

## INTRODUCTION

### Motivation

- Batch RL is important in risky/expensive real-world applications (medicine, robotics, etc.)
- Safe policy improvement (SPI) is critical to ensure new policies perform at least as well as the existing ones.
- Existing SPI methods: density-based policy regularization, pessimism-based approaches are too conservative, SPIBB[1] needs  $\pi_b$

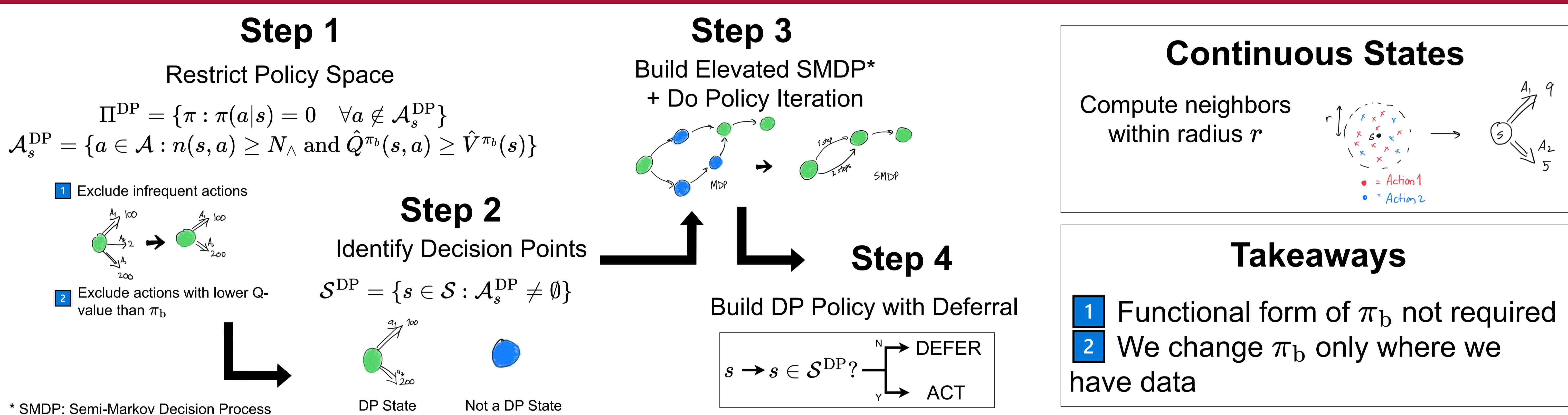
### Contributions

- Introduce Decision-Point RL (DPRL), limiting policy deviations to "decision points"
- DPRL does not require access to the behavior policy during training
- Tighter theoretical guarantees for SPI (dependent on visitation frequency)

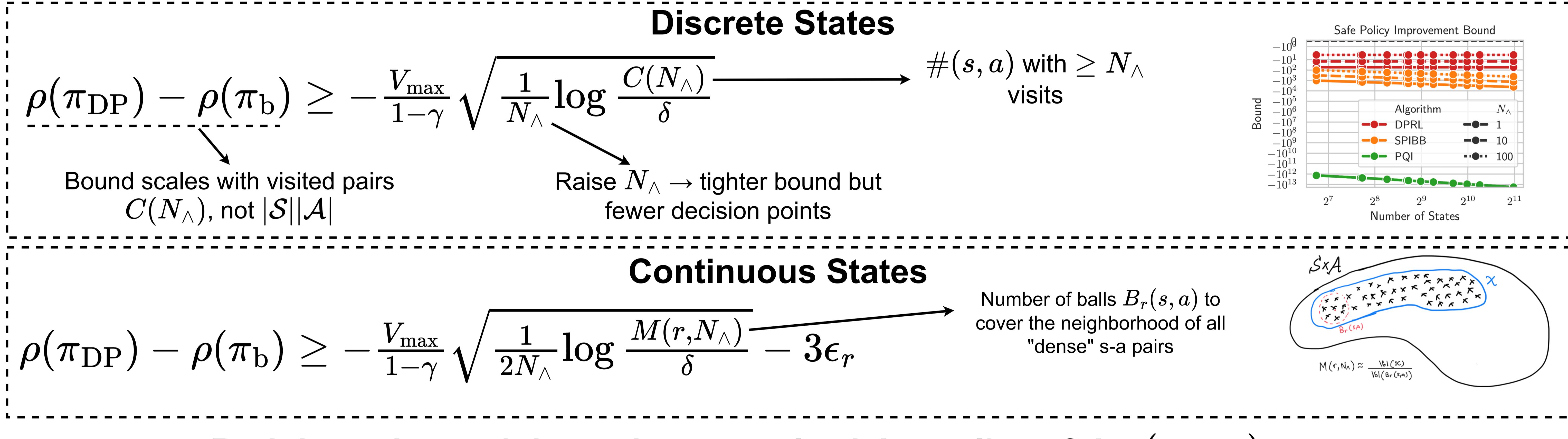
### Setup

- Assume MDP  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ , with discrete  $\mathcal{A}$  and bounded  $R \in [0, R_{\max}]$
- Given: Dataset  $\mathcal{D} = \{S_0^n, A_0^n, R_0^n, \dots, S_{T_n}^n, A_{T_n}^n, R_{T_n}^n\}_{n=1}^N$  with actions taking using Behavior Policy  $\pi_b$

## DECISION-POINT REINFORCEMENT LEARNING (DPRL)

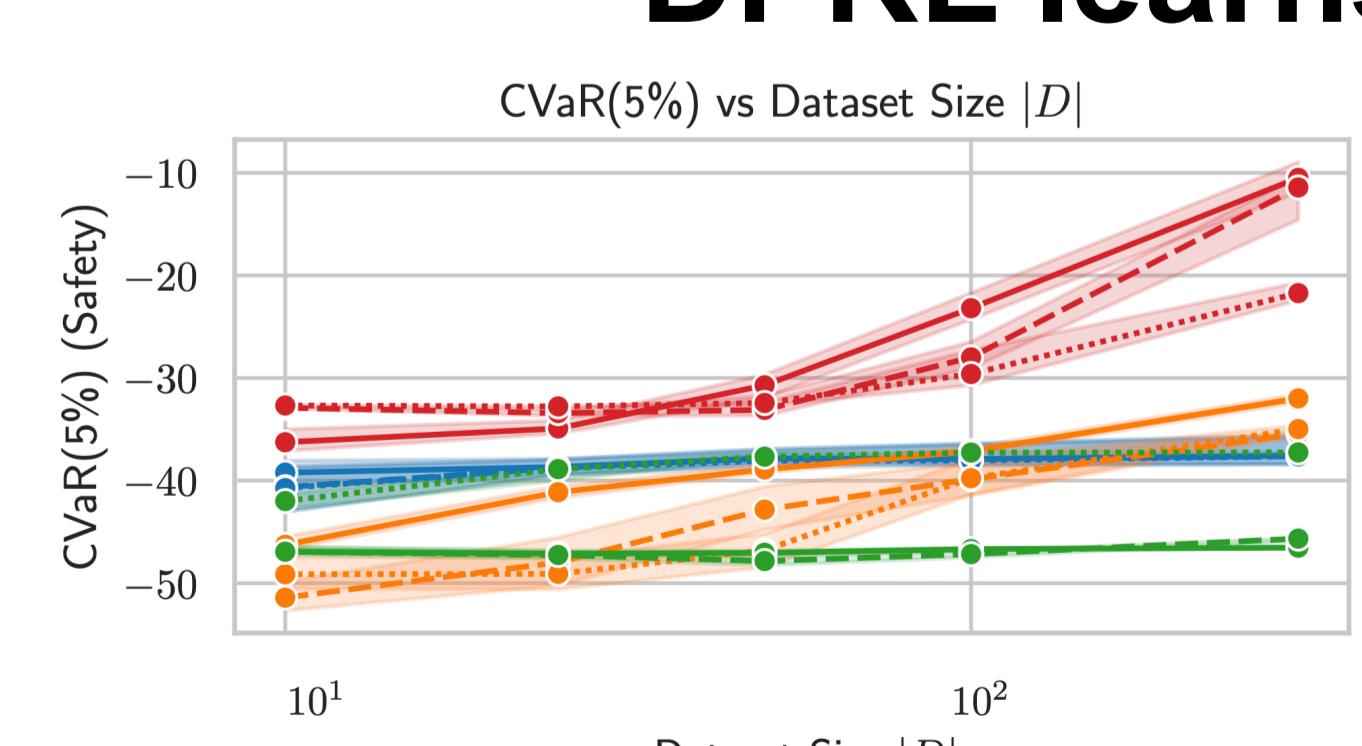


## THEORETICAL RESULTS



## EXPERIMENTAL RESULTS & CONCLUSION

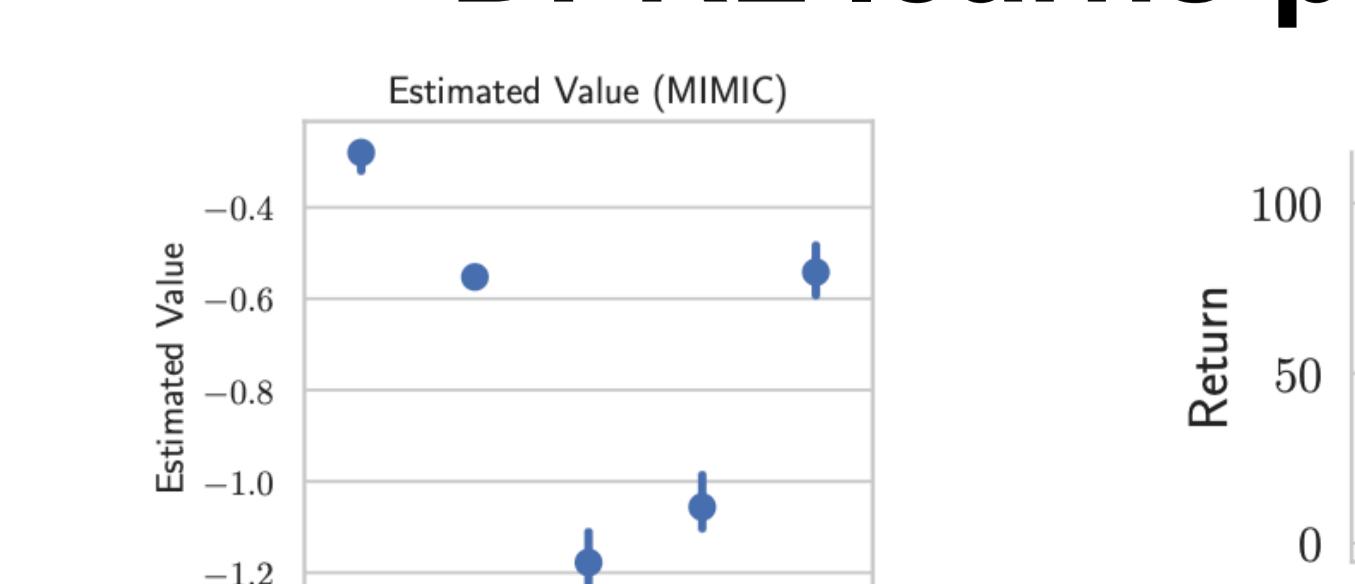
### DPRL learns safer policies



DPRL consistently safer (higher CVaR) than SPIBB[1], CQL[2], PQI[3]

Algorithm	$N_{\wedge}$	Algorithm	$N_{\wedge}$	Algorithm	$N_{\wedge}$	Algorithm	$N_{\wedge}$
DPRl	1	CQL	1	SPIBB	1	PQI	1
DPRl	5	CQL	5	SPIBB	5	PQI	5
DPRl	10	CQL	10	SPIBB	10	PQI	10
							0.001
							0.01
							0.1

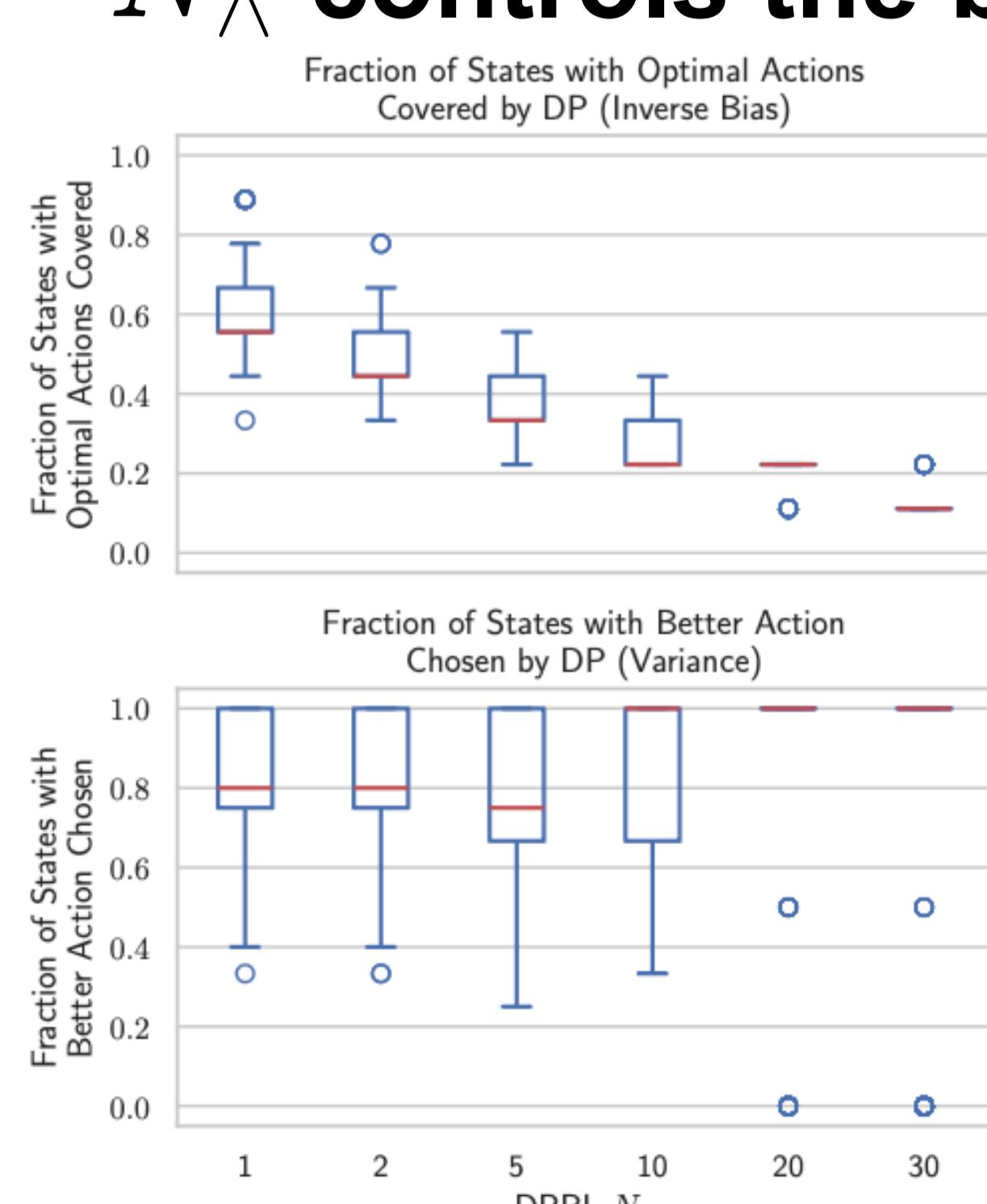
### DPRL learns performant policies



Atari Environments with suboptimal behavior (see paper for details)

Real-world MIMIC-IV clinical dataset with hypotensive patients

### $N_{\wedge}$ controls the bias-variance trade-off



Fraction of states where the optimal action is allowed decreases, leading to an increase in bias.

Fraction of states where an action with higher value is chosen increases.

This suggests reduced variance in value estimation of valid state-actions because we have more observations to estimate the value