

Importance Resampling:

Conclusions and Future Perspectives

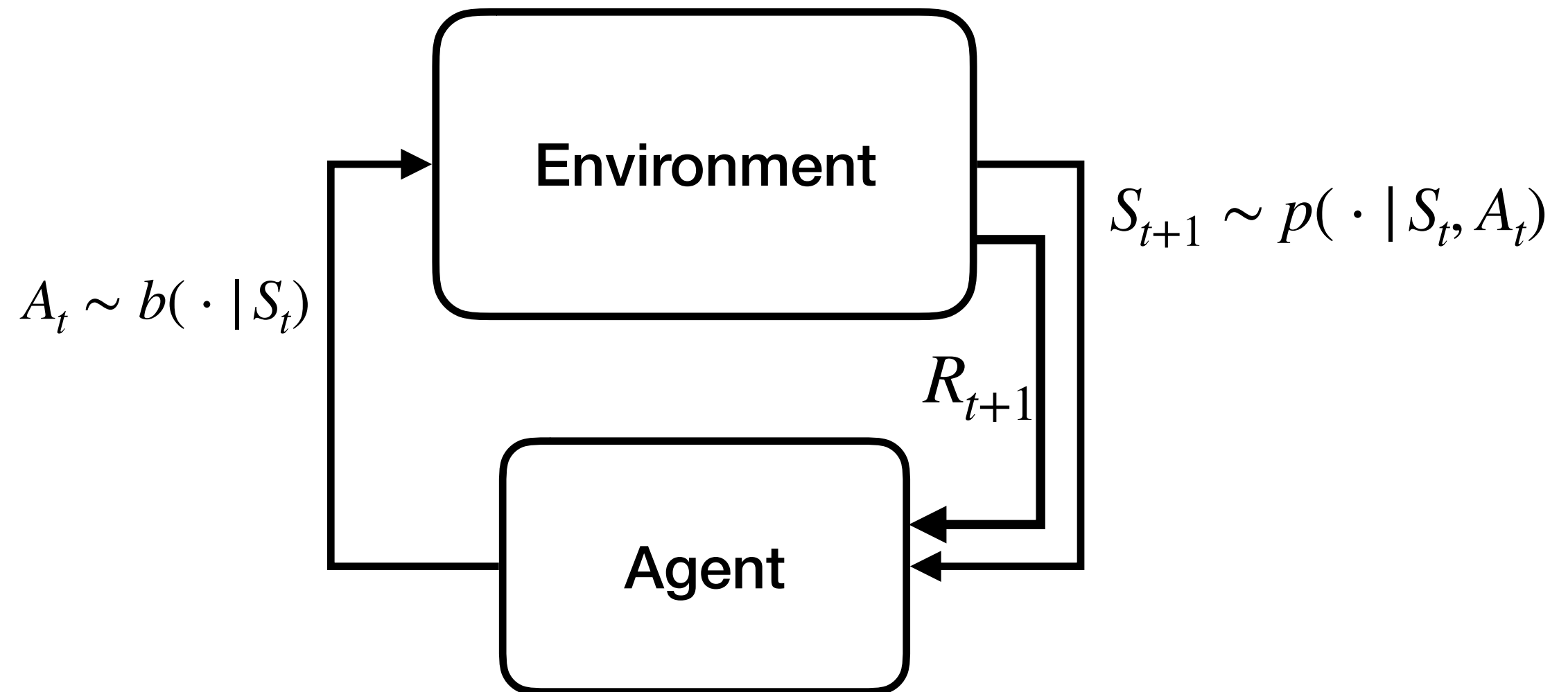
Matthew Schlegel, Wes Chung, Jian Qian,
Daniel Graves, and Martha White

Goal

Give an overview of our conclusions from exploring resampling for off-policy prediction.

Outline

- Background
- Reweighting Vs Resampling
- Empirical Results
- Conclusions
- **Future directions and perspectives**



Value Function

$$v(S_t) = \mathbb{E}_{\pi} \left[\sum_{j=t}^{\infty} \left(\prod_{i=t+1}^j \gamma(S_i) \right) R_{j+1} \right]$$

Expected Discounted Return

Target Policy: $A_{t:\infty} \sim \pi$

Behavior Policy: $A_{t:\infty} \sim \mu$

Discount: $\gamma(s_t) \in [0,1]$

Reward/Cumulant: $r(s_t) \in \mathbb{R}$

Off-policy Learning

Learn about a **target policy** π using data generated from a **behavior policy** b .

Off-policy Learning

Want $\mathbb{E} [\Delta_w(A) \mid A \sim \pi]$

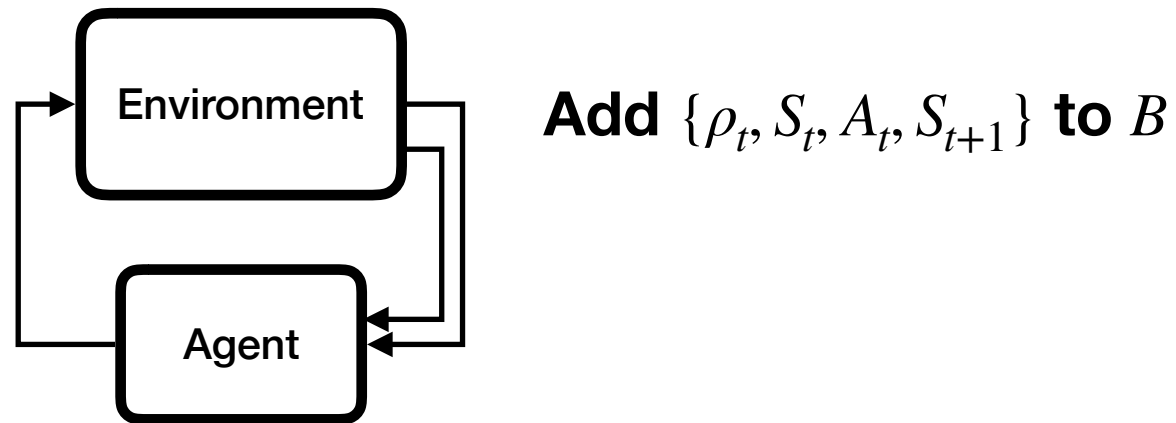
$= \mathbb{E} [\boxed{} \Delta_w(A) \mid A \sim b]$ **Have**

Off-policy Learning

$$\begin{aligned}\mathbb{E} [\Delta_w(A) \mid A \sim \pi] &= \sum_{a \in \mathcal{A}} \pi(a) \Delta_w(a) \\ &= \sum_{a \in \mathcal{A}} \pi(a) \frac{b(a)}{b(a)} \Delta_w(a) \\ &= \sum_{a \in \mathcal{A}} \frac{\pi(a)}{b(a)} b(a) \Delta_w(a) \\ &= \mathbb{E} [\rho(A) \Delta_w(A) \mid A \sim b]\end{aligned}$$

Reweighting

Interact with Environment:



Sample Minibatch:

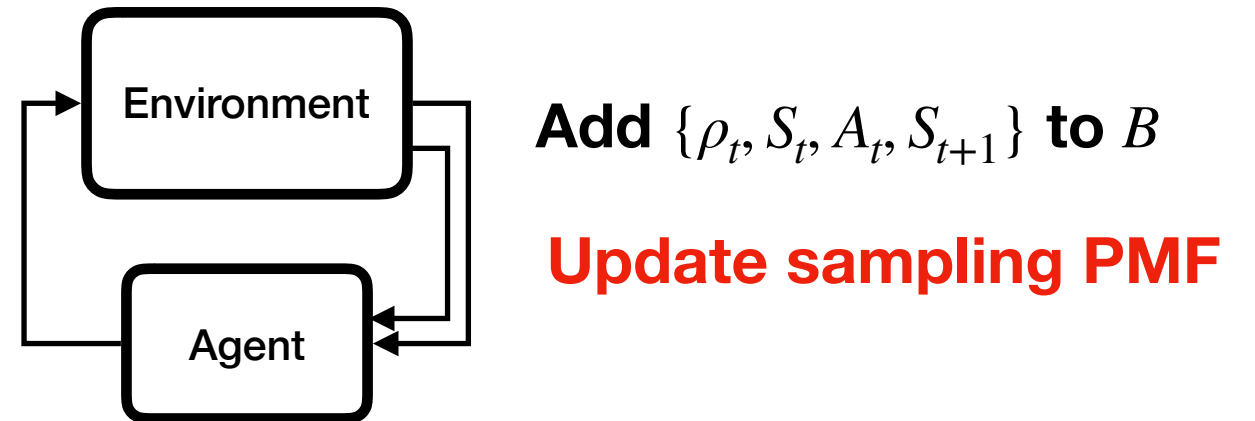
Sample transition $\{\rho_i, S_i, A_i, S'_i\}$ with $Pr \left\{ \frac{1}{|B|} \right\}$ (n times)

Calculate Updates:

$$\Delta_{IS} = \frac{1}{n} \sum_{i=1}^n \rho_i \delta_i \nabla_w V(s_i; w)$$

Resampling

Interact with Environment:



Sample Minibatch:

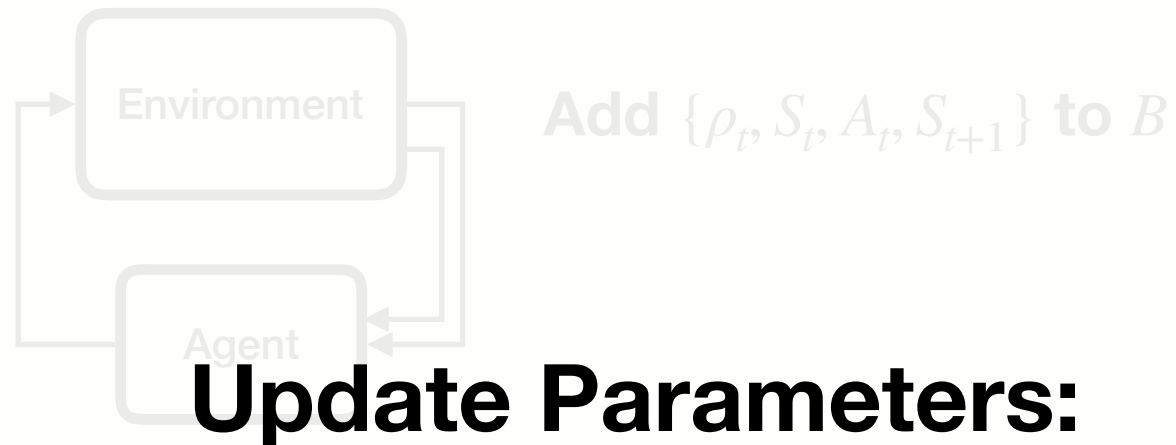
Sample transition $\{\rho_i, S_i, A_i, S'_i\}$ with $Pr \left\{ \frac{\rho_i}{\sum_j \rho_j} \right\}$ (n times)

Calculate Updates:

$$\Delta_{IR} = \frac{1}{n} \sum_{i=1}^n \delta_i \nabla_w V(s_i; w)$$

Reweighting

Interact with Environment:



Sample Minibatch:

$$\{ \rho_i, S_i, A_i, S'_i \} \text{ with } Pr \left\{ \frac{1}{|B|} \right\}$$

(n times)

$$w' = w - \alpha_t \Delta_*$$

Calculate Updates:

$$\Delta_{IS} = \frac{1}{n} \sum_{i=1}^n \rho_i \delta_i \nabla_w V(s_i; w)$$

Resampling

Interact with Environment:



Sample Minibatch:

$$\{ \rho_i, S_i, A_i, S'_i \} \text{ with } Pr \left\{ \frac{\rho_i}{\sum_j \rho_j} \right\}$$

(n times)

Calculate Updates:

$$\Delta_{IR} = \frac{1}{n} \sum_{i=1}^n \delta_i \nabla_w V(s_i; w)$$

Off-policy Learning

With a buffer of experience

Importance Sampling (IS):

$$\Delta_{IS} = \frac{1}{n} \sum_{i=1}^n \rho_i \delta_i \nabla_w V(s_i; w)$$

Importance Resampling (IR):

$$\Delta_{IR} = \frac{1}{n} \sum_{i=1}^n \delta_i \nabla_w V(s_i; w)$$

WIS-Minibatch:

$$\Delta_{WIS} = \frac{1}{\sum_{j=1}^n \rho_j} \sum_{i=1}^n \rho_i \delta_i \nabla_w V(s_i; w)$$

VTrace(0):

$$\bar{\rho}_i = \begin{cases} \rho_{clip} & \rho_i > \rho_{clip} \\ \rho_t & \textbf{o.w.} \end{cases}$$

$$\Delta_{VTrace} = \frac{1}{n} \sum_{i=1}^n \bar{\rho}_i \delta_i \nabla_w V(s_i; w)$$

Off-policy Learning

With a buffer of experience

Reweighting

Importance Sampling (IS)

VTrace(0)

WIS-Minibatch

Resampling

Importance Resampling (IR)

Off-policy Learning

WIS-TD(0)

Hypothesized Empirical Benefits

- IR reduces the **update variance** as compared with IS.
- IR can **update less to learn more** (sample efficiency).

Variance in Off-policy Prediction

Update Variance:

$$\text{Var} \{ \Delta_{IS} \} = \text{Var} \left\{ \left\| \frac{1}{n} \sum_{i=1}^n \rho_i \delta_i \nabla_w V(s_i; w) \right\|_1 \right\}$$

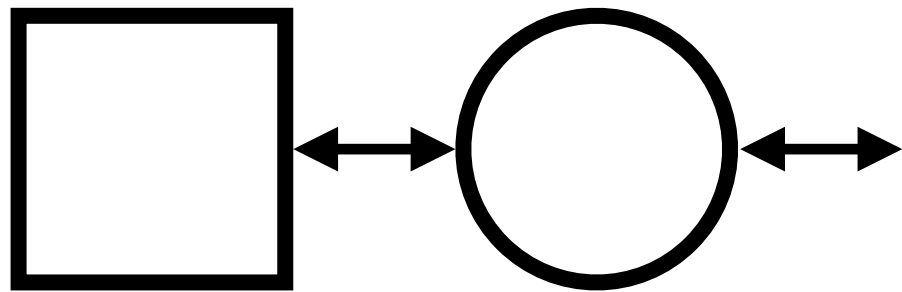
Benefits of reduced update variance:

- Reduced sensitivity to learning rate.
- Faster learning

Empirical Results

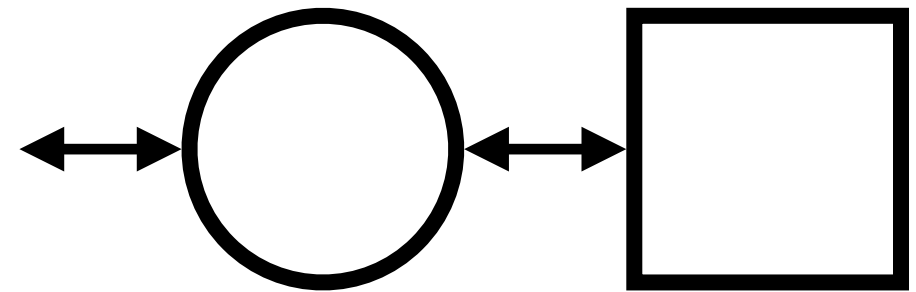
Markov Chain

$C = 0$



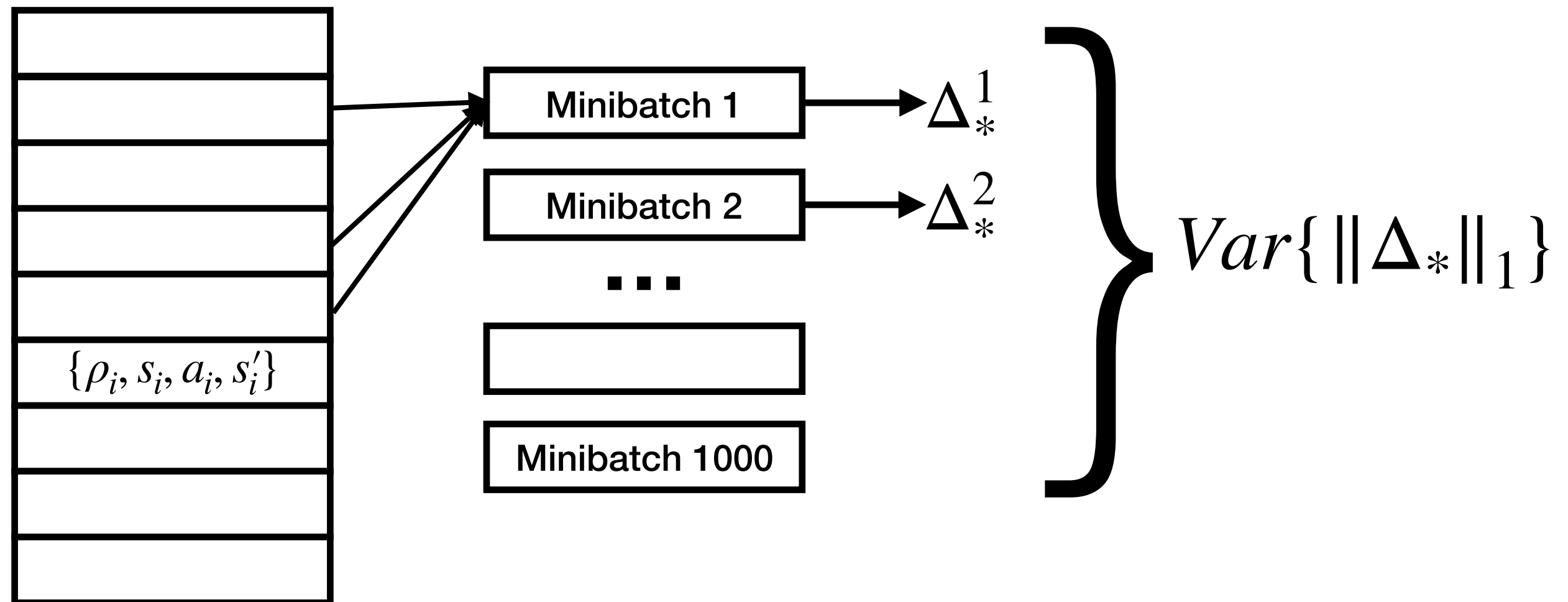
...

$C = 1$



Markov Chain

Estimating the Update Variance:



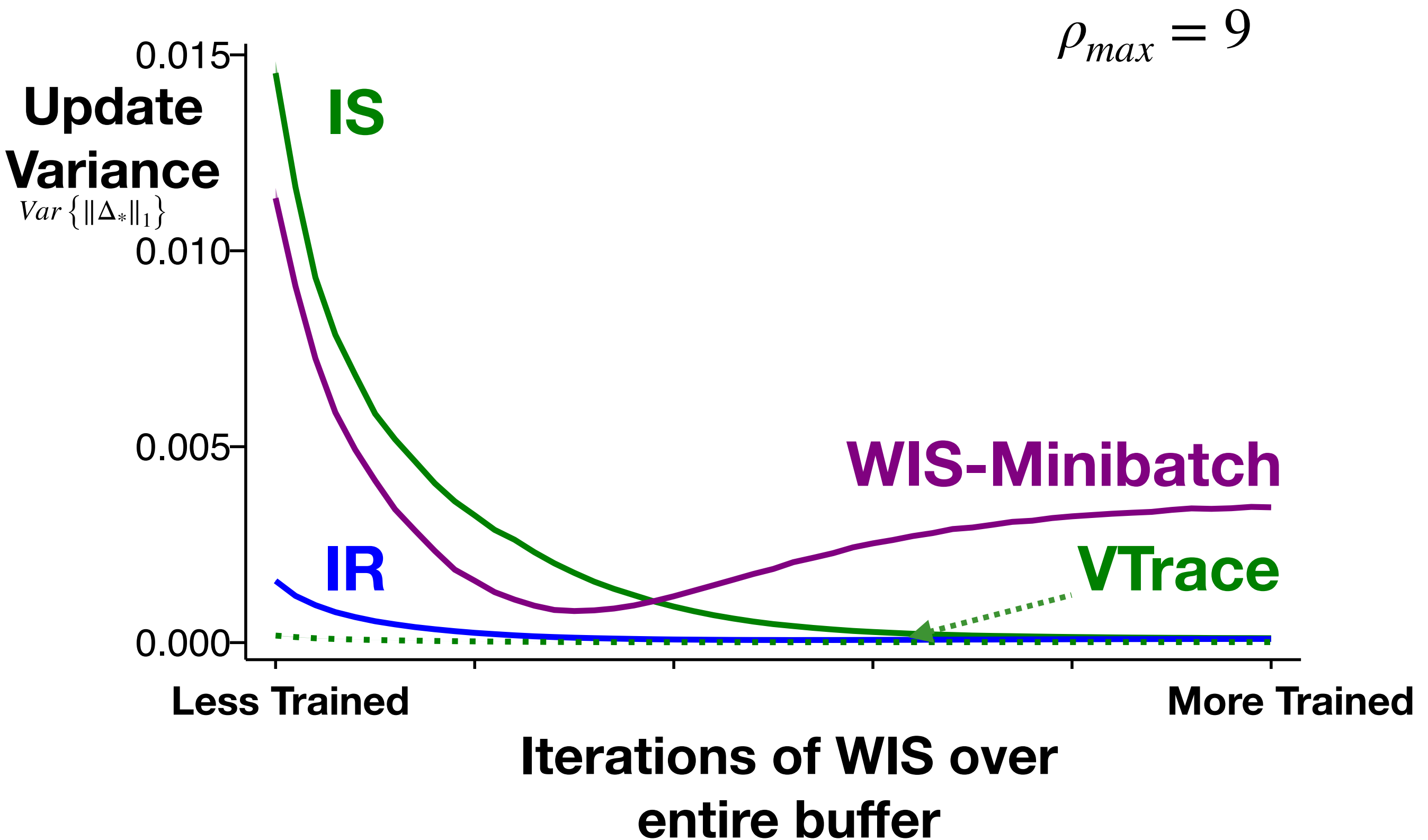
Behavior:

$$b(a | s) = \begin{cases} 0.9 & \text{if } a = \textit{left} \\ 0.1 & \text{if } a = \textit{right} \end{cases}$$

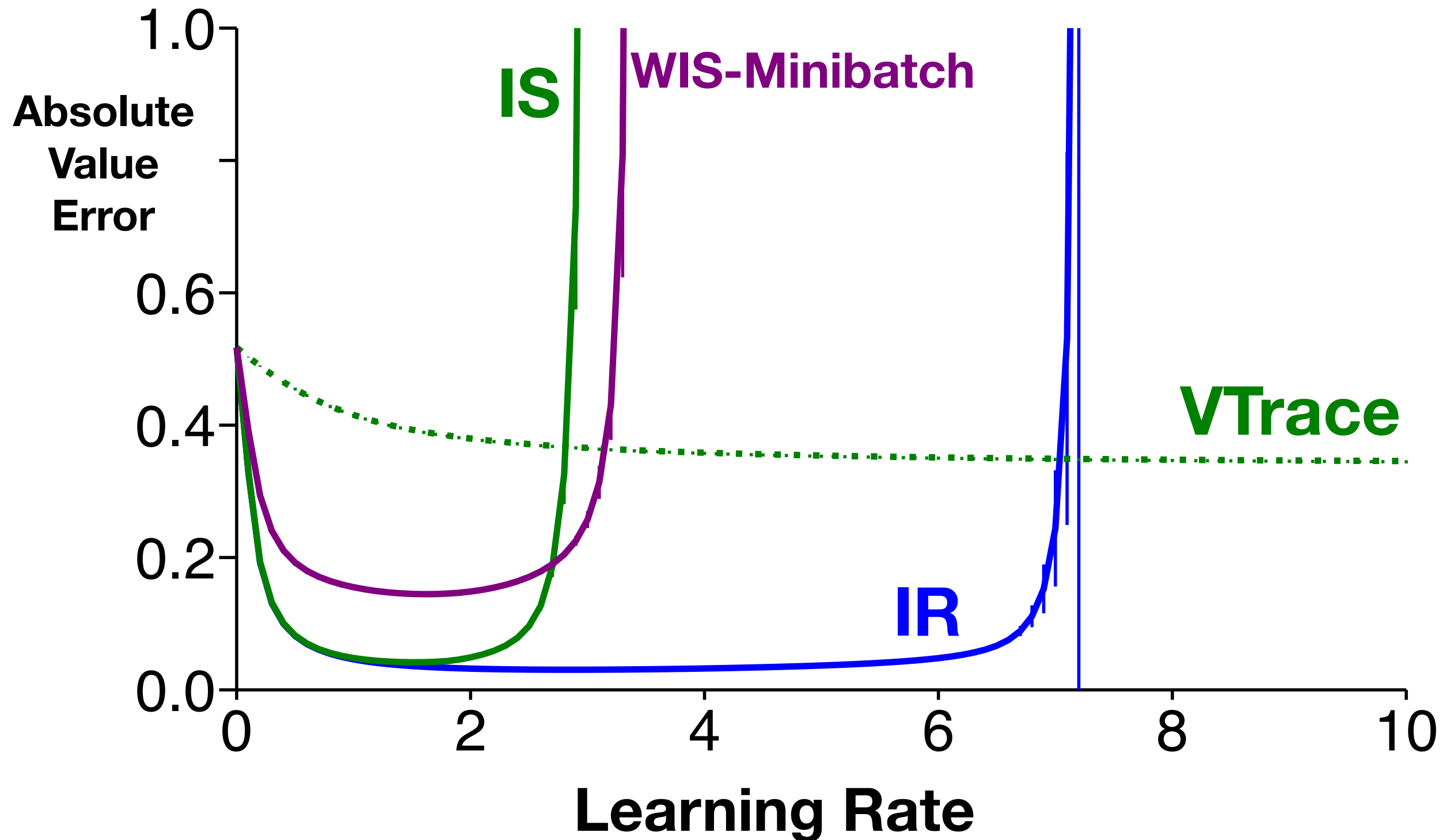
Target:

$$\pi(a | s) = \begin{cases} 0.1 & \text{if } a = \textit{left} \\ 0.9 & \text{if } a = \textit{right} \end{cases}$$

Markov Chain - Update Variance

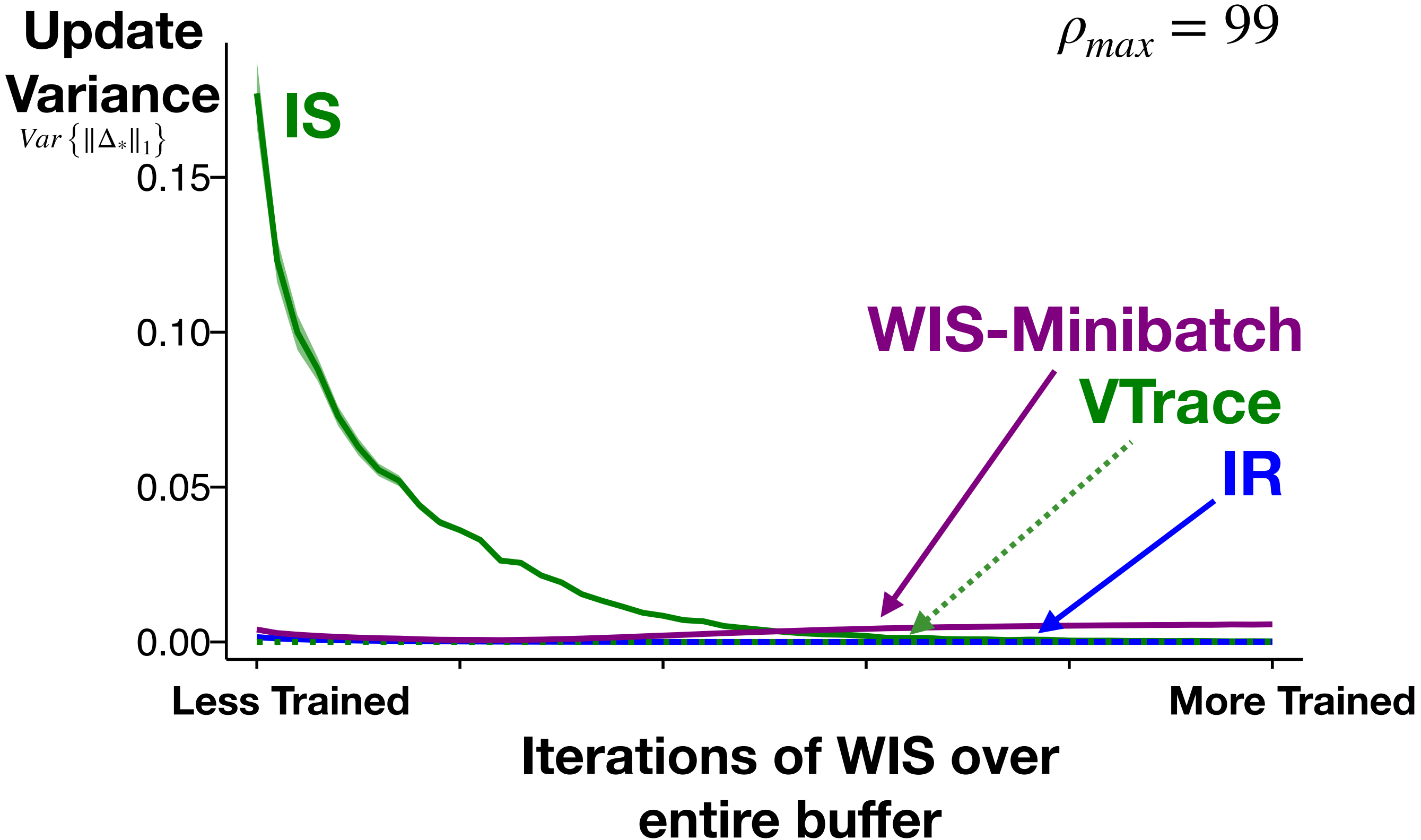


Markov Chain - Learning Rate Sensitivity

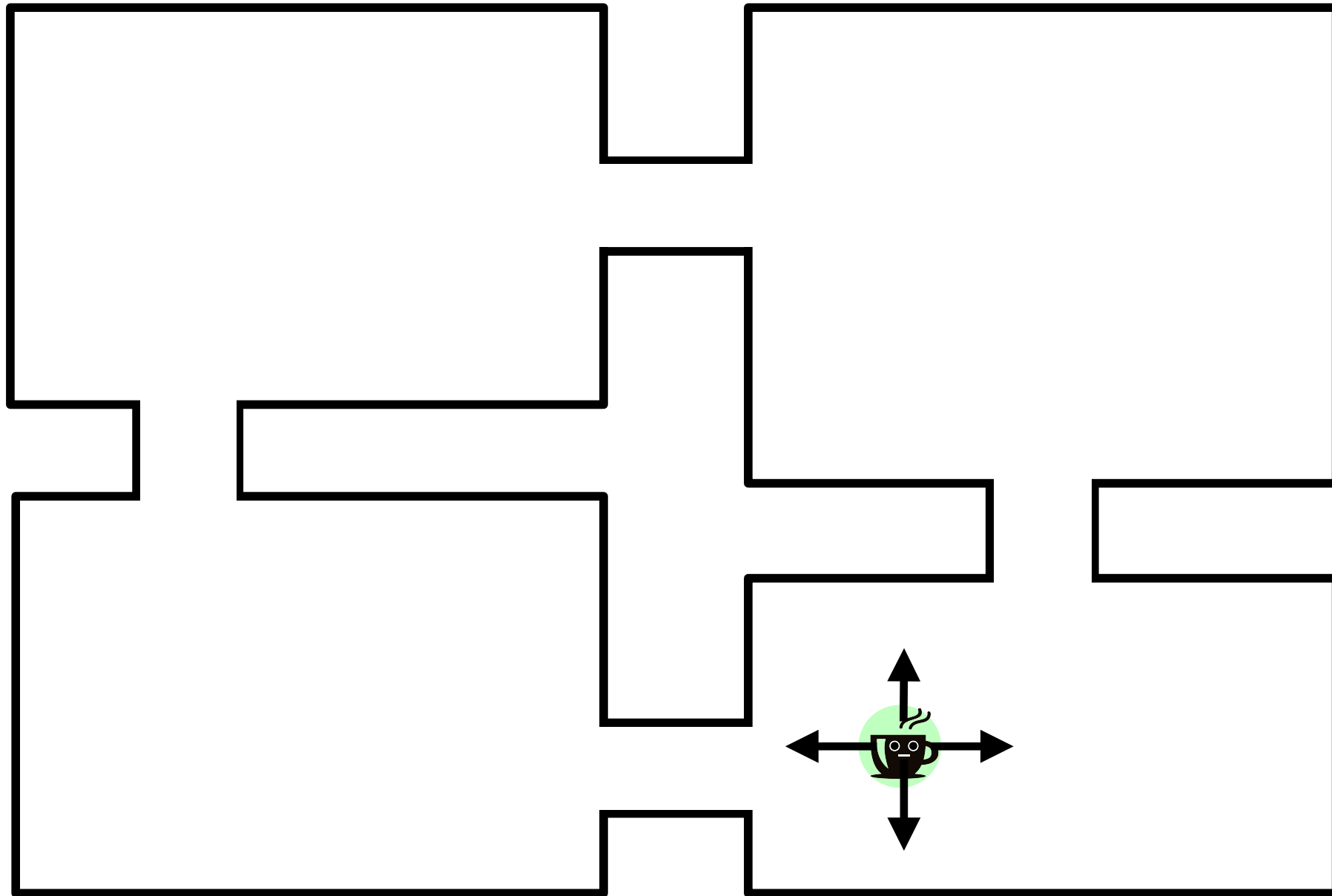


Markov Chain - Update Variance

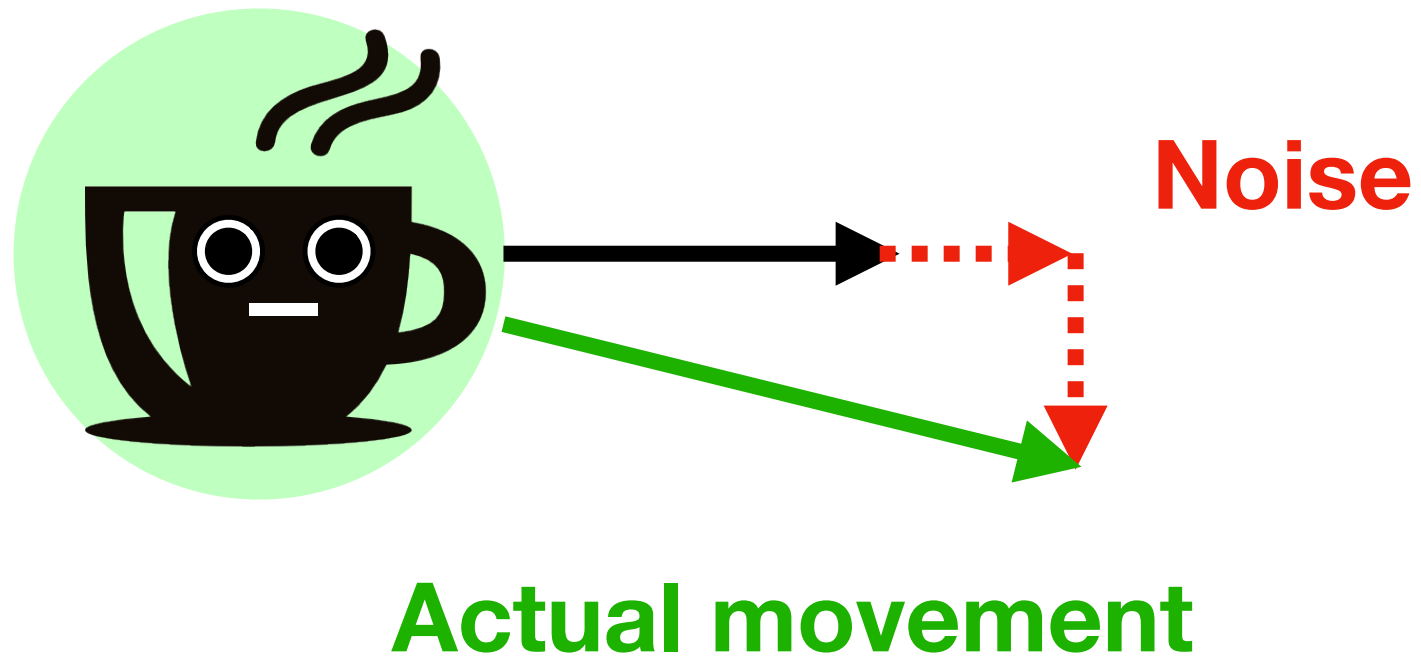
High Variance



Continuous Four Rooms



Continuous Four Rooms



Continuous Four Rooms

Evaluation:

- Sampled 1000 states from the stationary distribution of the behavior policy
- Estimated returns with 100 Monte Carlo rollouts

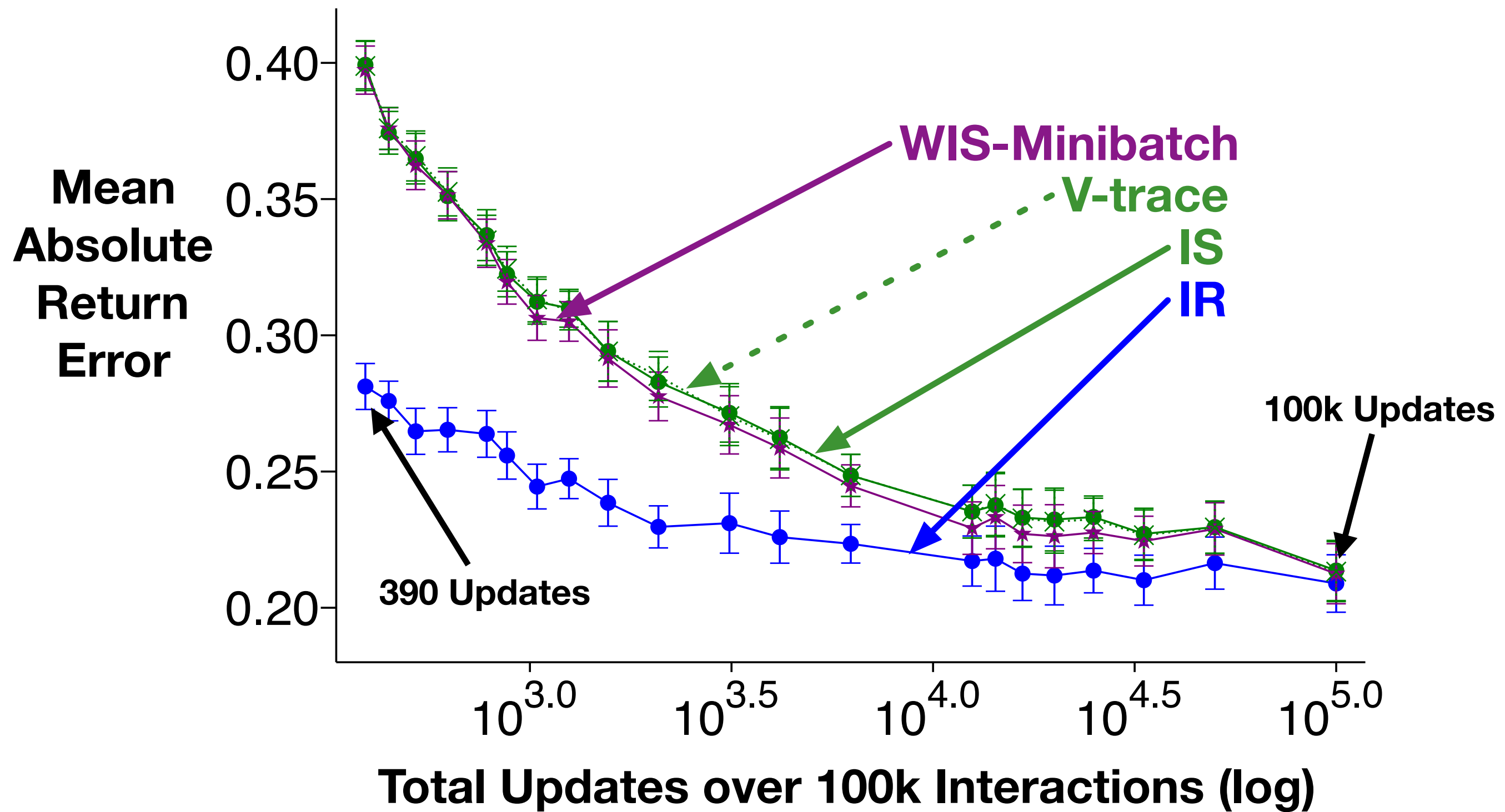
Behavior:

$$b(\cdot | s) = 0.25$$

Target:

$$\pi_1(a | s) = \begin{cases} 1 & \text{if } a = \text{down} \\ 0 & \text{o.w.} \end{cases}$$

Cont. Four Rooms - Total Updates



Conclusions

1. Resampling can have **lower variant updates** as compared to importance sampling.
2. Resampling generally needs **fewer updates** to reach comparable performance to importance sampling.
3. Resampling and importance sampling **perform comparably when many samples are used.**

Questions?



More Experiments!

Weird behavior of induced bias!!

Theory!

<https://arxiv.org/pdf/1906.04328.pdf>

Theory

Theoretical Properties of IR

- Biased and Consistent (with a small correction term)
- Consistent (with the correction term) under a changing buffer of experience.
- Under some conditions, guaranteed equal or lower variance compared to IS.

Bias and Consistency

Bias

Theorem 3.1. *[Bias for a fixed buffer of size n] Assume a buffer B of n transitions is sampled i.i.d., according to $d_\mu(s)\mu(a|s)P(s'|s, a)$. Let $X_{\text{WIS}^*} \stackrel{\text{def}}{=} \sum_{i=1}^n \frac{\rho_i}{\sum_{j=1}^n \rho_j} \Delta_i$ be the WIS-Optimal estimator of the update. Then,*

$$\mathbb{E}[X_{\text{IR}}] = \mathbb{E}[X_{\text{WIS}^*}]$$

and so the bias of X_{IR} is proportional to

$$\text{Bias}(X_{\text{IR}}) = \mathbb{E}[X_{\text{IR}}] - \mathbb{E}_\pi[\Delta] \propto \frac{1}{n} (\mathbb{E}_\pi[\Delta] \sigma_\rho^2 - \sigma_{\rho, \Delta} \sigma_\rho \sigma_\Delta) \quad (1)$$

Consistency

Theorem 3.2. *Let $B_i = \{X_{i-n+1}, \dots, X_i\}$ be the buffer of the most recent n transitions sampled by time i , i.i.d. as specified in Assumption 1. Let $X_{\text{BC}}^{(i)}$ be the bias-corrected IR estimator, with k samples from buffer B_i . Define the sliding-window estimator $X_t \stackrel{\text{def}}{=} \frac{1}{t} \sum_{i=1}^t X_{\text{BC}}^{(i)}$. Assume there exists a $c > 0$ such that $\text{Var}(X_{\text{BC}}^{(i)}) \leq c \forall i$. Then, as $t \rightarrow \infty$, X_t converges in probability to $\mathbb{E}_\pi[\Delta]$.*

Variance

Theorem 3.3. Assume that, for a given buffer B , $\|\Delta_j\|_2^2 > \frac{c}{\rho_j}$ for samples where $\rho_j \geq \bar{\rho}$, and that $\|\Delta_j\|_2^2 < \frac{c}{\rho_j}$ for samples where $\rho_j < \bar{\rho}$, for some $c > 0$. Then the BC-IR estimator has lower variance than the IS estimator: $\mathbb{V}(X_{\text{BC}} \mid B) < \mathbb{V}(X_{\text{IS}} \mid B)$.

Theorem 3.4. Assume ρ and the magnitude of the update $\|\Delta\|_2^2$ are independent

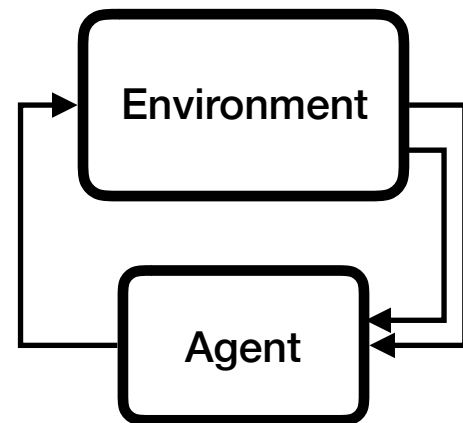
$$\mathbb{E}[\rho_j \|\Delta_j\|_2^2 \mid B] = \mathbb{E}[\rho_j \mid B] \mathbb{E}[\|\Delta_j\|_2^2 \mid B]$$

Then the BC-IR estimator will have equal or lower variance than the IS estimator.

Future Directions

Future Directions - Sampling Policy

Interact with Environment:



Add $\{\bar{\rho}_t, s_t, a_t, s_{t+1}\}$ **to** B

$$\bar{\rho}_t = \frac{d(s_t, a_t)}{b(s_t | a_t)}$$

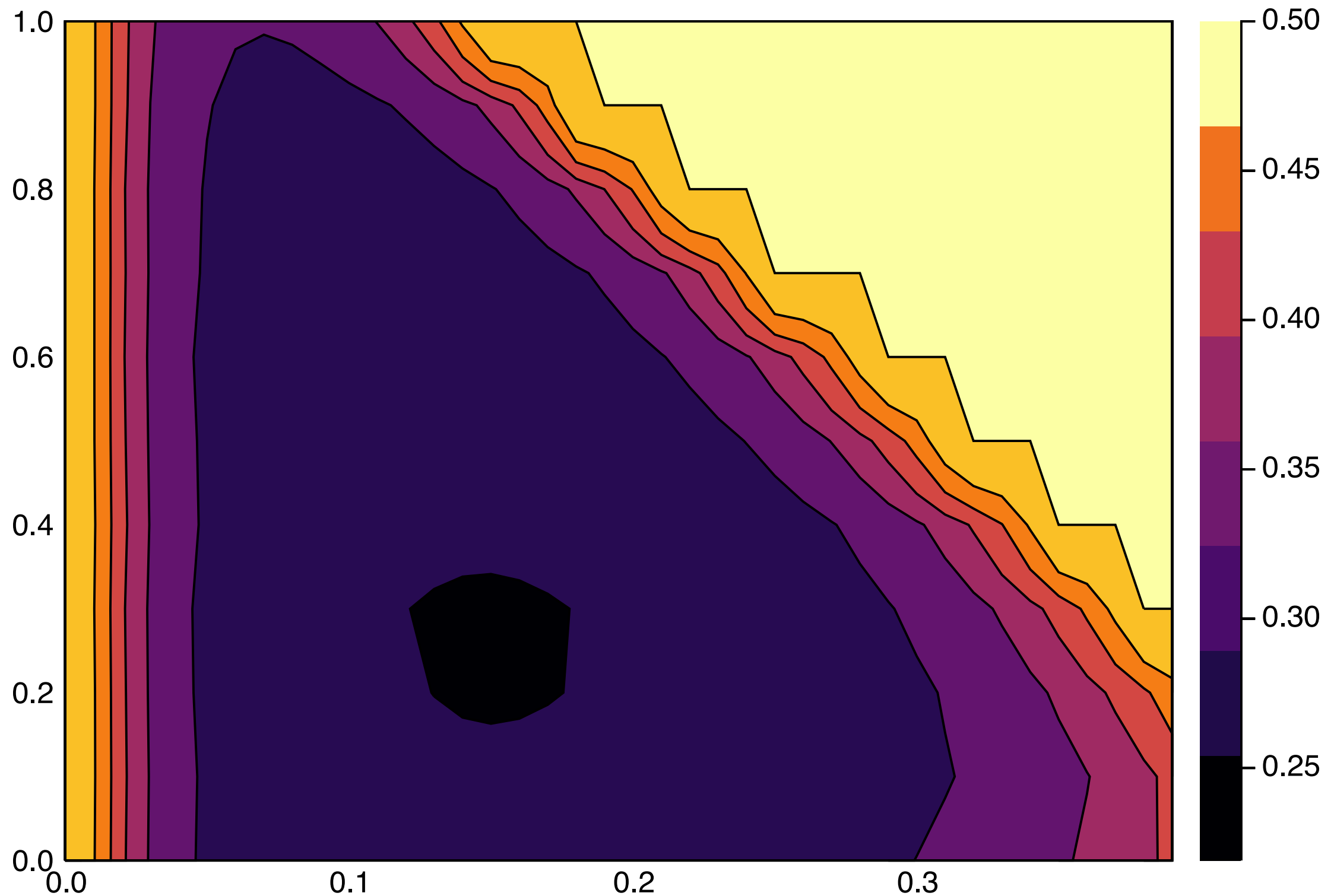
Sample Minibatch:

$$\{\pi_{\text{sample}}(a_i | s_i), s_i, a_i, s'_i\} \text{ with } \mathbb{P} \left\{ \frac{\bar{\rho}_i}{\sum_{j=1}^{|B|} \bar{\rho}_j} \mid B \right\}$$

Update Parameters:

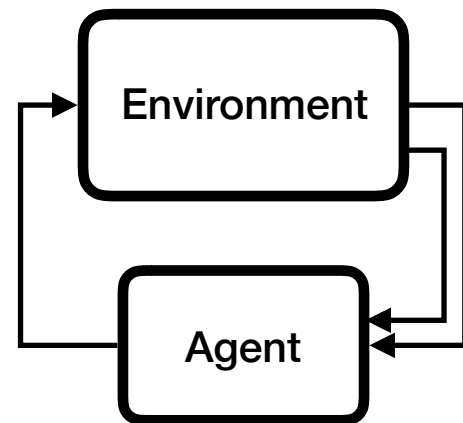
$$\Delta\theta = \frac{1}{n} \sum_{i=1}^n \frac{\pi(a_i | s_i)}{d(s_i, a_i)} \delta_i \nabla_{\theta} V(s_i; \theta)$$

Future Directions - Sampling Policy



Future Directions - Multi-step learning

Interact with Environment:



Add $\{\bar{\rho}_t, s_t, a_t, s_{t+1}\}$ **to** B

$$\bar{\rho}_t = \rho_{t:t+k} = \prod_{i=t}^{t+k} \rho(a_i | s_i)$$

Sample Minibatch:

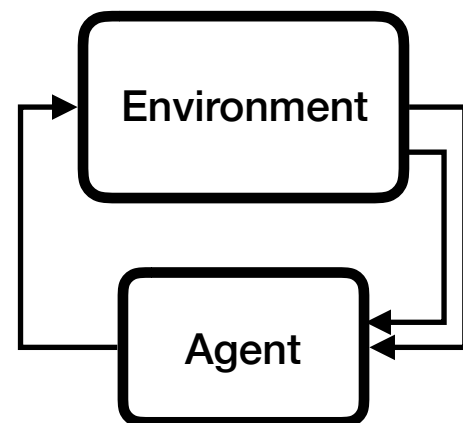
$$\{\bar{\rho}_i, s_i, a_i, s'_i\} \text{ with } \mathbb{P} \left\{ \frac{\bar{\rho}_i}{\sum_{j=1}^{|B|} \bar{\rho}_j} \mid B \right\}$$

Update Parameters:

$$\Delta\theta = \frac{1}{n} \sum_{i=1}^n \delta_i^k \nabla_{\theta} V(s_i; \theta)$$

Future Directions - state distributions

Interact with Environment:



Add $\{\bar{\rho}_t, s_t, a_t, s_{t+1}\}$ **to** B

$$\bar{\rho}_t = \frac{\pi(s_t | a_t) \hat{\mu}_\pi(s_t)}{b(s_t | a_t) \hat{\mu}_b(s_t)}$$

Sample Minibatch:

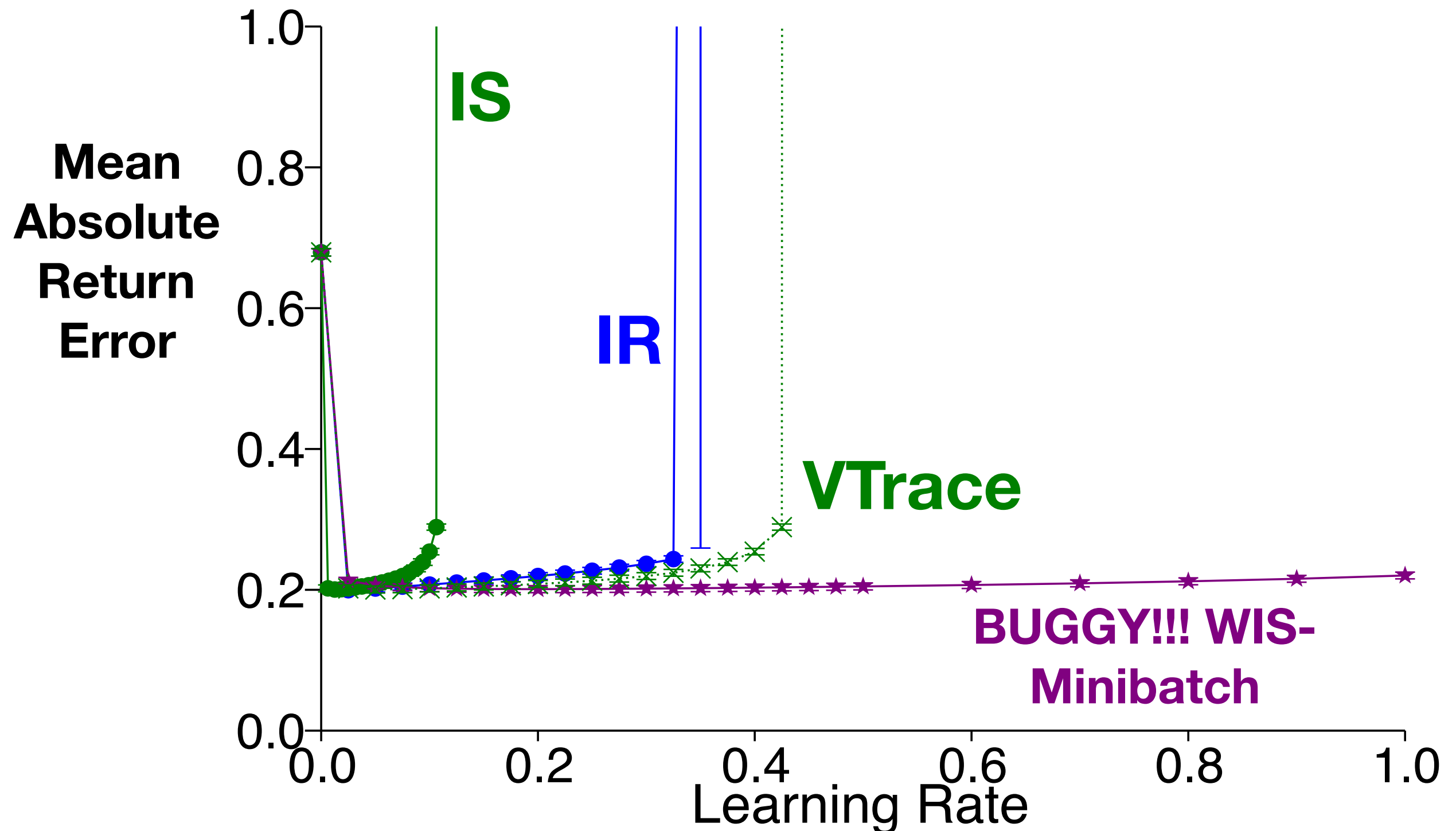
$$\{\bar{\rho}_i, s_i, a_i, s'_i\} \text{ with } \mathbb{P} \left\{ \frac{\bar{\rho}_i}{\sum_{j=1}^{|B|} \bar{\rho}_j} \mid B \right\}$$

Update Parameters:

$$\Delta\theta = \frac{1}{n} \sum_{i=1}^n \delta_i \nabla_{\theta} V(s_i; \theta)$$

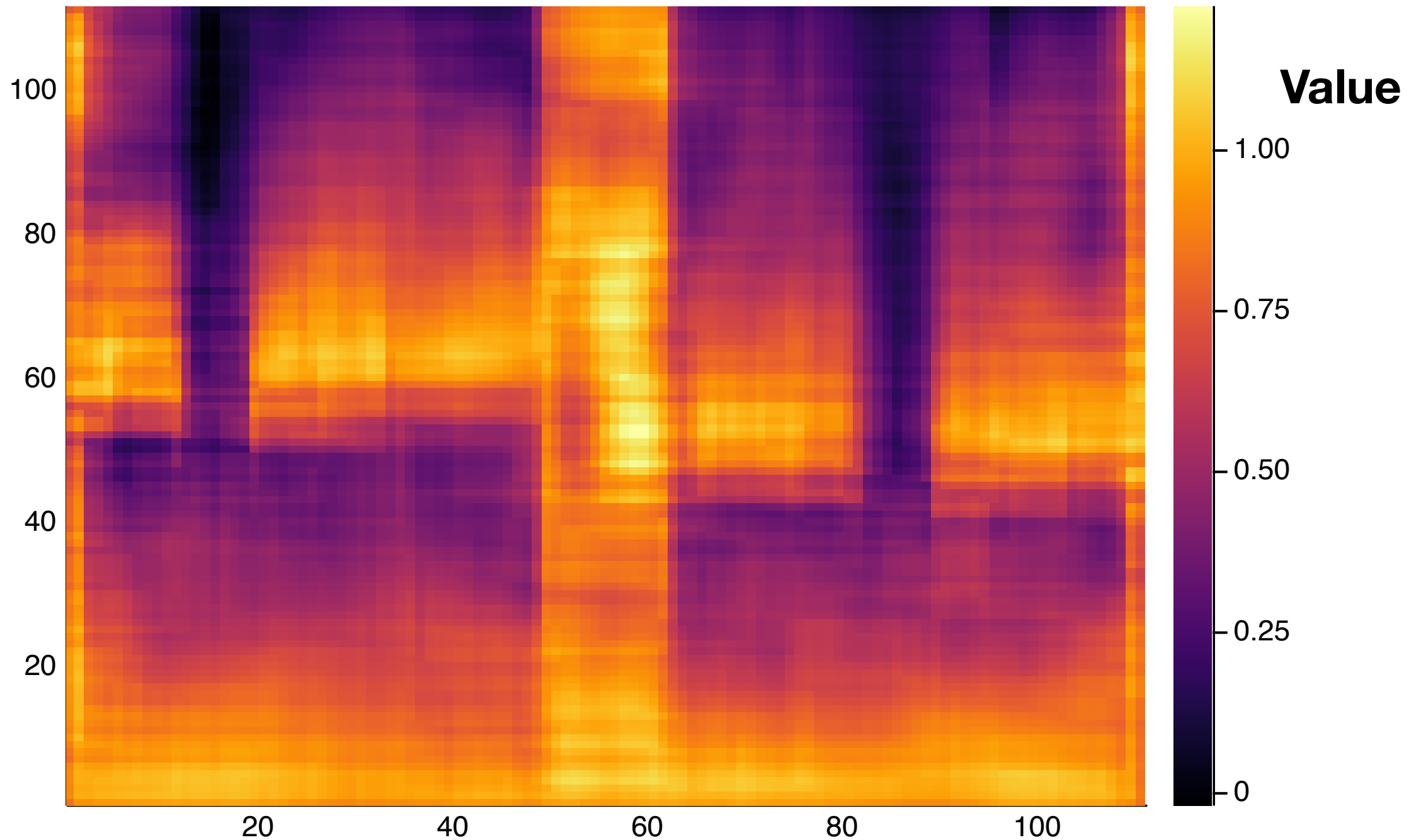
Extra Results

Four Rooms Cont LR sens



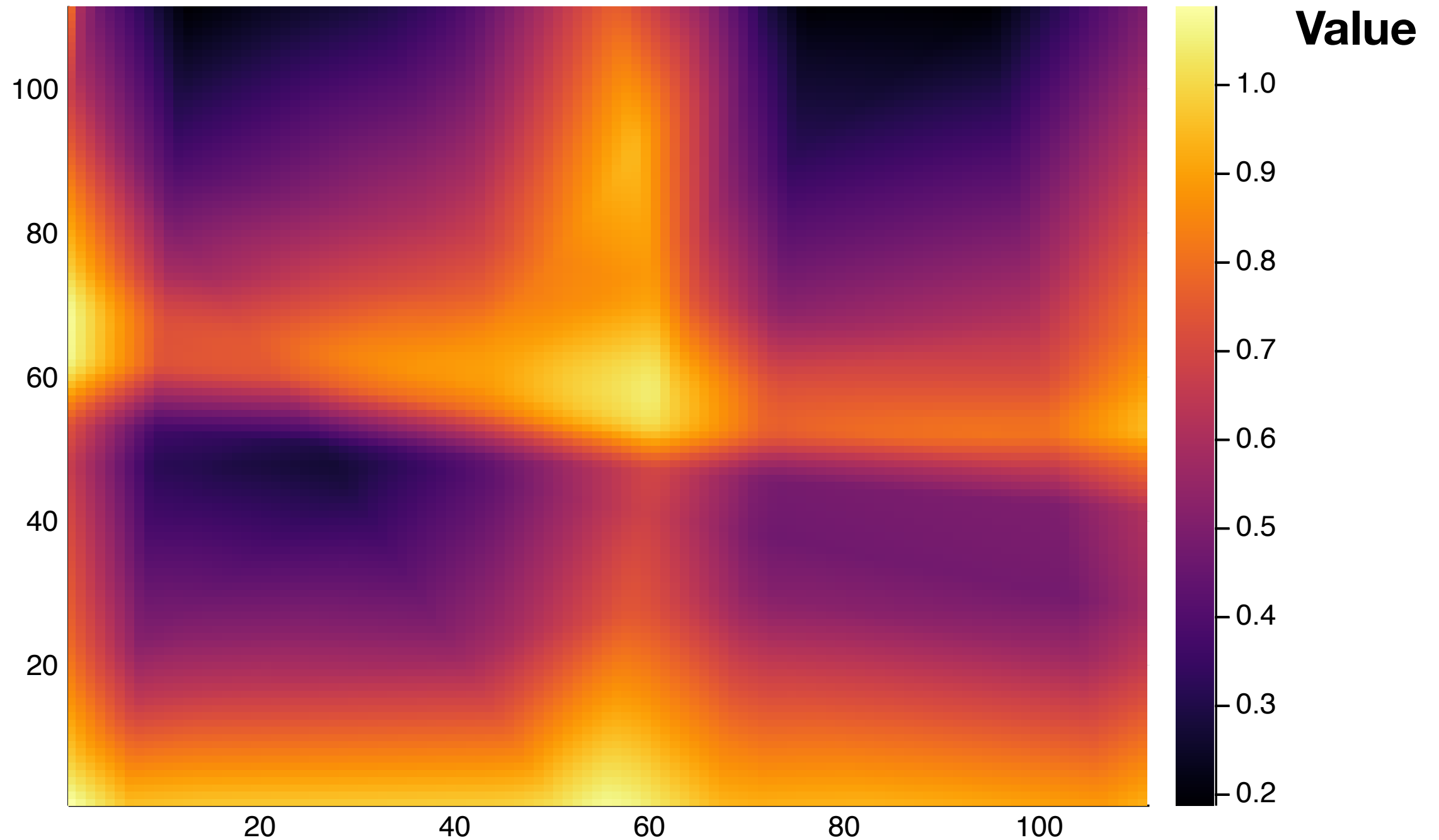
Learned Value Function CFR

IS (Tile Coded)

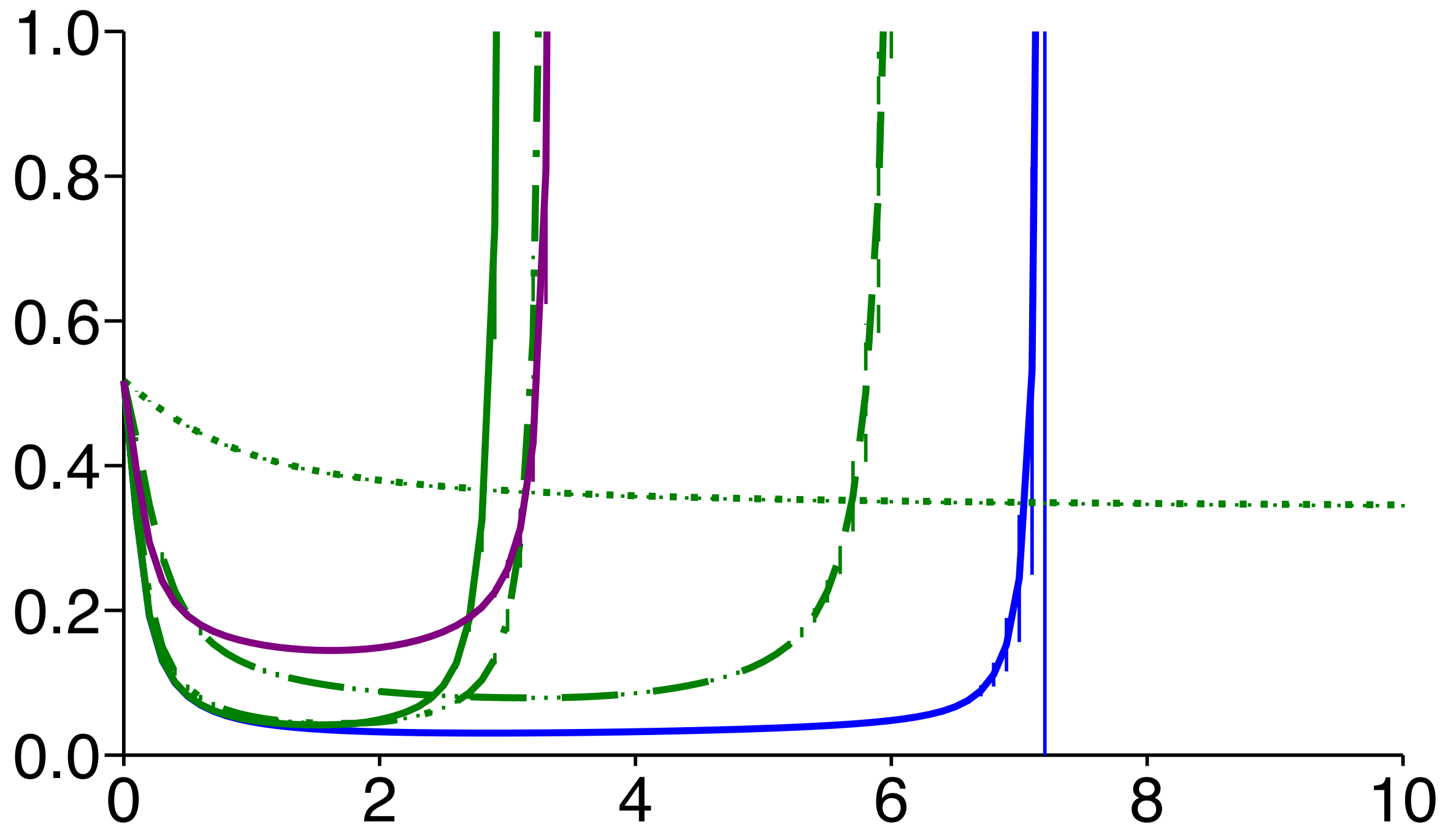


Learned Value Function CFR

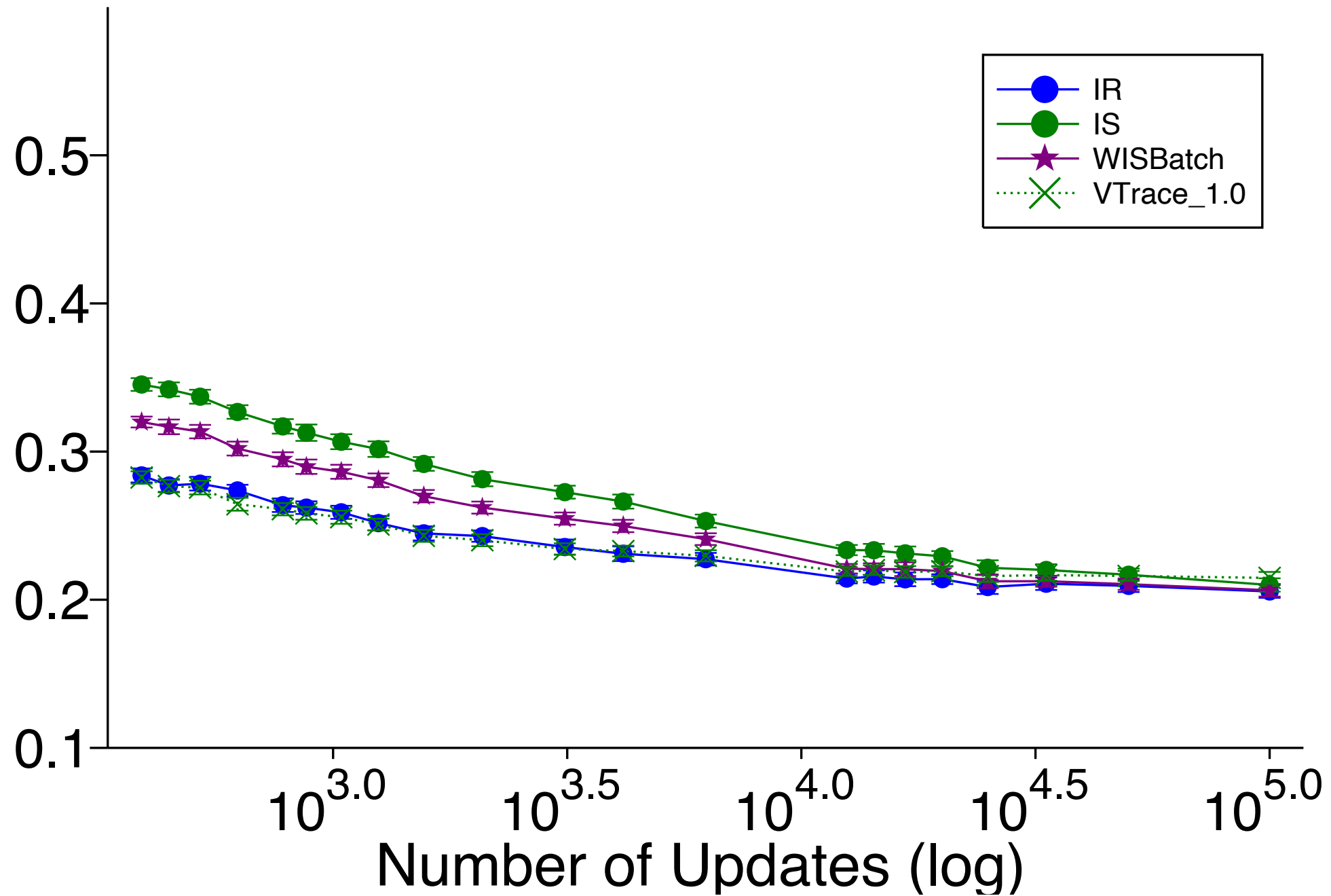
IS (Artificial Neural Network)



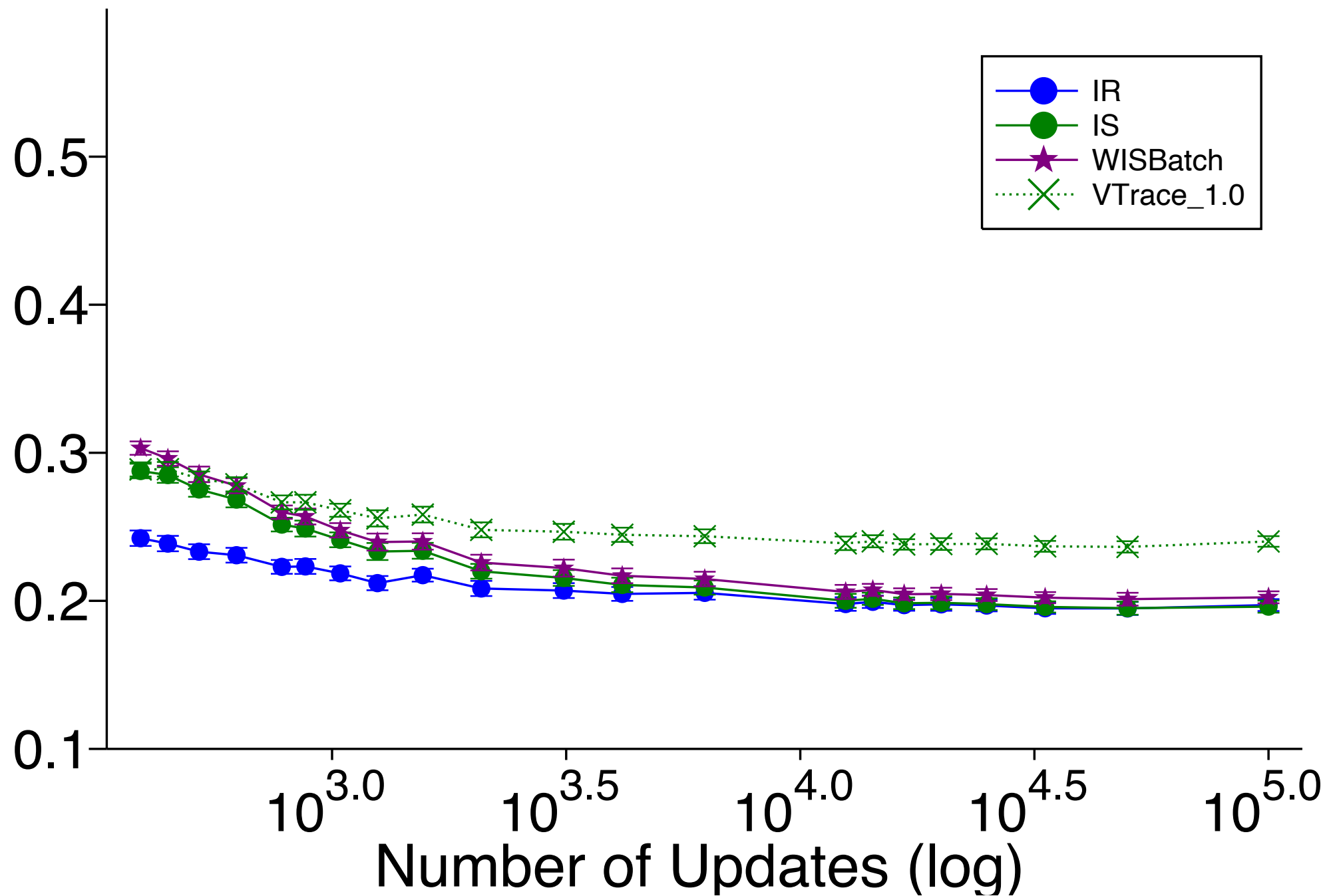
Future Directions - Odd Behavior of Bias



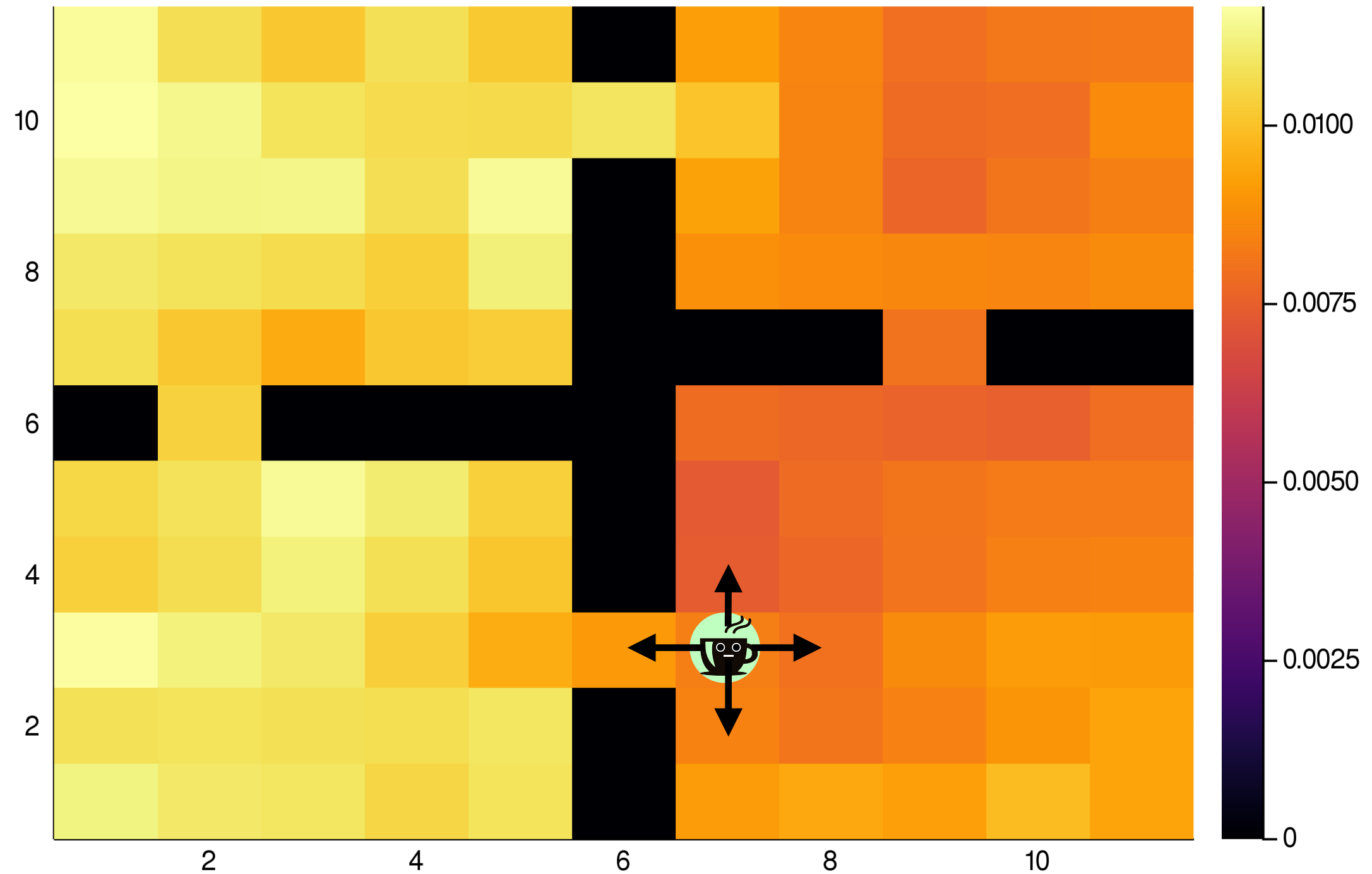
Future Directions - Odd Behavior of Bias



Odd Behavior of Bias (VTrace)



Four Rooms



Four Rooms

Evaluation:

- Estimate value function using dynamic programming

Behavior:

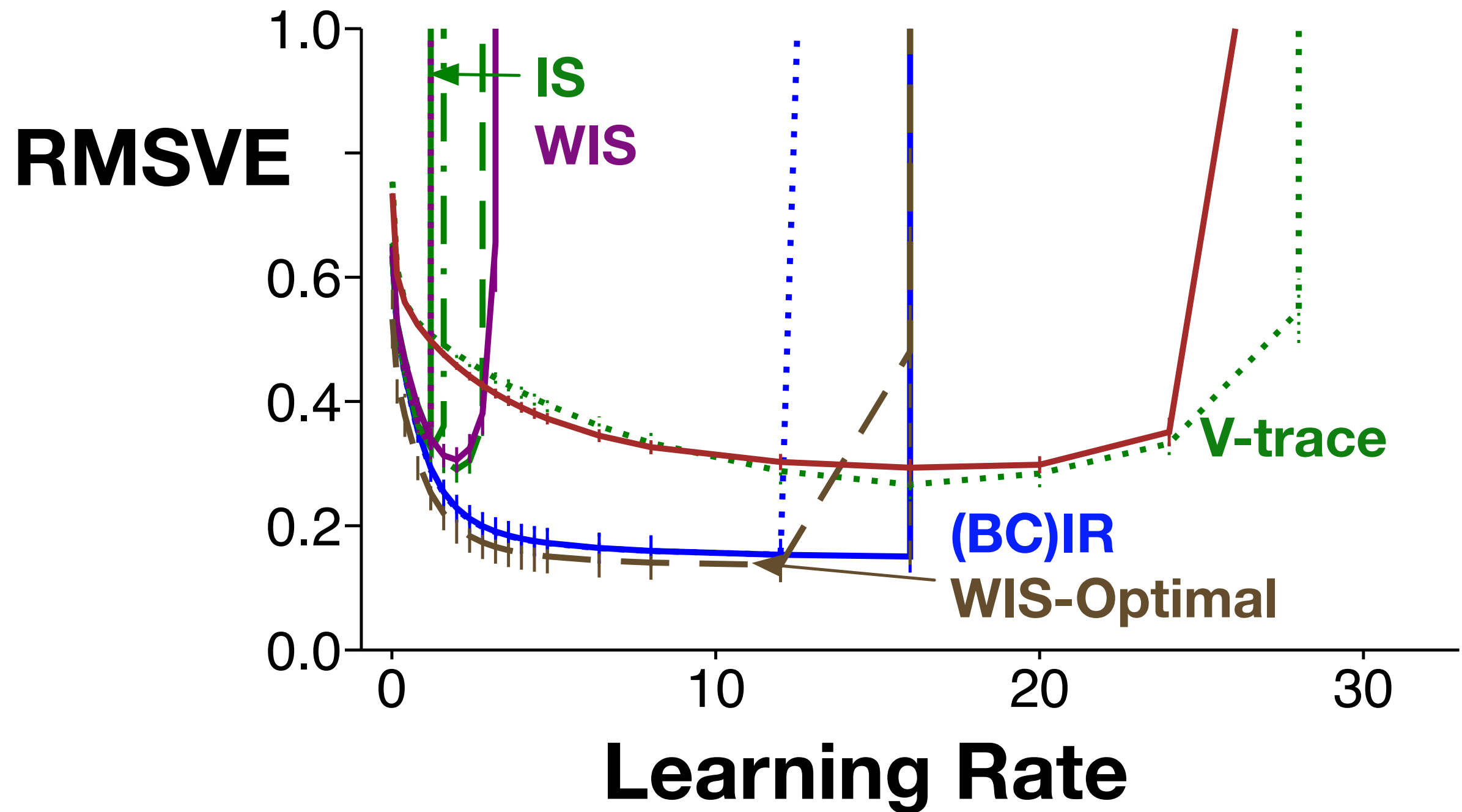
$$b(\cdot \mid s \notin S_{hv}) = 0.25$$

$$b(a \mid s \in S_{hv}) = \begin{cases} 0.1 & \text{if } a = \textit{down} \\ \frac{0.9}{3} & \text{o.w.} \end{cases}$$

Target:

$$\pi_1(a \mid s) = \begin{cases} 1 & \text{if } a = \textbf{down} \\ 0 & \text{o.w.} \end{cases}$$

Four Rooms - 31,250 updates over 500k interactions



Four Rooms - 500k updates

