}_{t=0}^{T-1})$ where

$$m_t(x_t, a_t; \zeta) = \zeta_1 + \zeta_2 q(x_t, a_t; \hat{\tau}) \quad (0 \leq t \leq T-1).$$

Then, we define an estimator $\hat{\zeta}$ as

$$\hat{\zeta} = \arg \min_{\zeta \in \mathbb{R}^2} E_n[v(\{m_t(x_t, a_t; \zeta)\}_{t=0}^{T-1})]. \quad (2)$$

Same as EMP.

Experiment

We compare common methods with the proposed methods.

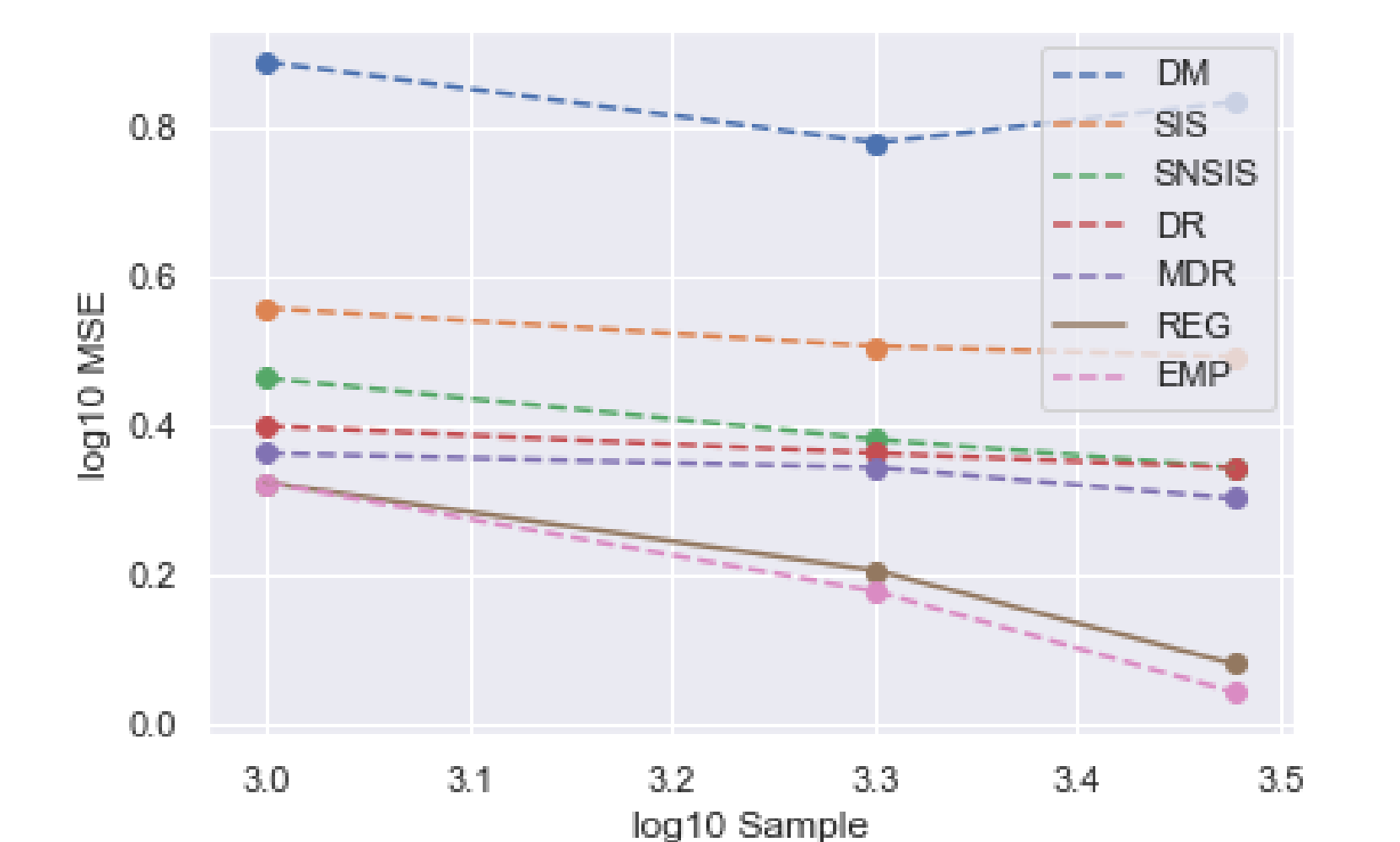


Figure 1:Cliff Walking

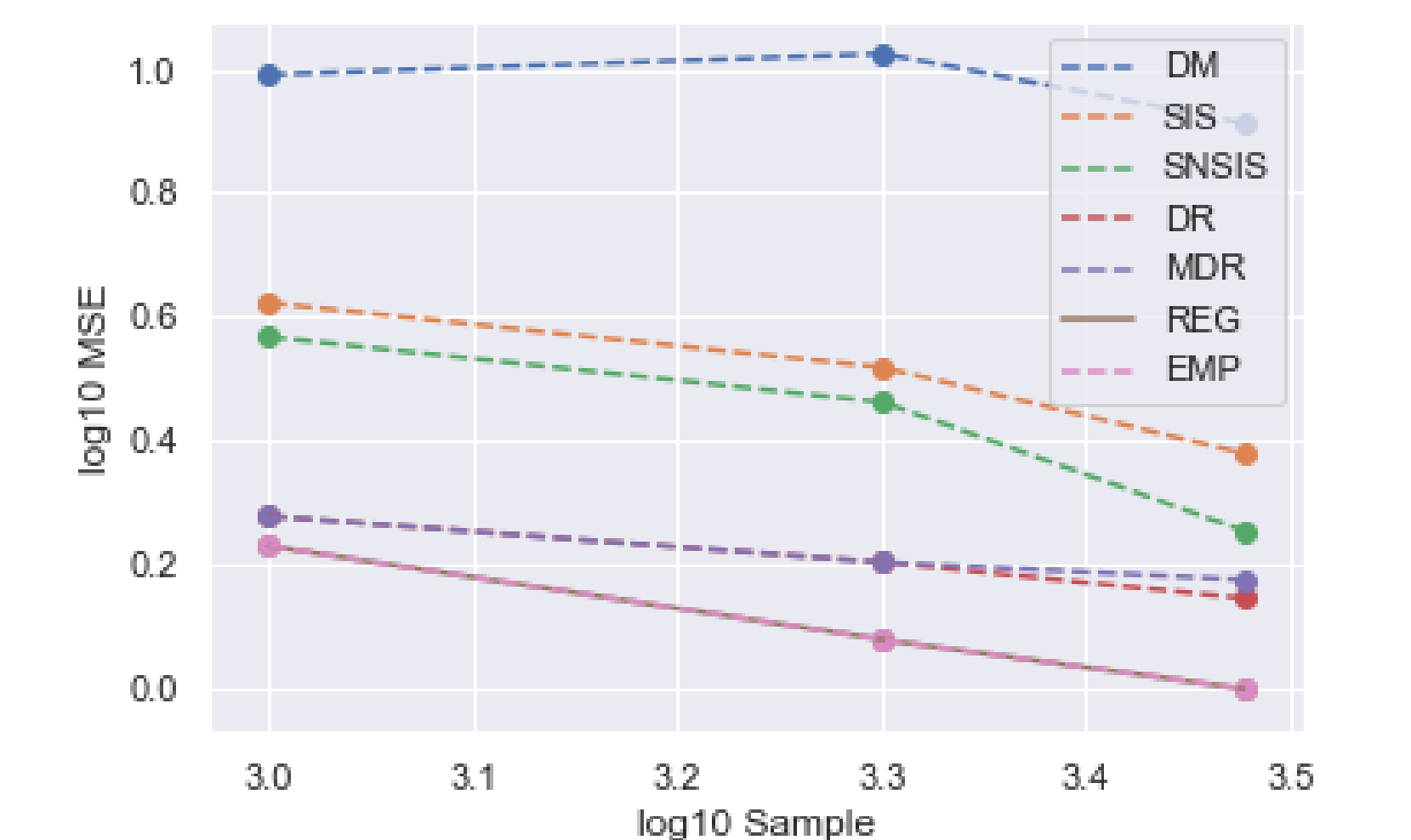


Figure 2:Mountain Car

Important Comparison

Table 1: Comparison of policy evaluation methods. The notation (*) means proposed estimator. The notation # means partially satisfied, as discussed in the text. (S)IS and SN(S)IS refer either to stepwise or non-stepwise.

	DM	(S)IS	SN(S)IS	DR	SNDR	MDR	REG(*)	SNREG(*)	EMP(*)
Consistency		✓	✓	✓	✓	✓	✓	✓	✓
Local efficiency	✓			✓	✓		✓	✓	✓
Intrinsic efficiency						#	✓	#	✓
Boundedness	1		1		2			2	1
Stability	#		✓		#				✓