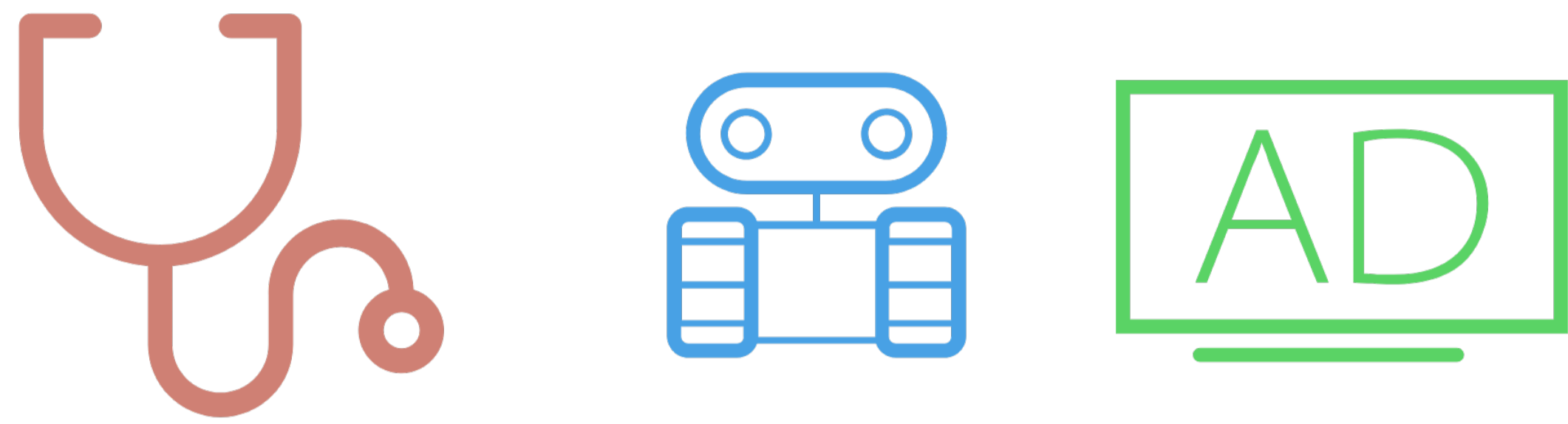


Infinite-Horizon Off-Policy Evaluation

- **Problem:** Evaluate a target policy π given data from behavior policy π_0 .
- **Wide Applications:** whenever evaluating new policies is costly or impossible, due to **high cost, risk, or ethics, legal concerns**.



Healthcare Robotic & Control Recommendation

- **Setup:** given behavior trajectories $\{s_i, a_i, s'_i, r_i\} \sim \pi_0$, we want to estimate the average discounted reward of π

$$R^\pi = \mathbb{E}_{\tau \sim \pi} \left[\frac{\sum_{t=0}^{\infty} \gamma^t r_t}{\sum_{t=0}^{\infty} \gamma^t} \right] = (1 - \gamma) \mathbb{E}_{\tau \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t r_t \right].$$

Bellman Equation

- **Value function:** Define $\mathcal{P}^\pi f(s) = \sum_{a,s'} \pi(a|s) T(s'|s, a) f(s')$, we have the Bellman equation for value function

$$V^\pi = r^\pi + \gamma \mathcal{P}^\pi V^\pi.$$

- **Density function:** Define $\mathcal{T}^\pi f(s') = \sum_{a,s} \pi(a|s) T(s'|s, a) f(s)$, we have the Bellman equation for density function [1]

$$d_\pi = (1 - \gamma) \mu_0 + \gamma \mathcal{T}^\pi d_\pi.$$

Basic Estimators

- **Estimation via Value Function (Direct Methods):**

$$R^\pi = (1 - \gamma) \mathbb{E}_{s \sim \mu_0} [V^\pi(s)].$$

Given a learned \hat{V} , we have the following value function estimator

$$\hat{R}_{\text{VAL}}^\pi[\hat{V}] = \frac{(1 - \gamma)}{n_0} \sum_{i=1}^{n_0} \hat{V}(s_0^{(i)}).$$

- **Estimation via State Density Function (IS methods)[1]:**

$$R^\pi = \mathbb{E}_{s,a \sim d_{\pi_0}} \left[w_{\pi/\pi_0}(s) \frac{\pi(a|s)}{\pi_0(a|s)} r(s, a) \right],$$

where $w_{\pi/\pi_0}(s) = \frac{d_\pi(s)}{d_{\pi_0}(s)}$ is the state density ratio function.

Given a learned density ratio \hat{w} , we have the IS estimator

$$\hat{R}_{\text{SIS}}^\pi[\hat{w}] = \frac{1}{n} \sum_{i=1}^n \hat{w}(s_i) \frac{\pi(a_i|s_i)}{\pi_0(a_i|s_i)} r_i.$$

Doubly Robust Estimator

- Combine the estimators $R_{\text{SIS}}^\pi[\hat{w}]$ and $R_{\text{VAL}}^\pi[\hat{V}]$, we get the "doubly robust" estimator

$$R_{\text{DR}}^\pi[\hat{V}, \hat{w}] = \underbrace{\sum_s r^\pi(s) d_{\pi_0}(s) \hat{w}(s)}_{R_{\text{SIS}}^\pi[\hat{w}]} + (1 - \gamma) \underbrace{\sum_s \hat{V}(s) \mu_0(s)}_{R_{\text{VAL}}^\pi[\hat{V}]} - \underbrace{\sum_s (\hat{V}(s) - \gamma \mathcal{P}^\pi \hat{V}(s)) d_{\pi_0}(s) \hat{w}(s)}_{R_{\text{bridge}}^\pi[\hat{V}, \hat{w}]}.$$

- **Theorem: Double Robustness**

$$R_{\text{DR}}^\pi[\hat{V}, \hat{w}] - R^\pi = \mathbb{E}_{s \sim d_{\pi_0}} [\varepsilon_{\hat{w}}(s) \varepsilon_{\hat{V}}(s)],$$

$$\varepsilon_{\hat{w}} = \frac{d_\pi(s)}{d_{\pi_0}(s)} - \hat{w}(s), \quad \varepsilon_{\hat{V}}(s) = \hat{V}(s) - r^\pi(s) - \gamma \mathcal{P}^\pi \hat{V}(s).$$

Double Robustness and Lagrangian Duality

Primal Problem

$$R^\pi = \min_V \left\{ (1 - \gamma) \sum_s \mu_0(s) V(s) \quad \text{s.t.} \quad V(s) \geq r^\pi(s) + \gamma \mathcal{P}^\pi V(s), \quad \forall s \right\}$$

↓ Lagrangian multiplier $\rho(s) \geq 0$

Lagrangian function

$$L(V, \rho) = \begin{cases} (1 - \gamma) \sum_s \mu_0(s) V(s) + \sum_s \rho(s) r^\pi(s) - \sum_s \rho(s) \cdot (I - \gamma \mathcal{P}^\pi) V(s) \\ \mathcal{P}^\pi \text{ and } \mathcal{T}^\pi \text{ are self-adjoint!} \\ \sum_s \rho(s) r^\pi(s) + (1 - \gamma) \sum_s \mu_0(s) V(s) - \sum_s (I - \gamma \mathcal{T}^\pi) \rho(s) \cdot V(s) \end{cases}$$

$L(V, \rho)$ is the DR estimator!

$$w_{\rho/\pi_0} = \rho(s)/d_{\pi_0}, \quad R_{\text{DR}}^\pi[V, w_{\rho/\pi_0}] = L(V, \rho)$$

↑ Lagrangian multiplier $V(s)$

Dual Problem

$$R^\pi = \max_{\rho \geq 0} \left\{ \sum_s \rho(s) r^\pi(s) \quad \text{s.t.} \quad \rho(s') = (1 - \gamma) \mu_0(s') + \gamma \mathcal{T}^\pi \rho(s'), \quad \forall s' \right\}$$

- **Application I: Extension to $\gamma = 1$ (Average Case)**

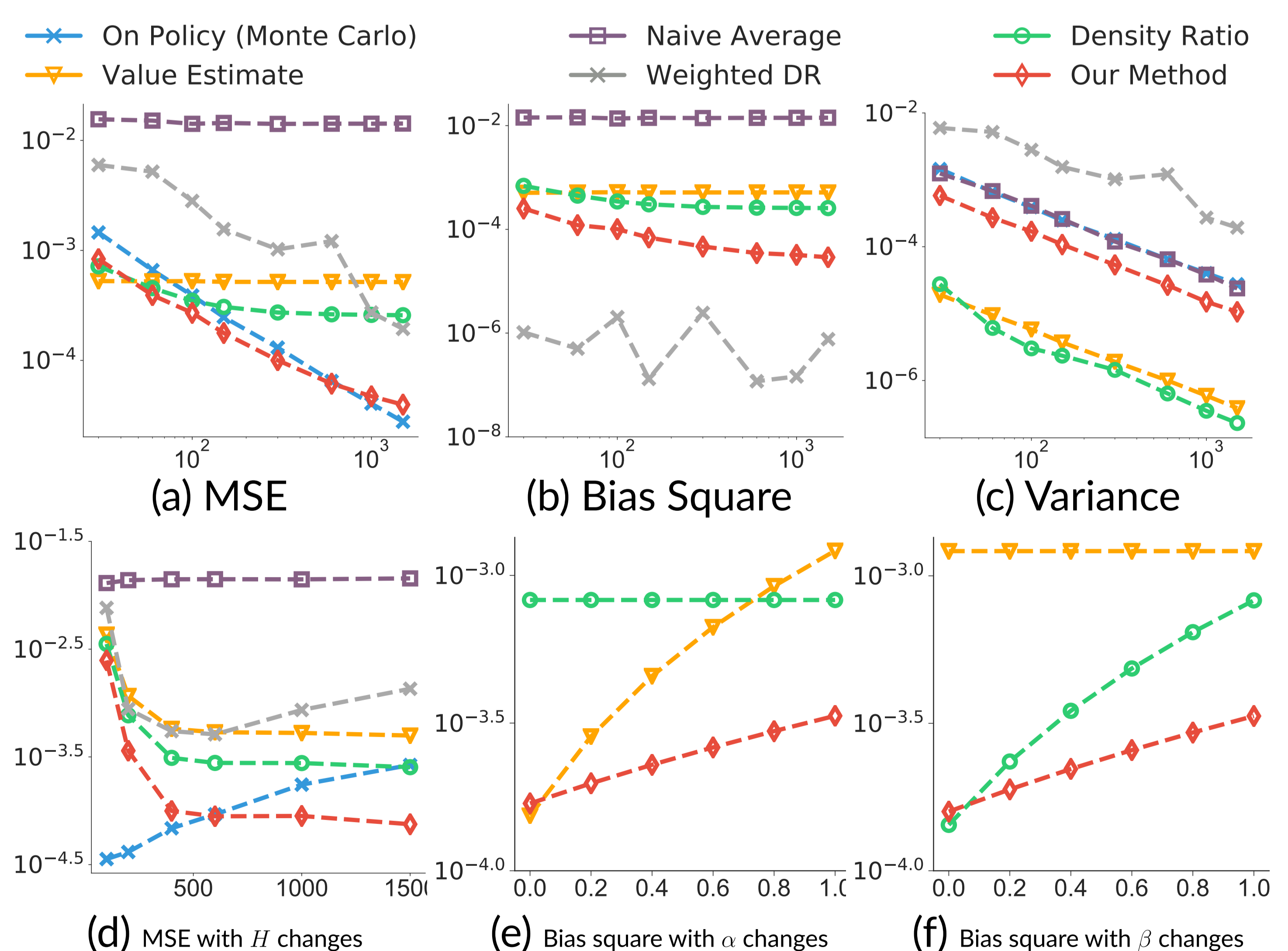
$$R_{\text{DR}}^\pi[\hat{V}, \hat{w}] := \frac{\mathbb{E}_{s \sim d_{\pi_0}} [\hat{w}(s) (r^\pi(s) - \hat{V}(s) + \mathcal{P}^\pi \hat{V}(s))]}{\mathbb{E}_{s \sim d_{\pi_0}} [\hat{w}(s)]}.$$

- **Application II: Extension to Q-value functions:**

$$R_{\text{DR}}^\pi[Q, x] = \underbrace{(1 - \gamma) \mathbb{E}_{s \sim \mu_0, a \sim \pi(\cdot|s)} [Q(s, a)]}_{R_{\text{VAL}}^\pi[Q]} + \underbrace{\mathbb{E}_{s, a \sim \mathcal{D}} [x(s, a) r(s, a)]}_{R_{\text{SAIS}}^\pi[x]} - \underbrace{\mathbb{E}_{s, a, s' \sim \mathcal{D}, a' \sim \pi(\cdot|s')} [x(s, a) (Q(s, a) - \gamma Q(s', a'))]}_{R_{\text{bridge}}^\pi[Q, x]}.$$

Experimental Results

Taxi Environment



- **Related DR method:** contextual bandit [2], finite-horizon RL[3], concurrent work [4].
- **Acknowledgment:** This work is supported in part by NSF CRII 1830161 and NSF CAREER 1846421.

