

Abstract

Motivation

- In model-based offline RL, transition probabilities must be **estimated from limited data**.
- The *Percentile Criterion* (PC) is used to learn policies that are robust to uncertainty in transition probabilities.
- PC is optimized by constructing a special *uncertainty set* and **optimizing the policy against the worst model** in this set.
- Existing works use *Bayesian credible regions* as uncertainty sets which can be **unnecessarily large** and result in **overly conservative** policies!

Our Contributions:

- A **dynamic programming framework** that optimizes a lower bound on the percentile criterion without explicitly constructing ambiguity sets.
- Finite-sample and asymptotic analysis of the performance loss** of VaR framework.
- Theoretical Comparison of VaR ambiguity sets** with Bayesian credible regions-based ambiguity sets.
- Empirically analysis of the efficacy of the framework in several domains.

Problem

Notation:

- MDP: Markov Decision Process $(\mathcal{S}, \mathcal{A}, P, r, p_0, \gamma)$
- \tilde{P} : Random transition probabilities with posterior distribution f .
- v : Value Function
- $w_{s,a} = r_{s,a} + \gamma v$: 1-step discounted value

Percentile Criterion: Computes a policy π that maximizes the returns corresponding to the worst δ -percentile model.

$$\arg \max_{\pi \in \Pi, y \in \mathbb{R}} \left\{ y \mid \mathbb{P}_{\tilde{P} \sim f} [\rho(\pi, \tilde{P}) \geq y] > 1 - \delta \right\}. \quad (1)$$

Percentile Criterion is non-convex. Difficult to solve!

Existing solution: Solve it using Robust MDPs with Bayesian Credible Regions-based ambiguity sets

$$\arg \max_{\pi \in \Pi} \min_{P \in \mathcal{P}^{\text{BCR}}} \rho(\pi, P).$$

$$\forall s \in \mathcal{S}, a \in \mathcal{A}, \mathcal{P}_{s,a}^{\text{BCR}} = \mathcal{P}_{s,a}(\mathbf{w}, \psi, q) = \left\{ \mathbf{p} \in \Delta^S \mid \|\mathbf{p}_{s,a} - \tilde{\mathbf{p}}_{s,a}\|_{q, \mathbf{w}} \leq \psi_{s,a} \right\},$$

where weights \mathbf{w} and size $\psi_{s,a} \forall s \in \mathcal{S}, a \in \mathcal{A}$ are optimized to minimize span of ambiguity sets while satisfying $\mathbb{P}[\rho(\pi, \tilde{P}) \geq \min_{P \in \mathcal{P}^{\text{BCR}}} \rho(\pi, P)] \geq 1 - \delta$.

These sets can be unnecessarily large and result in conservative policies!

Solution: VaR Framework

Define the *VaR* Bellman optimality operator \mathcal{T}_{VaR} as

$$(\mathcal{T}_{VaR} \mathbf{v})(s) = \max_{a \in \mathcal{A}} VaR_{\alpha} [\tilde{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a}].$$

where VaR_{α} of a bounded random variable \tilde{X} with a CDF function $F: \mathbb{R} \rightarrow [0, 1]$ is defined as $VaR_{\alpha}[\tilde{X}] = \inf \{z \in \mathbb{R} \mid F(z) > \alpha\}$.

Theoretically Proven Properties of VaR Framework

- \mathcal{T}_{VaR} is a γ -**Contraction mapping** when reward is independent of next-state and $\sqrt{\gamma}$ -**Contraction mapping** otherwise.
- VaR Framework **optimizes a lower bound on the percentile-criterion**.
- The optimal VaR policy $\hat{\pi}$ solves

$$\max_{\pi \in \Pi^D} \min_{P \in \mathcal{P}^{VaR, \hat{\pi}}} \rho(\pi, P).$$

where for any value function \mathbf{v} , the *VaR* ambiguity set $\mathcal{P}^{VaR, \mathbf{v}}$ is defined as

$$\mathcal{P}^{VaR, \mathbf{v}} = \times_{s \in \mathcal{S}, a \in \mathcal{A}} \mathcal{P}_{s,a}^{VaR, \mathbf{v}} \text{ where } \mathcal{P}_{s,a}^{VaR, \mathbf{v}} = \left\{ \mathbf{p}_{s,a} \in \Delta^S \mid \mathbf{p}_{s,a}^{\top} \mathbf{v} \geq VaR_{\alpha} [\tilde{\mathbf{p}}_{s,a}^{\top} \mathbf{v}] \right\},$$

and $\hat{\mathbf{v}}^{\pi}$ is the fixed point of the *VaR* Bellman evaluation operator \mathcal{T}_{VaR}^{π} for each $\pi \in \Pi^D$.

Performance Error bounds for VaR Framework

- $\hat{\mathbf{v}}$: fixed point of the *VaR* Bellman optimality operator \mathcal{T}_{VaR}
- π^* : optimal policy in (1)
- $\rho^* = VaR_{\alpha} [\rho(\pi^*, \tilde{P})]$: optimal percentile returns
- $\delta \in (0, 1)$ - confidence level
- $\alpha: \delta/(2SA)$ in \mathcal{T}_{VaR}
- $I(\mathbf{p}_{s,a}^*)^{-1}$: Fisher Information matrix of the true transition probabilities $\mathbf{p}_{s,a}^*$
- $\sigma_{\max} = \max_{s \in \mathcal{S}, a \in \mathcal{A}} \sqrt{\hat{\mathbf{w}}_{s,a}^{\top} I(\mathbf{p}_{s,a}^*)^{-1} \hat{\mathbf{w}}_{s,a}}$: maximum asymptotic standard deviation of the returns estimate $\tilde{\mathbf{p}}_{s,a}^{\top} \hat{\mathbf{w}}_{s,a}$ for state-action pair (s, a)

Theorem 1: [Performance Error] With probability at least $1 - \delta$, the performance loss with respect to ρ^* is

$$\rho^* - \hat{\rho} \leq \frac{1}{1 - \gamma} \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} (VaR_{1-\alpha} [\tilde{\mathbf{p}}_{s,a}^{\top} \hat{\mathbf{w}}_{s,a}] - VaR_{\alpha} [\tilde{\mathbf{p}}_{s,a}^{\top} \hat{\mathbf{w}}_{s,a}]). \quad (2)$$

Theorem 2: [Asymptotic Performance Error] With probability at least $1 - \delta$, the asymptotic performance of the *VaR* framework $\hat{\rho}$ w.r.t. the optimal percentile returns ρ^* satisfies

$$\lim_{N \rightarrow \infty} \sqrt{N}(\rho^* - \hat{\rho}) \leq \frac{1}{1 - \gamma} (2\Phi^{-1}(1 - \alpha)\sigma_{\max}) \leq \frac{1}{1 - \gamma} \sqrt{8 \ln(1/\alpha)} \sigma_{\max}.$$

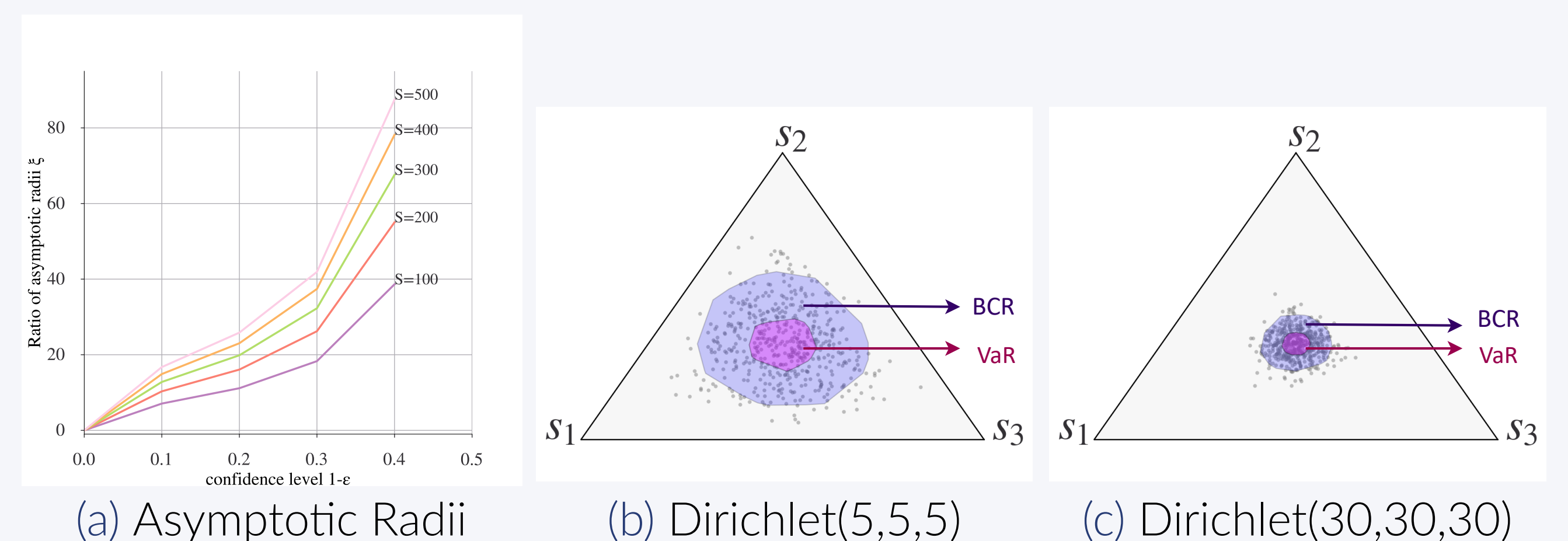
Comparison with Bayesian Credible Regions

- $\tilde{P} = (\tilde{\mathbf{p}}_{s,a})_{s \in \mathcal{S}, a \in \mathcal{A}}$: Maximum Likelihood Estimate (MLE) of transition probabilities computed from data \mathcal{D}
- $\Sigma = (I(\mathbf{p}_{s,a}^*)^{-1})_{s \in \mathcal{S}, a \in \mathcal{A}}$: covariance matrix of \tilde{P} .

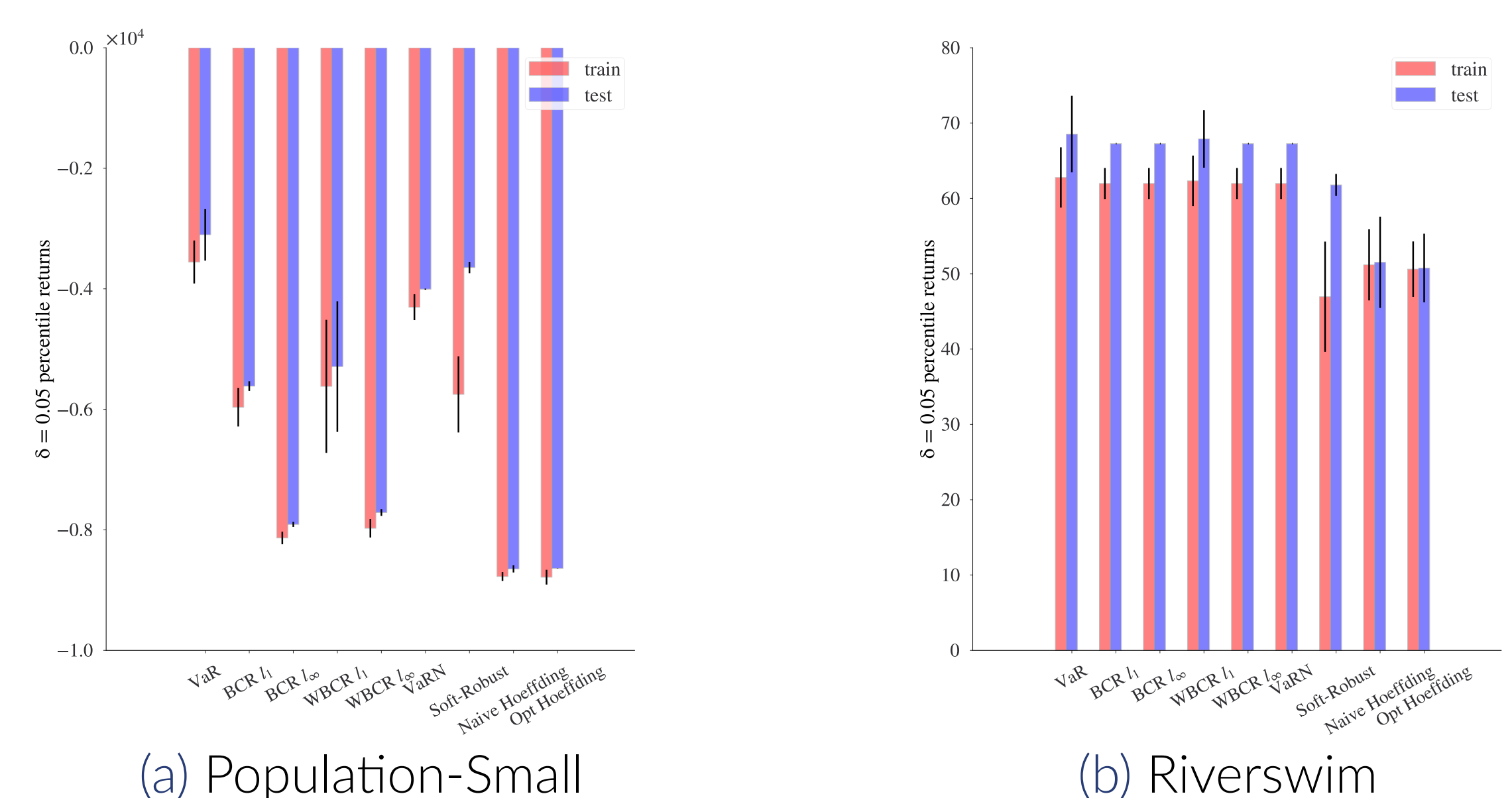
Theorem 3: [Asymptotic radius of VaR Ambiguity Sets] For all $s \in \mathcal{S}, a \in \mathcal{A}$,

$$\lim_{N \rightarrow \infty} \sqrt{N}(\mathcal{P}_{s,a}^{VaR} - \tilde{\mathbf{p}}_{s,a}) = \left\{ \mathbf{p}_{s,a} \in \Delta^S \mid \|\mathbf{p}_{s,a} - \tilde{\mathbf{p}}_{s,a}\|_{\Sigma_{s,a}^{-1}} \leq \Phi^{-1}(1 - \alpha) \right\} - \tilde{\mathbf{p}}_{s,a}.$$

Theorem 4: [Asymptotic radius of Bayesian Credible Region] For any state s and action a , if $\mathcal{P}_{s,a}^{\text{BCR}}$ be any Bayesian credible region and $\xi < \sqrt{\chi_{S,1-\alpha}^2 / \Phi^{-1}(1-\alpha)}$, then, $\forall s \in \mathcal{S}, a \in \mathcal{A}, \lim_{N \rightarrow \infty} \sqrt{N}(\mathcal{P}_{s,a}^{\text{BCR}} - \tilde{\mathbf{p}}_{s,a}) \not\subseteq \lim_{N \rightarrow \infty} \sqrt{N}\xi(\mathcal{P}_{s,a}^{VaR} - \tilde{\mathbf{p}}_{s,a})$.



Empirical Results on Population and Riverswim Domains



References

- Bahram Behzadian, Reazul Hasan Russel, Marek Petrik, and Chin Pang Ho. Optimizing percentile criterion using robust MDPs. In 24th, pages 1009–1017, 2021.
- R. Tyrrell Rockafellar. *Convex analysis*. Princeton University Press, 1970.