

Offline Policy Selection under Uncertainty

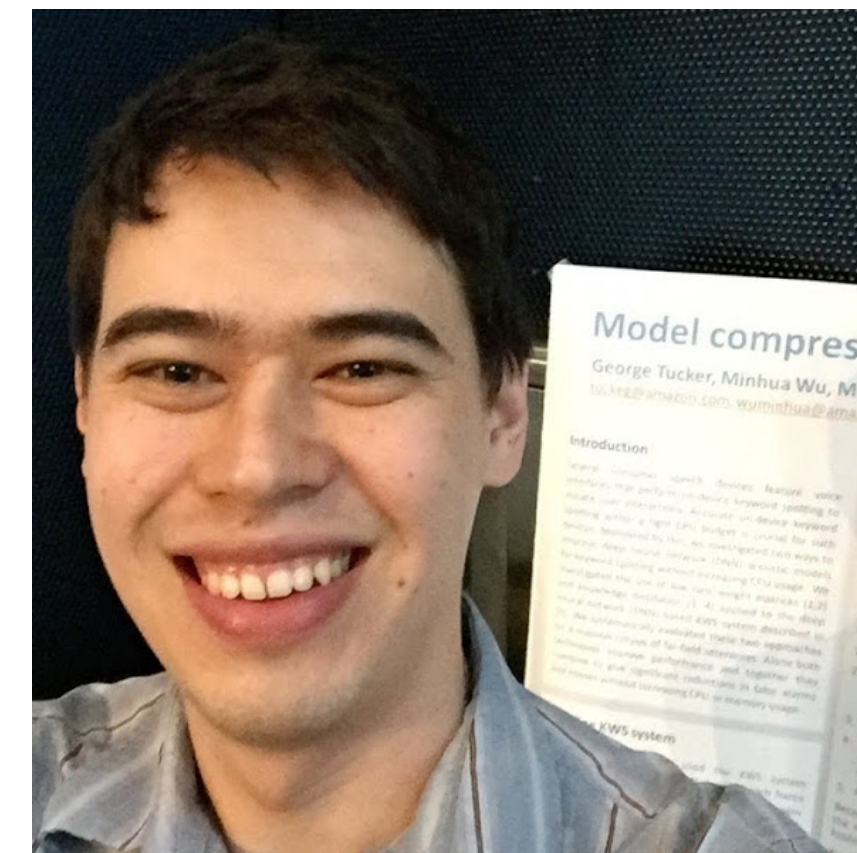
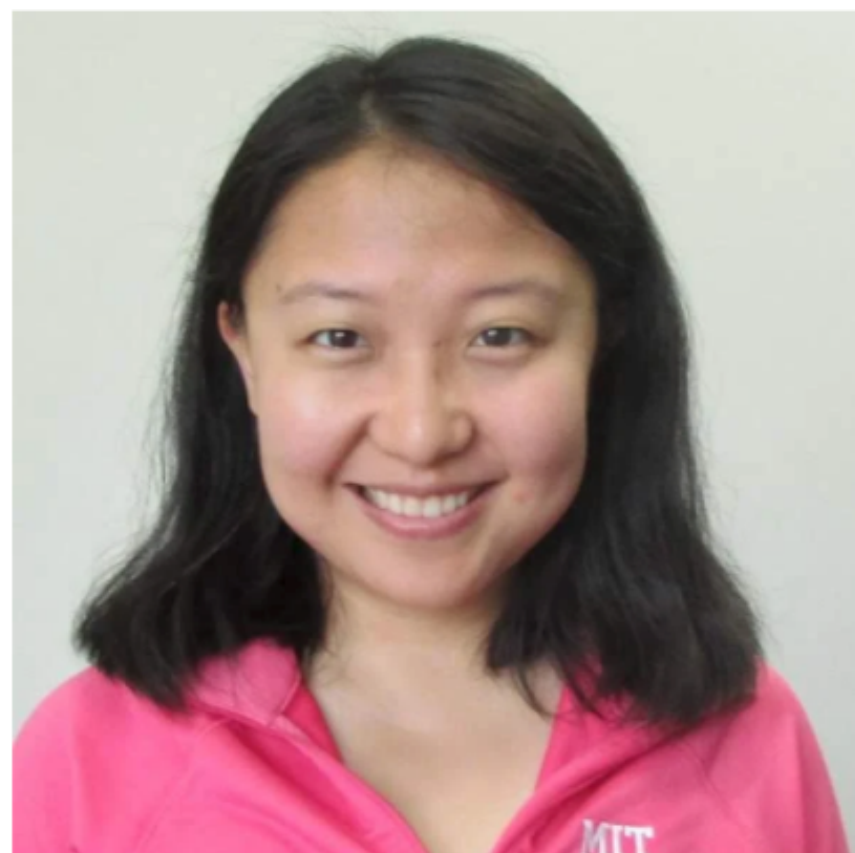
Sherry Yang*

Bo Dai*

Ofir Nachum*

George Tucker

Dale Schuurmans



Paper: <https://arxiv.org/abs/2012.06919>

Code: https://github.com/google-research/dice_rl

Offline Policy Selection

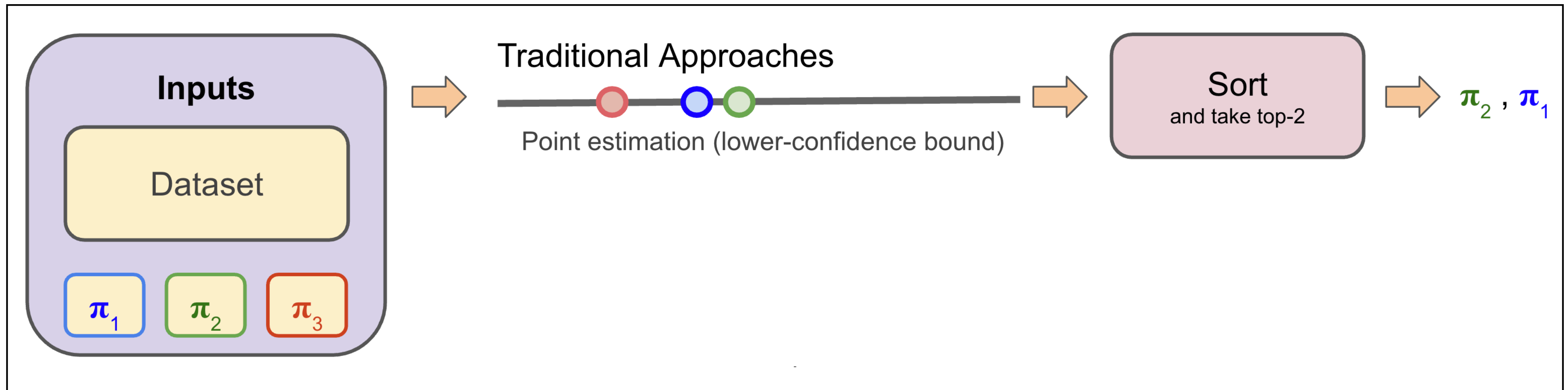
Offline policy selection:

- Compute a ranking $O \in \text{Perm}([1, N])$ over $\{\pi_i\}_{i=1}^N$ given a fixed dataset D according to some utility function u : $\mathcal{O} \leftarrow \text{ArgSortDescending}(\{u(\pi_i)\}_{i=1}^N)$

Offline Policy Selection

Offline policy selection:

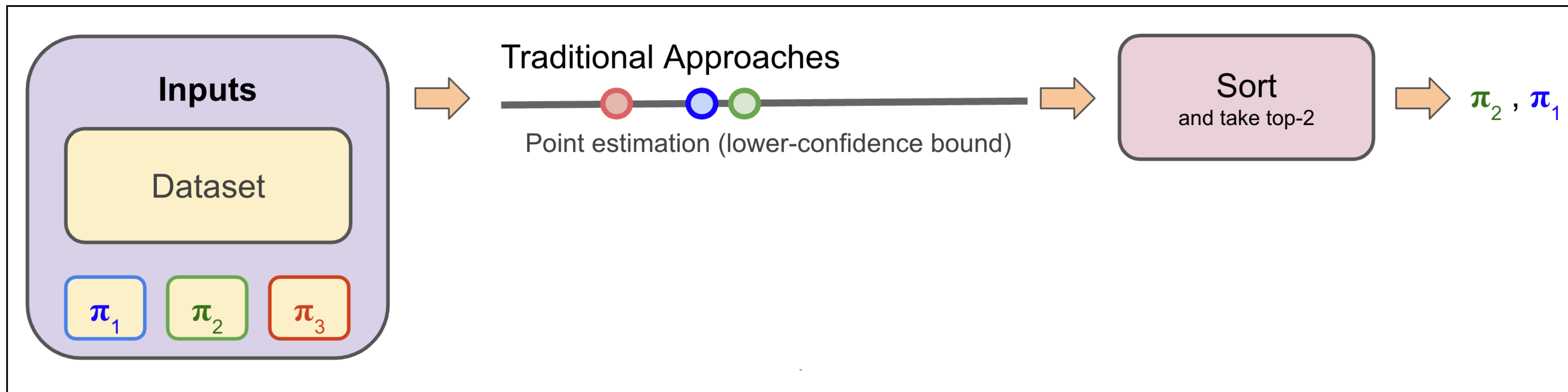
- Compute a ranking $O \in \text{Perm}([1, N])$ over $\{\pi_i\}_{i=1}^N$ given a fixed dataset D according to some utility function u : $\mathcal{O} \leftarrow \text{ArgSortDescending}(\{u(\pi_i)\}_{i=1}^N)$



Offline Policy Selection

Offline policy selection:

- Compute a ranking $O \in \text{Perm}([1, N])$ over $\{\pi_i\}_{i=1}^N$ given a fixed dataset D according to some utility function u : $\mathcal{O} \leftarrow \text{ArgSortDescending}(\{u(\pi_i)\}_{i=1}^N)$

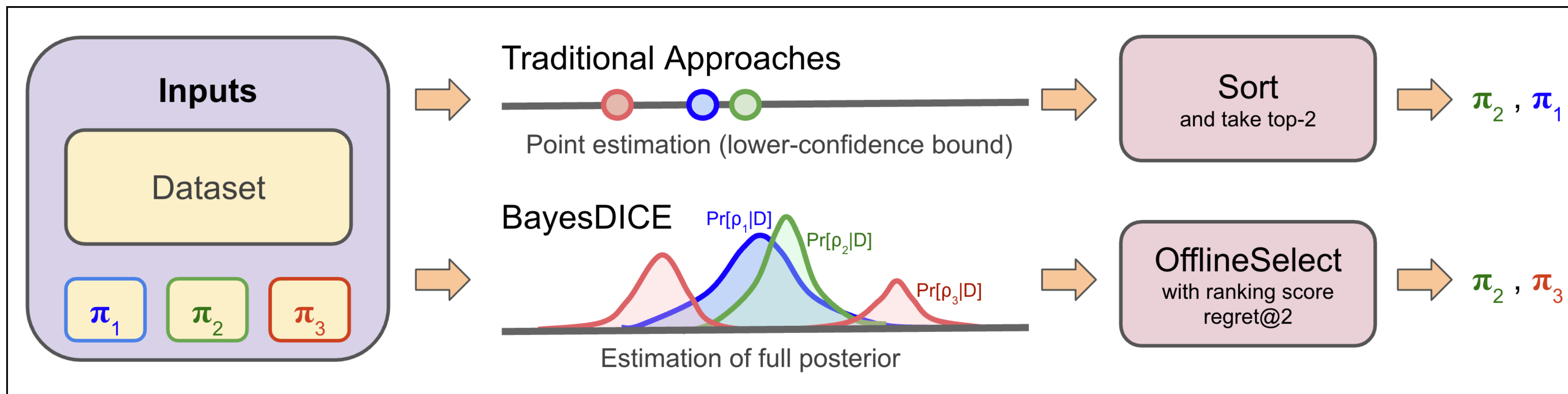


- Practical ranking criteria: top-k precision, top-k accuracy, top-k regret, top-k correlation, CVaR, ...

Offline Policy Selection

Offline policy selection:

- Compute a ranking $O \in \text{Perm}([1, N])$ over $\{\pi_i\}_{i=1}^N$ given a fixed dataset D according to some utility function u : $O \leftarrow \text{ArgSortDescending}(\{u(\pi_i)\}_{i=1}^N)$



- Practical ranking criteria: top-k precision, top-k accuracy, top-k regret, top-k correlation, CVaR, ...

BayesDICE

Recall off-policy evaluation:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\pi} [R(s,a)] \text{ where } \mathcal{P}_*^\pi d^\pi(s,a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$$

BayesDICE

Recall off-policy evaluation:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\pi} [R(s,a)] \text{ where } \mathcal{P}_*^\pi d^\pi(s,a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$$

DICE point estimator:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\mathcal{D}} [\zeta^*(s,a) \cdot R(s,a)] \text{ where } \zeta^*(s,a) := \frac{d^\pi(s,a)}{d^\mathcal{D}(s,a)}$$

[1] Nachum, et al. [Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections.](#)

[2] Yang, et al. [Off-policy evaluation via the regularized lagrangian.](#)

BayesDICE

Recall off-policy evaluation:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\pi} [R(s,a)] \text{ where } \mathcal{P}_*^\pi d^\pi(s,a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$$

DICE point estimator:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\mathcal{D}} [\zeta^*(s,a) \cdot R(s,a)] \text{ where } \zeta^*(s,a) := \frac{d^\pi(s,a)}{d^\mathcal{D}(s,a)}$$

BayesDICE learns $q(\zeta^\pi | \mathcal{D}) \propto p(\mathcal{D} | \zeta^\pi) p(\zeta^\pi)$:

- By optimizing $\min_{q \in \mathcal{P}} -\mathbb{E}_{q(\zeta^\pi)} [\log p(\mathcal{D} | \zeta^\pi)] + KL(q||p)$

[1] Nachum, et al. [Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections.](#)

[2] Yang, et al. [Off-policy evaluation via the regularized lagrangian.](#)

BayesDICE

Recall off-policy evaluation:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\pi} [R(s,a)] \text{ where } \mathcal{P}_*^\pi d^\pi(s,a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$$

DICE point estimator:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\mathcal{D}} [\zeta^*(s,a) \cdot R(s,a)] \text{ where } \zeta^*(s,a) := \frac{d^\pi(s,a)}{d^\mathcal{D}(s,a)}$$

BayesDICE learns $q(\zeta^\pi | \mathcal{D}) \propto p(\mathcal{D} | \zeta^\pi) p(\zeta^\pi)$:

- By optimizing $\min_{q \in \mathcal{P}} -\mathbb{E}_{q(\zeta^\pi)} [\log p(\mathcal{D} | \zeta^\pi)] + KL(q||p)$

constraints

[1] Nachum, et al. [Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections.](#)

[2] Yang, et al. [Off-policy evaluation via the regularized lagrangian.](#)

BayesDICE

Recall off-policy evaluation:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\pi} [R(s,a)] \text{ where } \mathcal{P}_*^\pi d^\pi(s,a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$$

DICE point estimator:

$$\rho(\pi) = \mathbb{E}_{(s,a) \sim d^\mathcal{D}} [\zeta^*(s,a) \cdot R(s,a)] \text{ where } \zeta^*(s,a) := \frac{d^\pi(s,a)}{d^\mathcal{D}(s,a)}$$

BayesDICE learns $q(\zeta^\pi | \mathcal{D}) \propto p(\mathcal{D} | \zeta^\pi) p(\zeta^\pi)$:

- By optimizing $\min_{q \in \mathcal{P}} -\mathbb{E}_{q(\zeta^\pi)} [\log p(\mathcal{D} | \zeta^\pi)] + KL(q||p)$
- Posterior regularization:

$$\min_q \xi + KL(q||p) \quad \text{s.t.} \quad q \in \mathcal{P} \cap \{\xi = -\mathbb{E}_{q(\zeta^\pi)} [\log p(\mathcal{D} | \zeta^\pi)]\}$$

constraints

[1] Nachum, et al. [Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections.](#)

[2] Yang, et al. [Off-policy evaluation via the regularized lagrangian.](#)

BayesDICE

Constraint embeddings:

- Density constraints: $\mathcal{P}_*^\pi d^\pi(s, a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$
Equivalently: $\Delta_d(s, a) := (1 - \gamma)\mu_0(s)\pi(a|s) + \gamma \cdot \mathcal{P}_*^\pi d(s, a) - d(s, a) = 0$

[1] Farias and Roy. [The linear programming approach to approximate dynamic programming.](#)

[2] Small, et al. [A Hilbert space embedding for distributions.](#)

[3] Nachum, et al. [Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections.](#)

BayesDICE

Constraint embeddings:

- Density constraints: $\mathcal{P}_*^\pi d^\pi(s, a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$
Equivalently: $\Delta_d(s, a) := (1 - \gamma)\mu_0(s)\pi(a|s) + \gamma \cdot \mathcal{P}_*^\pi d(s, a) - d(s, a) = 0$
- Introduce function space embedding ϕ :
$$\langle \phi, \Delta_d \rangle := \mathbb{E}_{(1-\gamma)\mu_0(s)\pi(a|s) + \gamma \cdot \mathcal{P}_*^\pi d(s, a)} [\phi(s, a)] - \mathbb{E}_{d(s, a)} [\phi(s, a)]$$

[1] Farias and Roy. [The linear programming approach to approximate dynamic programming.](#)

[2] Small, et al. [A Hilbert space embedding for distributions.](#)

[3] Nachum, et al. [Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections.](#)

BayesDICE

Constraint embeddings:

- Density constraints: $\mathcal{P}_*^\pi d^\pi(s, a) := \pi(a|s) \sum_{\tilde{s}, \tilde{a}} T(s|\tilde{s}, \tilde{a}) d^\pi(\tilde{s}, \tilde{a})$
Equivalently: $\Delta_d(s, a) := (1 - \gamma)\mu_0(s)\pi(a|s) + \gamma \cdot \mathcal{P}_*^\pi d(s, a) - d(s, a) = 0$
- Introduce function space embedding ϕ :
$$\langle \phi, \Delta_d \rangle := \mathbb{E}_{(1-\gamma)\mu_0(s)\pi(a|s) + \gamma \cdot \mathcal{P}_*^\pi d(s, a)} [\phi(s, a)] - \mathbb{E}_{d(s, a)} [\phi(s, a)]$$
- BayesDICE objective:

$$\min_{q \in \mathcal{P}} -\mathbb{E}_{q(\zeta^\pi)} [\log p(\mathcal{D}|\zeta^\pi)] + KL(q||p)$$

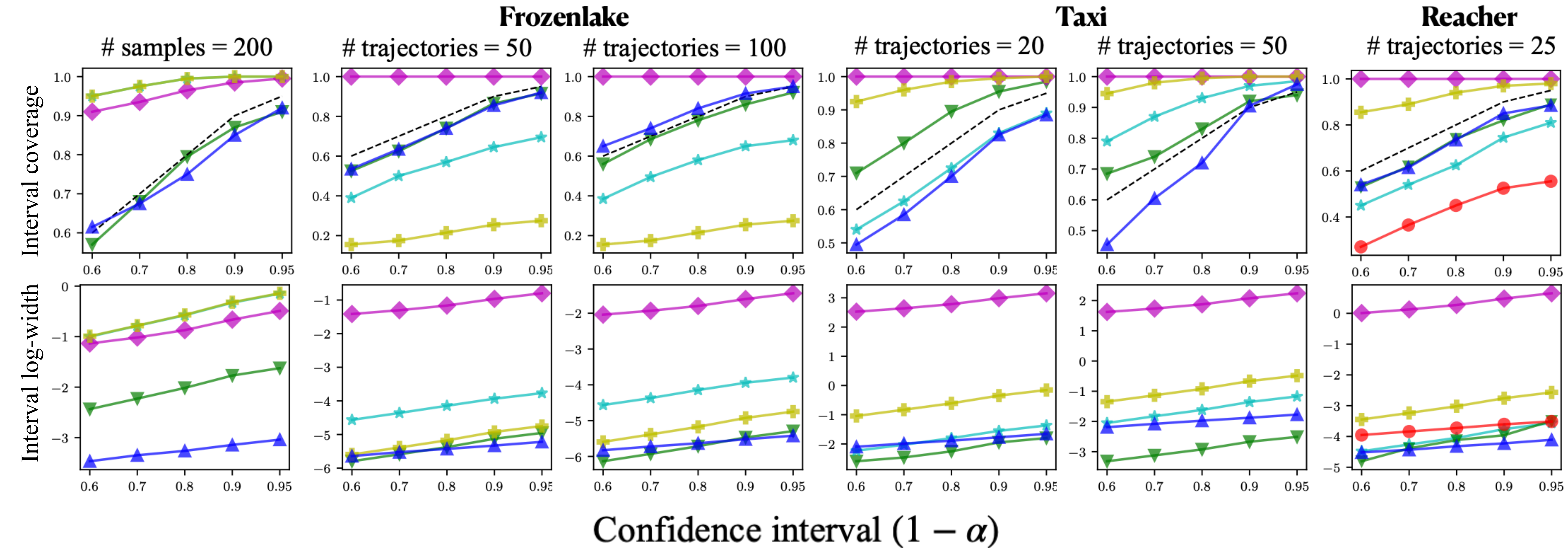
$$\min_q \frac{\lambda}{\epsilon} \mathbb{E}_q [\ell(\zeta, \mathcal{D})] + KL(q||p) \text{ where } \ell(\zeta, \mathcal{D}) = \langle \phi, \Delta_d \rangle^\top \langle \phi, \Delta_d \rangle$$
$$= \max_{\beta \in \mathcal{H}_\phi} \beta^\top \langle \phi, \Delta_d \rangle - \beta^\top \beta$$

[1] Farias and Roy. [The linear programming approach to approximate dynamic programming.](#)

[2] Small, et al. [A Hilbert space embedding for distributions.](#)

[3] Nachum, et al. [Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections.](#)

Experiments: CI Estimation



Experiments: Policy Selection

$\hat{\mathcal{S}} =$ ● Mean ▲ Mean - Std ▼ Mean + Std ---- Bayes-Optimal + \mathcal{S} ■ \mathcal{S}

Bandit

$\mathcal{S} = \text{Acc. Top-2}$

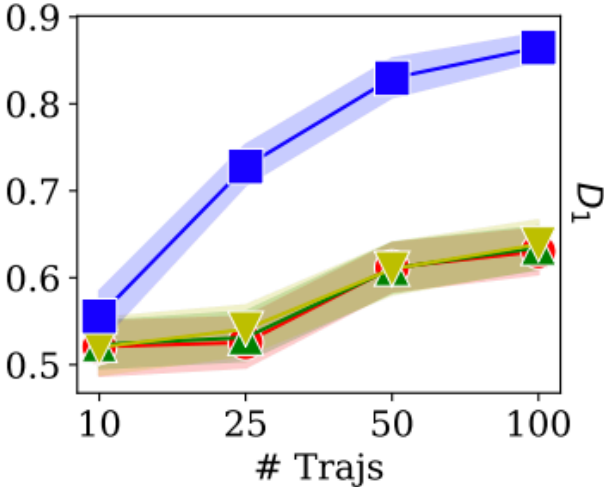
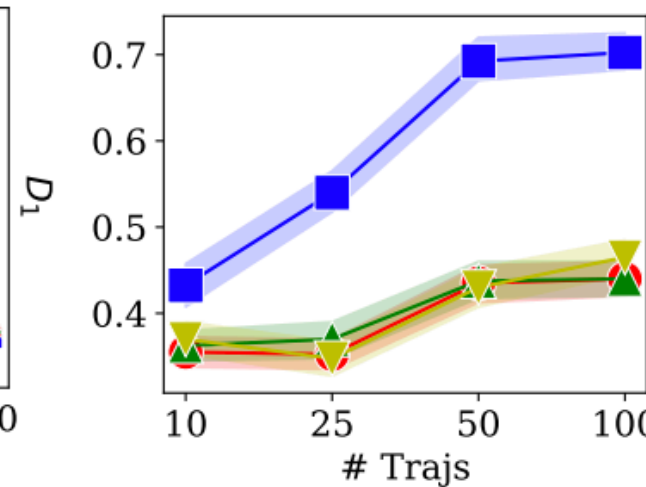
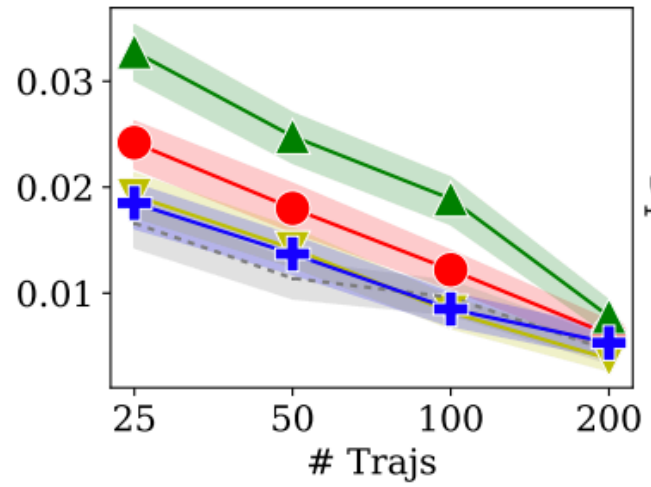
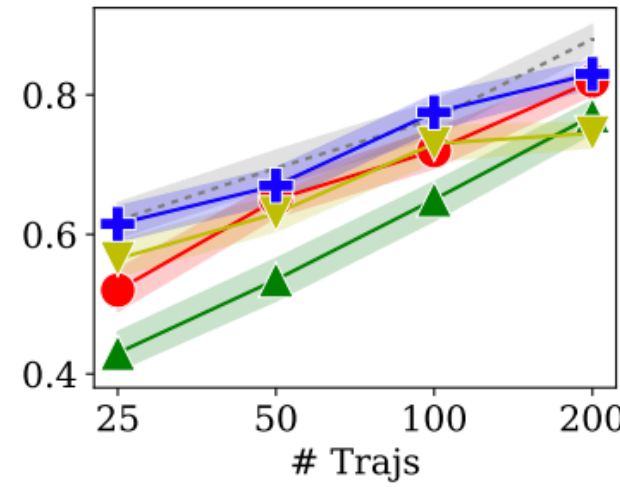
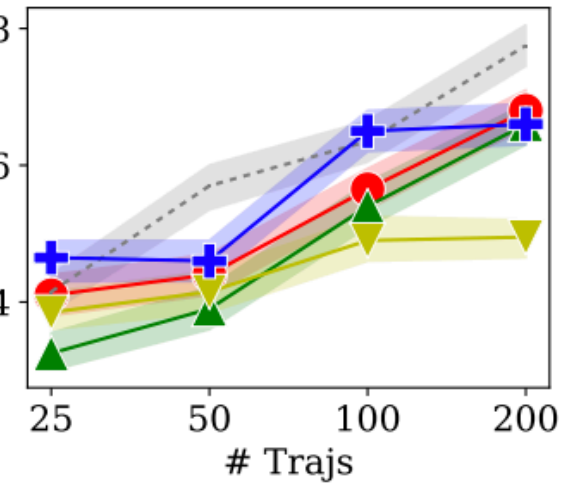
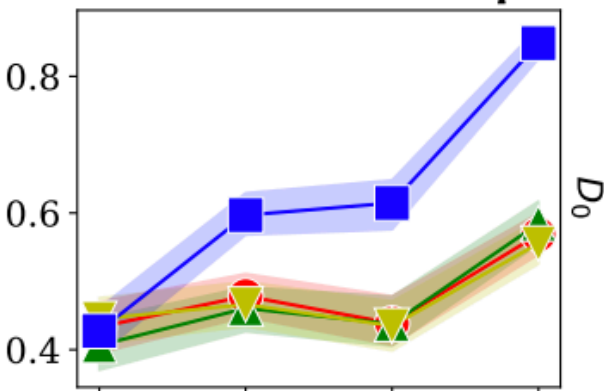
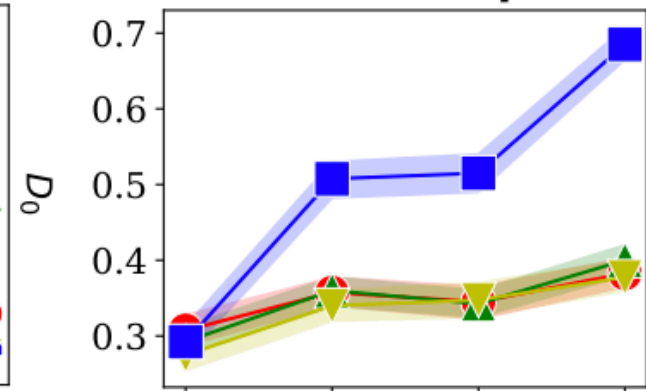
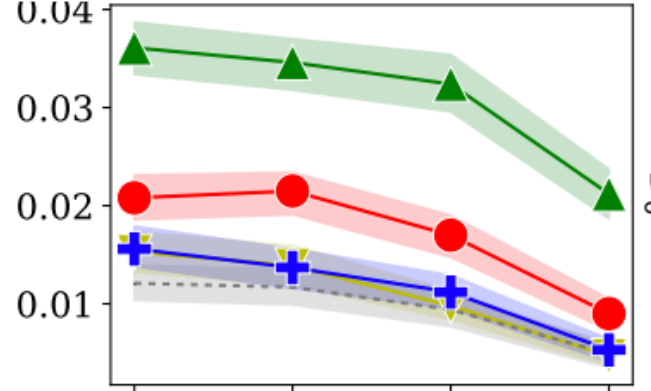
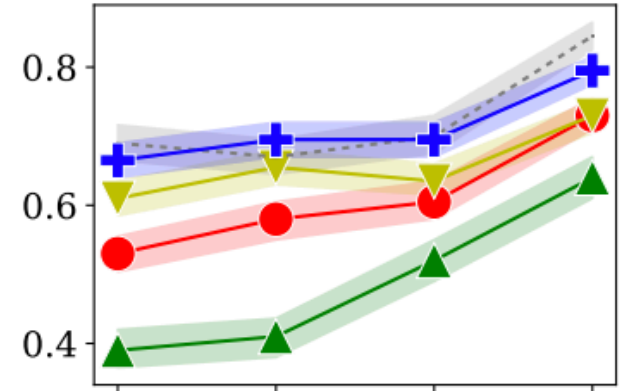
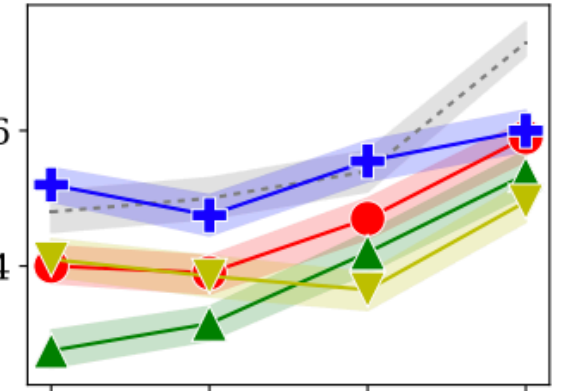
$\mathcal{S} = \text{Precision Top-2}$

$\mathcal{S} = \text{Regret Top-2}$

$\mathcal{S} = \text{Acc. Top-4}$

$\mathcal{S} = \text{Correlation Top-4}$

Reacher



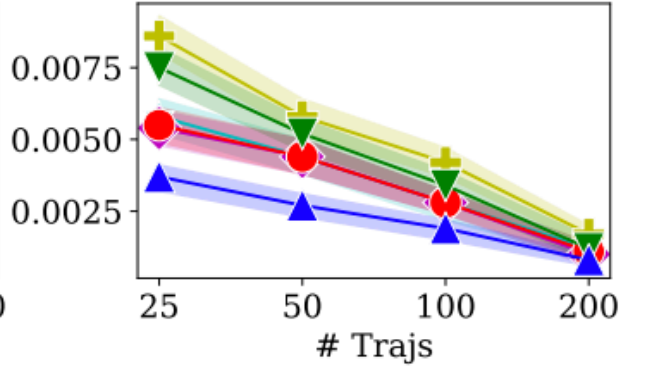
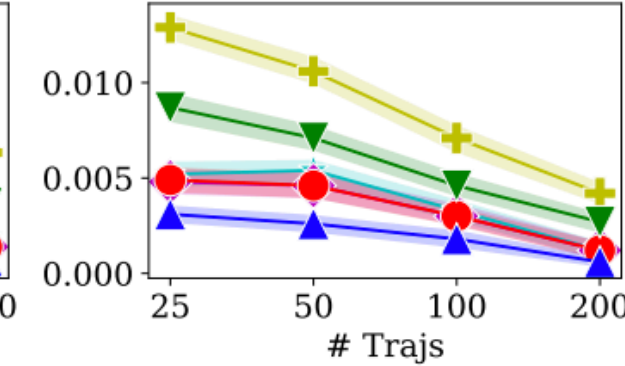
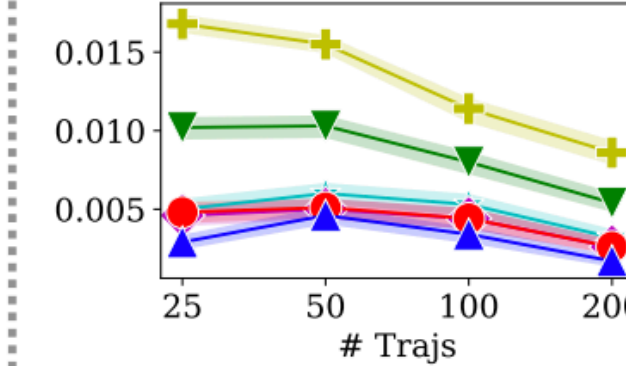
★ Bootstrap ◆ Bernstein + Student-t
● DualDICE ▼ CoinDICE ▲ BayesDICE

Bandit

Regret Top-3, $\alpha = 0.2$

Regret Top-3, $\alpha = 0.4$

Regret Top-3, $\alpha = 0.6$

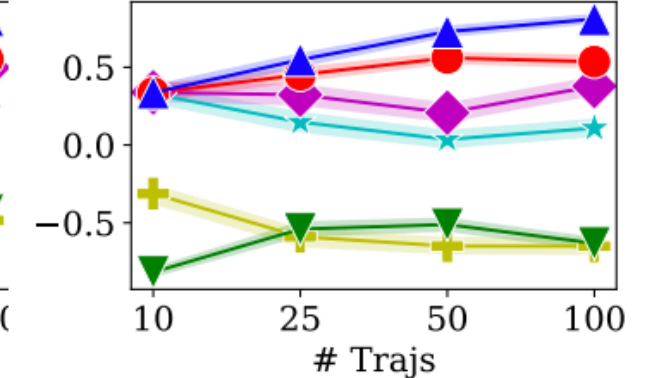
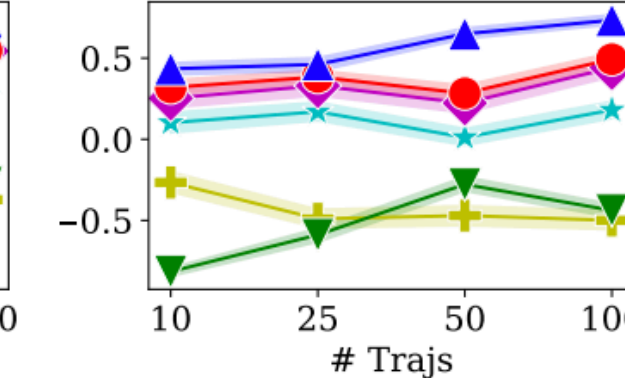
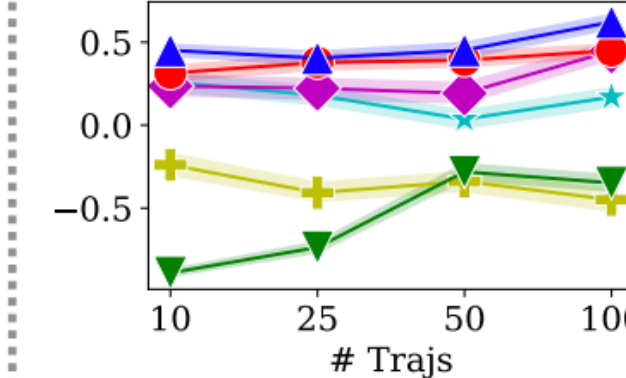


Reacher

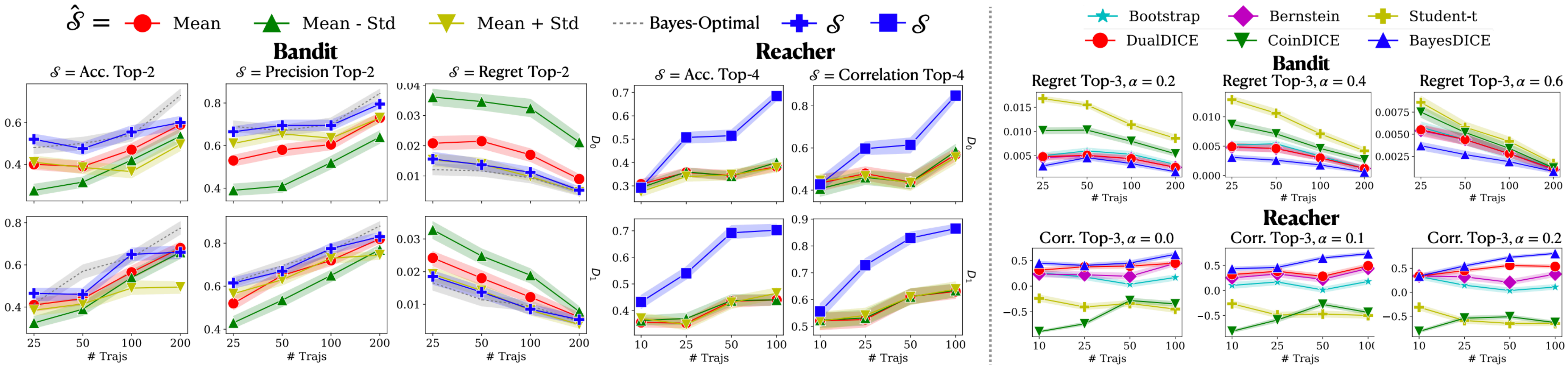
Corr. Top-3, $\alpha = 0.0$

Corr. Top-3, $\alpha = 0.1$

Corr. Top-3, $\alpha = 0.2$



Experiments: Policy Selection



Thank you. Checkout

Paper: <https://arxiv.org/pdf/2012.06919.pdf>

Code: https://github.com/google-research/google-research/tree/master/rl_repr