
Percentile Criterion Optimization in Offline Reinforcement Learning

Cyrus Cousins

Department of Computer Science
University of Massachusetts Amherst
cbcousins@umass.edu

Elita A. Lobo

Department of Computer Science
University of Massachusetts Amherst
elobo@umass.edu

Marek Petrik

Department of Computer Science
University of New Hampshire
mpetrik@cs.unh.edu

Yair Zick

Department of Computer Science
University of Massachusetts Amherst
yzick@umass.edu

Abstract

In reinforcement learning, robust policies for high-stakes decision-making problems with limited data are usually computed by optimizing the *percentile criterion*. The percentile criterion is optimized by constructing an *uncertainty set* that contains the true model with high probability and optimizing the policy for the worst model in the set. Since the percentile criterion is non-convex, constructing uncertainty sets is often challenging. Existing works use *Bayesian credible regions* as uncertainty sets, but they are often unnecessarily large and result in learning overly conservative policies. To overcome these shortcomings, we propose a novel Value-at-Risk based dynamic programming algorithm to optimize the percentile criterion without explicitly constructing any uncertainty sets. Our theoretical and empirical results show that our algorithm implicitly constructs much smaller uncertainty sets and learns less conservative robust policies.

1 Introduction

Batch Reinforcement Learning (Batch RL) [18] is popularly used for solving sequential decision-making problems using limited data. These algorithms are crucial in high-stakes domains where exploration is either infeasible or expensive, and policies must be learned from limited data. In model-based Batch RL algorithms, transition probabilities are learned from the data as well. Due to insufficient data, these transition probabilities are often imprecise. Errors in transition probabilities can accumulate, resulting in low-performing policies that fail when deployed.

To account for the uncertainty in transition probabilities, prior works use Bayesian models [10, 29] to model uncertainty and optimize the policy to maximize the returns corresponding to the worst α -percentile transition probability model. These policies guarantee that the true expected returns will be at least as large as the optimal returns with high confidence. This technique is commonly referred to as the *percentile-criterion* optimization. Unfortunately, the percentile criterion is NP-hard to optimize. Thus, current work uses *Robust Markov Decision Processes* (RMDPs) to optimize a lower bound on the percentile criterion. An RMDP takes as input an uncertainty set that contains the true transition probability model with high confidence and finds a policy that maximizes the returns of the worst model in the uncertainty set.

Unfortunately, the percentile criterion is non-convex, and therefore, constructing uncertainty sets is a challenging problem. Existing work uses Bayesian credible regions (BCR) [29] as uncertainty sets. However, these uncertainty sets are often unnecessarily large [12, 29] and result in learning

conservative robust policies. Some of the recent work approximates uncertainty sets using various heuristics [3, 29], but we show that they remain too conservative. Thus, the question of the *best possible near-optimal set that can be constructed for the percentile criterion* remains unanswered.

Our Contributions In this paper, we answer two important questions: a) *Are Bayesian credible regions the most optimal ambiguity sets for optimizing the percentile criterion?* b) *Can we obtain a less conservative solution to the percentile criterion without explicitly constructing ambiguity sets?* Our theoretical findings show that Bayesian credible regions can grow significantly with the number of states and therefore, tend to be unnecessarily large, resulting in highly conservative policies. As our main contribution, we provide a dynamic programming framework (Section 3), which we name the VaR framework, for optimizing a lower bound on the percentile criterion without explicitly constructing ambiguity sets. Specifically, we propose a new robust Bellman operator, the Value at Risk (VaR) Bellman operator, for optimizing the percentile criterion. We show that it is a valid contraction mapping that optimizes a tighter lower bound on the percentile criterion, compared to RMDPs with BCR ambiguity sets (Section 3). We theoretically analyze and bound the performance loss of our framework (Section 3.1). We also show that there exists directions in which the Bayesian credible regions can grow unnecessarily large with the number of states in the MDP and possibly result in a conservative solution. On the other hand, the ambiguity sets implicitly optimized by the VaR Bellman operator tend to be smaller and are independent of the number of states (Section 4). Finally, we empirically demonstrate the efficacy of our framework in several domains (Section 5).

1.1 Related Work

Several works propose different methods for solving the percentile criterion, as well as other robust measures for handling uncertainty in the transition model estimates. Russel and Petrik [29] and Behzadian et al. [3] propose various heuristics for minimizing the size of the ambiguity sets constructed for the percentile-criterion. Russel and Petrik [29] propose a method that interleaves robust value iteration with ambiguity set size optimization. Behzadian et al. [3] propose an iterative algorithm that optimizes the weights of ℓ_1 and ℓ_∞ ambiguity sets while optimizing the robust policy. However, these methods still construct Bayesian credible sets which can be unnecessarily large and result in conservative policies, as we show in Section 5.

Other works consider partial correlations between uncertain transition model parameters to mitigate the conservativeness of learned policies [4, 12, 13, 21, 22]. These approaches mitigate the conservativeness of S- and SA-rectangular ambiguity sets by capturing correlations between the uncertainty and by limiting the number of times the uncertain parameters deviate from the mean parameters. Despite these heuristics, most of these works [2, 14, 29, 35] either rely on weak statistical concentration bounds to construct frequentist ambiguity sets, or use Bayesian credible regions as ambiguity sets. These sets still tend to be unnecessarily large [12, 29], resulting in conservative policies. Robust RL work [2, 11, 14, 20, 24, 37] proposes other robust measures for handling uncertainty in transition models; however, these approaches do not provide probabilistic guarantees on the expected returns, and compute overly conservative policies.

2 Preliminaries

In the standard reinforcement learning setting, a sequential decision task is modeled as a Markov Decision Process (MDP) [26, 33]. An MDP is a tuple $\langle \mathcal{S}, \mathcal{A}, P, r, d_0, \gamma \rangle$ that consists of (a) a set of states $\mathcal{S} = \{1, 2, \dots, S\}$, (b) a set of actions $\mathcal{A} = \{1, 2, \dots, A\}$, (c) a deterministic reward function $R: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$, (d) a transition probability function $P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta^S$, (e) an initial state distribution $\mathbf{p}_0 \in \Delta^S$, where Δ^S represents the S -dimensional probability simplex, and (f) a discount factor $\gamma \in [0, 1]$. We use $\mathbf{p}_{s,a}$ to denote the vector of transition probabilities $P(s, a, \cdot)$ corresponding to the state s and action a . Similarly, for any reward function R , we use $\mathbf{r}_{s,a}$ to denote the vector of rewards $R(s, a, \cdot)$ corresponding to state s and action a . A Markovian policy $\pi: \mathcal{S} \rightarrow \Delta^A$ maps each state s to a distribution over actions \mathcal{A} . In a general RL setting, the goal is to compute a policy π that maximizes the expected discounted return $\rho(\pi, P)$ over an infinite horizon,

$$\max_{\pi \in \Pi} \rho(\pi, P) = \max_{\pi \in \Pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) \mid s_0 \sim \mathbf{p}_0, a_t \sim \pi(s_t), s_{t+1} \sim P(s_t, a_t, \cdot) \right].$$

The value of a policy π at any state s is the discounted sum of rewards received by an RL agent, if it starts from state s , i.e., $v^\pi(s) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, s_{t+1}) | s_0 = s, a_t \sim \pi, s_{t+1} \sim P(s_t, a_t, \cdot)]$. We assume a *batch reinforcement learning* setting [18] where the reward function is known, but the true transition probabilities P^* are unknown. Following prior work on robust Bayesian RL [7, 10, 29, 36], we use parametric Bayesian models to represent uncertainty over the true transition probabilities P^* . Given a batch of sample data \mathcal{D} , one can derive a posterior distribution over the random variable representing transition probabilities \tilde{P} .

We use the tilde to indicate that \tilde{P} is a random variable. To avoid unnecessary computational technicalities, we assume that \tilde{P} is a discrete random variable taking on values $\tilde{P}(\omega), \omega \in \Omega$ for some $\Omega = \{1, \dots, N\}$ with a distribution f . That is, the random variable \tilde{P} represents a discrete approximation of the true, possibly continuous posterior, as is common in methods like Sample Average Approximation (SAA) [31].

Percentile Criterion The α -percentile criterion is popularly used to derive robust policies under model uncertainty [10]. It aims to compute a policy π that maximizes the returns corresponding to the worst α -percentile model:

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

The value y lower-bounds the true expected discounted returns with confidence $1 - \alpha$ where $\alpha \in (0, 0.5)$. Optimizing the percentile criterion is equivalent to optimizing the Value at Risk (VaR_α) of expected discounted returns when there exists uncertainty in transition probabilities \tilde{P} and the expected returns function ρ is lower-semicontinuous. The optimization in (1) is equivalent to

$$\max_{\pi \in \Pi} \text{VaR}_\alpha [\rho(\pi, \tilde{P})], \quad (2)$$

where VaR_α of a bounded random variable \tilde{X} with a CDF function $F: \mathbb{R} \rightarrow [0, 1]$ is defined as [27]

$$\text{VaR}_\alpha[\tilde{X}] = \inf \{z \in \mathbb{R} | F(z) > \alpha\}. \quad (3)$$

A lower value of α in (1) indicates a higher confidence in the returns achieved in expectation. For example, $\text{VaR}_{0.05}[\rho(\pi, \tilde{P})] = x$ indicates that the true returns will be at least equal to the robust returns x for 95% of the transition probability models. When clear from context, we use VaR to denote the Value at Risk at confidence level α . Unfortunately, the optimization problem in (1) is NP-hard to optimize and is usually approximately solved using Robust MDPs.

Robust MDPs Robust MDPs (RMDPs) generalize MDPs to account for uncertainty, or ambiguity, in the transition model. An *ambiguity set* for an RMDP is constructed such that it contains the true model with high confidence. The optimal policy of a Robust MDP π^* maximizes the returns of the worst model in the ambiguity set: $\pi^* = \arg \max_{\pi \in \Pi} \min_{P \in \mathcal{P}} \rho(\pi, P)$. General RMDPs are NP-hard to solve [35], but they are tractable for broad classes of ambiguity sets. The simplest such type is the SA-rectangular ambiguity set [25, 35], defined as

$$\mathcal{P} = \{P \in (\Delta^S)^{S \times A} \mid \mathbf{p}_{s,a} \in \mathcal{P}_{s,a}, \forall s \in \mathcal{S}, \forall a \in \mathcal{A}\},$$

for a given $\mathcal{P}_{s,a} \subseteq \Delta^S, s \in \mathcal{S}, a \in \mathcal{A}$. SA-rectangular ambiguity sets [3, 29] assume that the transition probabilities corresponding to each state-action pair are independent. Similarly to MDPs, the optimal robust value function $\mathbf{v}^* \in \mathbb{R}^S$ for an SA-rectangular RMDP is the unique fixed point of the robust Bellman optimality operator $\mathcal{T}: \mathbb{R}^S \rightarrow \mathbb{R}^S$ defined as $(\mathcal{T}\mathbf{v})(s) = \max_{a \in \mathcal{A}} \min_{\mathbf{p}_{s,a} \in \mathcal{P}_{s,a}} \mathbf{p}_{s,a}^\top (\mathbf{r}_{s,a} + \gamma \cdot \mathbf{v})$.

To optimize the percentile criterion, an SA-rectangular ambiguity set \mathcal{P} is constructed such that it contains the true model with high probability, and thus, the following equation holds.

$$\Pr \left[\rho(\pi, \tilde{P}) \geq \min_{P \in \mathcal{P}} \rho(\pi, P) \right] \geq 1 - \alpha.$$

Although RMDPs have been used to approximately solve the percentile criterion [3], the quality of the robust policies it computes depends mainly on the size of the ambiguity sets. The larger the ambiguity sets, the more conservative the robust policy [22]. SA-rectangular ambiguity sets are most

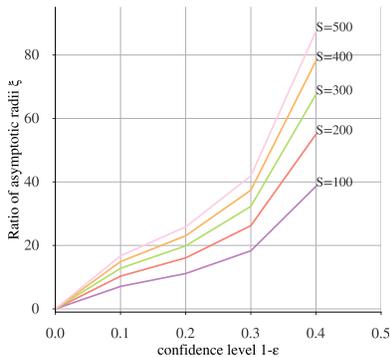


Figure 1: Compares the asymptotic radius of BCR ambiguity sets to VaR ambiguity sets. The asymptotic radius of the BCR ambiguity sets significantly grows with the number of states.

commonly studied; thus we focus our attention on SA-rectangular Robust MDPs. We investigate whether Bayesian credible regions are optimal ambiguity sets for optimizing the percentile criterion. We simply refer to SA-rectangular RMDPs and SA-rectangular ambiguity sets as Robust MDPs and ambiguity sets respectively.

Our work focuses on Bayesian (rather than frequentist) ambiguity sets. Bayesian ambiguity sets are usually constructed from Bayesian credible regions (BCR) [3, 29]. Given a state s and an action a , let $\psi_{s,a}$ represent the size of the BCR ambiguity sets; $\mathcal{P}_{s,a}^{\text{BCR}}$ and $\bar{p}_{s,a}$ represent the mean transition model. The set $\mathcal{P}_{s,a}^{\text{BCR}}$ is constructed as

$$\mathcal{P}_{s,a}^{\text{BCR}} = \mathcal{P}_{s,a}(\mathbf{b}, \psi, q) = \{ \mathbf{p}_{s,a} \in \Delta^S \mid \| \mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a} \|_{q,\mathbf{b}} \leq \psi_{s,a} \}, \quad (4)$$

where $q \in \{1, \infty\}$ represents the norm of the weighted ball in (4) and $\mathbf{b} \in \mathbb{R}_+^S$ is a weight vector. Here, \mathbf{b} is jointly optimized with $\psi \in \mathbb{R}$ to minimize the span of the ambiguity sets such that the true model is contained in the ambiguity set with high confidence, i.e., $\Pr(\tilde{\mathbf{p}}_{s,a} \in \mathcal{P}_{s,a}(\mathbf{b}, \psi, q)) \geq 1 - \alpha$. We refer to BCR ambiguity sets with non-uniform weights as weighted BCR ambiguity sets. We refer to the Robust Bellman optimality operator with BCR ambiguity sets \mathcal{T}_{BCR} as the BCR Bellman optimality operator, and to RMDPs with BCR ambiguity sets as BCR RMDPs. For any $\delta \in (0, 1)$, setting the confidence level α in \mathcal{T}_{BCR} to δ/SA for all state-action pairs yields $1 - \delta$ confidence on the returns of the optimal robust policy [3]. However, we show that even span-optimized BCR RMDPs can be sub-optimal for optimizing the percentile criterion.

We use the shorthand $\mathbf{w}_{s,a}$ for any $s \in \mathcal{S}$, $a \in \mathcal{A}$ to denote the vector of values associated with value $\mathbf{v} \in \mathbb{R}^S$ and the one-step transition from state s and action a , i.e., $\mathbf{w}_{s,a} = \mathbf{r}_{s,a} + \gamma \mathbf{v}$. We use $\bar{\mathbf{p}}_{s,a} \in \mathbb{R}^S$ and $\Sigma_{s,a} \in \mathbb{R}^{S \times S}$ for any $s \in \mathcal{S}$, $a \in \mathcal{A}$ to represent the empirical mean and covariance of transition model $\tilde{\mathbf{p}}_{s,a}$ estimated from \mathcal{D} . We use $\phi(\cdot)$ and $\Phi(\cdot)$ to represent the probability distribution function (PDF) and cumulative distribution function (CDF) respectively of the normal distribution with mean 0 and variance 1. The Z -Minkowski norm $\| \mathbf{x} \|_Z$ for a vector \mathbf{x} given some positive-definite matrix Z is defined as $\| \mathbf{x} \|_Z = \sqrt{\mathbf{x}^\top Z^{-1} \mathbf{x}}$.

Example 2.1. Consider an MDP with four states $\{s_0, s_1, s_2, s_3\}$ and a single action $\{a_0\}$. The state s_0 is the initial state and the states s_1, s_2, s_3 are terminal states with zero rewards. The posterior of the transition probability $\tilde{\mathbf{p}}_{s_0,a_0}$ follows a Dirichlet distribution $Dir(10, 10, 1)$ with mean $[0.48, 0.48, 0.04]$. The rewards for transitions from state 0 are given by $\mathbf{r}_{s_0,a_0} = [0.25, 0.25, -1]$.

We wish to optimize the percentile criterion with confidence level $\delta = 0.2$. Following the sampling procedure proposed by Russel and Petrik [29] to construct a uniformly weighted BCR ambiguity set for $\tilde{\mathbf{p}}_{s_0,a_0}$ with 100 posterior samples, yields an ambiguity set $\mathcal{P}_{s_0,a_0}^{\text{BCR}} = \{ \mathbf{p} \in \Delta^S \mid \| \mathbf{p} - \tilde{\mathbf{p}}_{s_0,a_0} \|_1 \leq 0.277 \}$. In this case, the reward estimate against the worst model in the ambiguity set $\mathbf{p} = [0.50, 0.32, 0.18]$ is $\rho^{\text{BCR}} = 0.025$. Since we have a single non-terminating state in the MDP, the percentile returns are given by $\rho^{\text{VaR}_\alpha} = \text{VaR}_{0.2}[\tilde{\mathbf{p}}_{s_0,a_0}^\top \mathbf{r}_{s_0,a_0}]$. Computing ρ^{VaR_α} for Dirichlet distribution $Dir(10, 10, 1)$, we get $\rho^{\text{VaR}_\alpha} = 0.17 > \rho^{\text{BCR}}$. Thus, this example shows that BCR ambiguity sets can be unnecessarily large and, thus, result in conservative policies.

3 VaR Framework

We introduce the VaR Bellman optimality operator \mathcal{T}_{VaR} for approximately solving the percentile criterion. We show that \mathcal{T}_{VaR} is a valid Bellman operator: it is a contraction mapping and lower bounds the percentile criterion. For any value function $\mathbf{v} \in \mathbb{R}^{\mathcal{S}}$, state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$, we define the VaR Bellman optimality operator \mathcal{T}_{VaR} as

$$(\mathcal{T}_{\text{VaR}}\mathbf{v})(s) = \max_{a \in \mathcal{A}} \text{VaR}_{\alpha} [\tilde{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a}] . \quad (5)$$

For each state s , \mathcal{T}_{VaR} maximizes the value corresponding to the worst α -percentile model. In contrast to the BCR Bellman optimality operator \mathcal{T}_{BCR} , computing \mathcal{T}_{VaR} does not require constructing ambiguity sets from confidence regions; it can simply be estimated from samples of the model posterior distribution, as we later show.

In Proposition A.4, we formally prove that \mathcal{T}_{VaR} is a contraction mapping, and thus has a unique fixed point. This fixed point is the value of the optimal policy that maximizes a tight lower bound on the percentile criterion.

We now show that the VaR Bellman optimality operator \mathcal{T}_{VaR} optimizes a lower bound on the percentile criterion. Given a policy π , a state $s \in \mathcal{S}$, and a transition model P , let

$$(\mathcal{T}^{\pi}\mathbf{v})(s) = \sum_{a \in \mathcal{A}} \pi(s, a) \tilde{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a} ,$$

represent the Bellman evaluation operator for transition model P . Furthermore, let

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

represent the VaR Bellman evaluation operator for random transition model \tilde{P} . We use $\hat{\mathbf{v}}^{\pi}$ and \mathbf{v}^{π} to denote the fixed points of $\mathcal{T}_{\text{VaR}}^{\pi}$ and \mathcal{T}^{π} respectively, and $\tilde{\mathbf{v}}^{\pi}$ to represent the random value function of policy π computed using a realization \tilde{P} of the posterior distribution of the transition model P .

Proposition 3.1 (Lower Bound Percentile Criterion). *For any $\delta \in (0, 1)$, if we set the confidence level α in the operator $\mathcal{T}_{\text{VaR}}^{\pi}$ to δ/s , then for every policy $\pi \in \Pi : \Pr_{\tilde{P}} [\hat{\mathbf{v}}^{\pi} \preceq \tilde{\mathbf{v}}^{\pi} | \mathcal{D}] \geq 1 - \delta$, where \tilde{P} is a realization of the posterior distribution of the transition model P conditioned on observed transitions \mathcal{D} .*

See Appendix B.1 for the proof. Proposition 3.1 shows that for any policy π and state s , the VaR value at state s , $\hat{\mathbf{v}}^{\pi}(s)$ lower bounds the true value $\tilde{\mathbf{v}}^{\pi}(s)$ with high confidence. Comparing Proposition 3.1 with the definition of the percentile-criterion in (1), it is easy to see that the percentile-criterion requires confidence guarantees only on the returns computed from the initial states, whereas the equation in Proposition 3.1 provides confidence guarantees on the value of every state. Therefore, for any policy π , the value $\mathbf{p}_0^{\top} \hat{\mathbf{v}}^{\pi}$ is a lower bound on the percentile-criterion objective $\text{VaR}_{\delta}[\rho(\pi, \tilde{P})]$. Since \mathcal{T}_{VaR} finds a policy π that maximizes the value $\mathbf{p}_0^{\top} \hat{\mathbf{v}}^{\pi}$, it follows [26] that \mathcal{T}_{VaR} optimizes a lower bound on the percentile criterion in (1).

Proposition 3.2. *Suppose that $\tilde{\mathbf{p}}_{s,a}$ for any state s and action a , is a multivariate sub-Gaussian with mean $\bar{\mathbf{p}}_{s,a}$ and covariance factor $\Sigma_{s,a}$, i.e., $\mathbb{E} \left[\exp \left(\lambda (\tilde{\mathbf{p}}_{s,a} - \bar{\mathbf{p}}_{s,a})^{\top} \mathbf{w} \right) \right] \leq \exp \left(\lambda^2 \mathbf{w}^{\top} \Sigma_{s,a} \mathbf{w} / 2 \right)$, $\forall \lambda \in \mathbb{R}, \forall \mathbf{w} \in \mathbb{R}^{\mathcal{S}}$. Then, for any state $s \in \mathcal{S}$, \mathcal{T}_{VaR} satisfies*

$$(\mathcal{T}_{\text{VaR}}\mathbf{v})(s) \geq \max_{a \in \mathcal{A}} \left(\bar{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a} - \sqrt{2 \ln(1/\alpha)} \sqrt{\mathbf{w}_{s,a}^{\top} \Sigma_{s,a} \mathbf{w}_{s,a}} \right) .$$

See Appendix B.2 for the proof. Proposition 3.2 shows that by assuming that the transition probabilities are sub-Gaussian, we can easily compute a lower bound of the VaR Bellman update $(\mathcal{T}_{\text{VaR}}\mathbf{v})$ for a given value function using only the mean and the covariance matrix of \tilde{P} .

In the following proposition, we provide the exact form of the VaR_{α} Bellman optimality operator \mathcal{T}_{VaR} under normal conditions, which is a special case of Proposition 3.2.

Proposition 3.3. *Suppose that \tilde{P} is normally distributed, i.e., for any state s and action a , $\tilde{\mathbf{p}}_{s,a} \sim \mathcal{N}(\bar{\mathbf{p}}_{s,a}, \Sigma_{s,a})$.*

Then, \mathcal{T}_{VaR} for any state $s \in \mathcal{S}$ takes the form

$$(\mathcal{T}_{\text{VaR}}\mathbf{v})(s) = \max_{a \in \mathcal{A}} \left(\bar{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a} - \Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^{\top} \Sigma_{s,a} \mathbf{w}_{s,a}} \right) .$$

3.1 Performance Guarantees

We now derive finite-sample and asymptotic bounds on the loss of the VaR framework.

Theorem 3.4 (Performance). *Let \hat{v} be the fixed point of the VaR Bellman optimality operator \mathcal{T}_{VaR} , and π^* be the optimal policy in (1). Let $\rho^* = \text{VaR}_\alpha \left[\rho(\pi^*, \tilde{P}) \right]$ denote the optimal percentile returns and $\hat{\rho} = \mathbf{p}_0^\top \hat{v}$ denote the lower bound on the percentile returns computed using the Bellman operator \mathcal{T}_{VaR} . For any $\delta \in (0, 1)$, we set the confidence level $\alpha = \delta/(2SA)$ in \mathcal{T}_{VaR} . Then, with probability at least $1 - \delta$, the performance loss with respect to ρ^* is*

$$\rho^* - \hat{\rho} \leq \frac{1}{1 - \gamma} \max_{s \in S} \max_{a \in A} (\text{VaR}_{1-\alpha}[\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}] - \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}]) . \quad (6)$$

See Appendix B.4 for the proof. Theorem 3.4 bounds the finite sample performance loss of the VaR framework. The loss varies proportionally to the maximum difference between the α and $1 - \alpha$ percentile of the one-step Bellman update for the optimal robust value function \hat{v} . As expected, the VaR framework performs better when the uncertainty in the transition models is small.

Theorem 3.5 (Asymptotic Performance). *For any $\delta \in (0, 1)$, set $\alpha = \delta/(2SA)$ in \mathcal{T}_{VaR} . Let $I(\mathbf{p}_{s,a}^*)^{-1}$ for any state s and action a , be the Fisher Information matrix corresponding to the true transition probabilities $\mathbf{p}_{s,a}^*$. Furthermore, let $\sigma_{\max} = \max_{s \in S, a \in A} \sqrt{\hat{\mathbf{w}}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \hat{\mathbf{w}}_{s,a}}$ represent the maximum asymptotic standard deviation of the returns estimate $\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}$ for any state-action pair (s, a) . Then, with probability at least $1 - \delta$, the asymptotic performance of the VaR framework $\hat{\rho}$ w.r.t. the optimal percentile returns ρ^* satisfies*

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

See Appendix B.5 for the proof. Theorem 3.5 shows that almost surely, the asymptotic loss in performance of the VaR framework convergence to 0, i.e., $\lim_{N \rightarrow \infty} (\rho^* - \hat{\rho}) = 0$.

3.2 Dynamic Programming Algorithm

We provide a detailed description of the VaR value iteration algorithm (Algorithm 3.1) below. We also bound the number of samples required to estimate the VaR Bellman update $(\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})(s)$ for any given policy π and state s with high confidence $1 - \zeta$.

Algorithm 3.1: Generalized VaR Value Iteration Algorithm

Input: Confidence α , Posterior distribution f , target Bellman residual ε

Output: Robust policy π , lower bound \mathbf{v}

- 1 Initialize π with arbitrary π_0 , robust value-function \mathbf{v} with arbitrary \mathbf{v}_0 , $k = 0$;
 - 2 Sample N models $\tilde{P}(\omega_1), \tilde{P}(\omega_2), \dots, \tilde{P}(\omega_N)$ from posterior f ;
 - 3 **repeat**
 - 4 **for** $s \leftarrow 1$ **to** S **do**
 - 5 Initialize $\mathbf{q} \leftarrow []$;
 - 6 **for** $a \leftarrow 1$ **to** A **do**
 - 7 $\mathbf{q}[a] \leftarrow \widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top (\mathbf{r}_{s,a} + \gamma \mathbf{v}_k)]$ OR ;
 - 8 $\mathbf{q}[a] \leftarrow \tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a}}$ (under normal Assumptions);
 - 9 **end**
 - 10 $\mathbf{v}_k(s) \leftarrow \max(\mathbf{q})$; $\pi_k(s) \leftarrow \arg \max(\mathbf{q})$;
 - 11 **end**
 - 12 $k \leftarrow k + 1$;
 - 13 **until** $\|\mathbf{v}_k - \mathbf{v}_{k-1}\|_\infty \leq \varepsilon$;
 - 14 **return** π_k, \mathbf{v}_k ;
-

In each iteration of Algorithm 3.1, we compute the one-step VaR Bellman update $\mathcal{T}_{\text{VaR}}(\mathbf{v})$ using the current value function \mathbf{v} . When \tilde{P} is not normally distributed, we use the Quick Select algorithm [16] to efficiently compute the empirical estimate of the α -percentile of returns for any state s and action a , i.e., $\widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top (\mathbf{r}_{s,a} + \gamma \mathbf{v})]$ in $\mathcal{O}(SAN)$ time (Proposition 3.6). On the other hand, when \tilde{P} is

normally distributed, we compute the VaR Bellman update $\mathcal{T}_{\text{VaR}}(v)$ using the empirical estimate of mean $(\bar{\mathbf{p}}_{s,a})_{s \in \mathcal{S}, a \in \mathcal{A}}$ and covariance $(\Sigma_{s,a})_{s \in \mathcal{S}, a \in \mathcal{A}}$ of the transition probabilities derived from the data \mathcal{D} (Proposition 3.3). We repeat these steps until convergence.

Proposition 3.6 (Time Complexity). *The time complexity of a single iteration of the loop in line 3 of the VaR Value Iteration (Algorithm 3.1) is in $\mathcal{O}(SAN)$, where N is the number of samples of the posterior samples of \tilde{P} .*

Proposition 3.7 (Empirical Error Bound). *For any state s , action a and value function v , let $\widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ represent the empirical estimate of α -percentile of returns $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ and Φ_f represent the cumulative density function (CDF) of the random estimate of returns $\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}$. Suppose that Φ_f is differentiable at the point $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ and let $m = \Phi'(\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}])$ represents the density of estimate of returns at point $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$. Let N^* be the number of posterior samples required to obtain empirical error $\varepsilon \in \mathbb{R}$, with confidence $1 - \zeta$, i.e., $\Pr \left[\left| \widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] - \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] \right| > \varepsilon \right] \leq \zeta$. Then, $\lim_{\varepsilon \rightarrow 0} N^* \varepsilon^2 = \ln(2/\zeta)/2m^2$.*

4 Comparison with Bayesian Credible Regions

We are now ready to answer the question: *Are Bayesian credible regions the optimal ambiguity sets for optimizing the percentile criterion?* For this, we compare the VaR framework with BCR Robust MDPs. First, we derive the robust form of the VaR framework and show that in contrast to the BCR Bellman operator, the VaR Bellman optimality operator implicitly constructs value function dependent ambiguity sets and thus, these sets tend to be smaller (Proposition 4.1). Then, we show that the solution given by the VaR operator \mathcal{T}_{VaR} is never worse than the solution given by the BCR operator \mathcal{T}_{BCR} (Proposition 4.2). Finally, we compare the asymptotic radii of the BCR ambiguity sets and the VaR ambiguity sets implicitly constructed by \mathcal{T}_{VaR} . For any given confidence level α , the radius of the VaR ambiguity sets are asymptotically smaller than BCR ambiguity sets (Theorem 4.4). Precisely, the ratio of the radii of VaR ambiguity sets to BCR ambiguity sets is at least $\sqrt{\chi_{S,1-\alpha}^2 / \Phi^{-1}(1-\alpha)}$, where $\chi_{S,1-\alpha}^2$ is the CDF inverse of $1 - \alpha$ percentile of Chi-squared distribution with degree of freedom S and $\Phi^{-1}(1 - \alpha)$ is the $1 - \alpha$ percentile of $\mathcal{N}(0, 1)$. This implies that there exists directions in which the BCR ambiguity sets is atleast $\Omega(\sqrt{S})$ larger than VaR ambiguity sets. Thus, we prove that VaR framework is better suited for optimizing the percentile criterion than BCR RMDPs.

For any value function v , define the VaR ambiguity set $\mathcal{P}^{\text{VaR},v}$ as

$$\mathcal{P}^{\text{VaR},v} = \bigtimes_{s \in \mathcal{S}, a \in \mathcal{A}} \mathcal{P}_{s,a}^{\text{VaR},v} \text{ where } \mathcal{P}_{s,a}^{\text{VaR},v} = \{ \mathbf{p}_{s,a} \in \Delta^S \mid \mathbf{p}_{s,a}^\top \mathbf{v} \geq \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{v}] \}. \quad (7)$$

Proposition 4.1 (Equivalence). *Let \hat{v}^π be the fixed point of the VaR Bellman evaluation operator \mathcal{T}_{VaR} for each $\pi \in \Pi^D$, i.e., $\hat{v}^\pi = (\mathcal{T}_{\text{VaR}}^\pi \hat{v}^\pi)$, where Π^D is the set of all deterministic policies. Then, the optimal VaR policy $\hat{\pi}$ solves*

$$\max_{\pi \in \Pi^D} \min_{P \in \mathcal{P}^{\text{VaR},\hat{v}^\pi}} \rho(\pi, P). \quad (8)$$

See Appendix B.8 for the proof. Proposition 4.1 shows that the VaR Bellman optimality operator optimizes a unique robust MDP whose ambiguity sets are SA-rectangular and policy dependent. Notice that for any state s and action a , the ambiguity set is a half-plane $\{ \mathbf{p}_{s,a} \in \mathbb{R}^S : \mathbf{p}_{s,a}^\top \mathbf{v}^\pi \leq \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{v}^\pi] \}$ dependent on the value function \mathbf{v}^π of the current policy π . In contrast, BCR ambiguity sets are independent of any policy or value function and are constructed such that they provide high-confidence guarantees on returns of all policies simultaneously. As a result, BCR ambiguity sets tend to be unnecessarily large.

Proposition 4.2. *For any policy π , the fixed point of the VaR policy evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ dominates the fixed point of the Bellman evaluation operator $\mathcal{T}_{\text{BCR}}^\pi$, i.e., $\mathcal{T}_{\text{BCR}}^\pi \cdots \mathcal{T}_{\text{BCR}}^\pi \mathbf{v} \preceq \mathcal{T}_{\text{VaR}}^\pi \cdots \mathcal{T}_{\text{VaR}}^\pi \mathbf{v}$ for any \mathbf{v} . Similar results hold for policy optimization operators \mathcal{T}_{VaR} and \mathcal{T}_{BCR} , i.e., $\mathcal{T}_{\text{BCR}} \cdots \mathcal{T}_{\text{BCR}} \mathbf{v} \preceq \mathcal{T}_{\text{VaR}} \cdots \mathcal{T}_{\text{VaR}} \mathbf{v}$ for any \mathbf{v} .*

Proposition 4.2 proves that the lower bound on the percentile-criterion optimized by the VaR Bellman operator \mathcal{T}_{VaR} is never worse than the lower-bound optimized by the BCR Bellman operator \mathcal{T}_{BCR} .

Although this is a weaker claim, in the results that follow, we show that asymptotically, the BCR ambiguity sets tend to grow unnecessarily large with the number of states, and therefore, BCR RMDPs are more susceptible to generating overly conservative solutions.

We now compute the asymptotic radii of BCR ambiguity sets and VaR ambiguity sets.

The Bernstein Von-Mises Theorem [34] establishes the asymptotic properties of the posterior distribution of \tilde{P} constructed from N independent samples. We assume that this theorem holds for transition probabilities \tilde{P} for computing the asymptotic radii of the VaR and BCR ambiguity sets.

Theorem 4.3 (Asymptotic Radii of VaR Ambiguity Sets). *Let $\bar{P} = (\bar{p}_{s,a})_{s \in \mathcal{S}, a \in \mathcal{A}}$ be the Maximum Likelihood estimate of transition probabilities computed from data \mathcal{D} and $\Sigma = (I(\mathbf{p}_{s,a}^*))^{-1}_{s \in \mathcal{S}, a \in \mathcal{A}}$ be the corresponding covariance matrix. Then, $\forall s \in \mathcal{S}, a \in \mathcal{A}$,*

$$\lim_{N \rightarrow \infty} \sqrt{N}(\mathcal{P}_{s,a}^{\text{VaR}} - \bar{p}_{s,a}) = \left\{ \mathbf{p}_{s,a} \in \Delta^S \mid \|\mathbf{p}_{s,a} - \bar{p}_{s,a}\|_{\Sigma_{s,a}^{-1}} \leq \Phi^{-1}(1 - \alpha) \right\} - \bar{p}_{s,a} . \quad (9)$$

Theorem 4.3 shows that the asymptotic form of the VaR ambiguity set is an ℓ_2 ellipsoid ball with radius $\frac{\Phi^{-1}(1-\alpha)}{\sqrt{N}}$. It is important to notice that, in contrast to the finite-sample VaR ambiguity set $\mathcal{P}_{s,a}^{\text{VaR}, v^\pi}$ in problem (7), the asymptotic VaR ambiguity set $\mathcal{P}_{s,a}^{\text{VaR}}$ is independent of the value function \mathbf{v} . This is because $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{v}^\pi]$ is convex in the value function \mathbf{v} , [28] and as a consequence, the asymptotic ambiguity set $\mathcal{P}_{s,a}^{\text{VaR}}$ is simply the intersection of closed half-planes in $\mathcal{P}_{s,a}^{\text{VaR}, v^\pi}$, computed over all value functions $\mathbf{v} \in \mathbb{R}^S$ [8]. It is also worth noting that the radius of the asymptotic VaR ambiguity set $\mathcal{P}_{s,a}^{\text{VaR}}$ is a constant. In contrast, the asymptotic radius of the BCR ambiguity sets grows with the number of states, as we show in the following proposition.

Theorem 4.4 (Asymptotic Radius of Bayesian Credible Regions). *For any state s and action a , let $\mathcal{P}_{s,a}^{\text{BCR}}$ be any Bayesian credible region and $\bar{\mathbf{p}}_{s,a}$ be the maximum likelihood estimator based on data \mathcal{D} . Furthermore, $\xi < \sqrt{\chi_{S,1-\alpha}^2}/\Phi^{-1}(1-\alpha)$. Then, $\forall s \in \mathcal{S}, a \in \mathcal{A}$,*

$$\lim_{N \rightarrow \infty} \sqrt{N}(\mathcal{P}_{s,a}^{\text{BCR}} - \bar{\mathbf{p}}_{s,a}) \not\subseteq \lim_{N \rightarrow \infty} \sqrt{N}\xi(\mathcal{P}_{s,a}^{\text{VaR}} - \bar{\mathbf{p}}_{s,a}) . \quad (10)$$

See Appendix B.11 for the proof. We note that Theorem 4.4 is an adaptation of Theorem 10 in [15] in RL which proves that there exists directions in which the Bayesian credible regions is at least $\xi = \sqrt{\chi_{S,1-\alpha}^2}/\Phi^{-1}(1-\alpha)$ larger than VaR ambiguity sets. Since the value of ξ only grows with the number of states, we conclude that BCR ambiguity sets are sub-optimal for optimizing the percentile criterion. Figure 1 shows the growth in the ratio of radius of BCR to VaR ambiguity sets with an increasing number of states.

5 Experiments

We now empirically analyze the robustness of the VaR framework in three different domains. Specifically, we investigate if the VaR_α framework learns robust policies that are less conservative than BCR Robust MDPs.

Riverswim: The Riverswim MDP [32] consists of five states and two actions. The state represents the coordinates of the swimmer in the river and action represents the direction of the swim. The task of the agent is to learn a policy that would take the swimmer to the other end of the river.

Population Growth Model: The Population Growth MDP [17] models the population growth of pests and consists of 50 states and 5 actions. The states represent the pest population and actions represent the pest control measures. In our experiments, we use two different instantiations of the Population Growth Model: Population-Small and Population, which vary in the number of posterior samples.

Inventory Management: The Inventory Management MDP [38] models the classical inventory management problem and consists of 30 states and 30 actions. States represent the inventory level and actions represent the inventory to be purchased. The sale price s , holding cost c and purchase costs p are 3.99, 0.03, and 2.219. The demand is normally distributed with mean= $s/4$ and standard deviation $s/6$.

Implementation details: For each domain in our experiments, we sample a dataset \mathcal{D} consisting of n tuples of the form $\{s, a, r, s'\}$, corresponding to the state s , the action taken a , the reward r and

the next state s' . We construct a posterior distribution over the models using \mathcal{D} , assuming Dirichlet priors over the model parameters. Using MCMC sampling, we construct two datasets containing M and K transition models, respectively.

We construct ten train datasets by randomly sampling 80% of the models from the first dataset ten times. For any given confidence level δ , we train one RL agent per train dataset and method.

For evaluation, we consider two instances of the VaR framework: one (denoted by $VaRN$) that assumes that \tilde{P} is a multivariate normal, and another (denoted by VaR) that does not assume any structure over \tilde{P} . We use four baseline methods for evaluating the robustness of our framework. They are: BCR Robust MDPs with weighted ℓ_1 ambiguity sets ($WBCR \ell_1$), weighted ℓ_∞ ambiguity sets ($WBCR \ell_\infty$), unweighted ℓ_1 ambiguity sets ($BCR \ell_1$) and unweighted ℓ_∞ ambiguity sets ($BCR \ell_\infty$). See Appendix C in the appendix for more details.

We compare the mean and standard deviation of the robust performance (δ -percentile of expected returns) of the VaR framework on the test dataset with that of other baselines for different values of δ .

Methods	Riverswim	Inventory	Population-Small	Population
VaR	68.54 ± 2.54	457.95 ± 0.37	-3102.48 ± 214.85	-4578.84 ± 69.76
BCR ℓ_1	67.27 ± 0.0	369.67 ± 0.0	-5614.95 ± 40.14	-5971.81 ± 618.04
BCR ℓ_∞	67.27 ± 0.0	199.41 ± 19.51	-7908.92 ± 20.8	-9020.42 ± 51.11
WBCR ℓ_1	67.9 ± 1.91	454.1 ± 2.08	-5290.38 ± 542.13	-5350.25 ± 68.34
WBCR ℓ_∞	67.27 ± 0.0	199.4 ± 19.51	-7712.43 ± 27.98	-8378.0 ± 63.03
VaRN	67.27 ± 0.0	452.78 ± 0.01	-4005.53 ± 4.38	-4576.65 ± 0.0

Table 1: shows the mean and standard deviation of the robust (percentile) returns at $\delta = 0.05$ achieved by VaR , $VaRN$, $BCR \ell_1$, $BCR \ell_\infty$, $WBCR \ell_1$ and $WBCR$ in Riverswim, Inventory, Population-Small, and Population domain.

Experimental Results Table 1 summarizes the performance of the VaR framework and the baselines for confidence level $\delta = 0.05$ (Table 2 in the appendix summarizes the results for $\delta = 0.3$). We observe that for both confidence levels, $\delta = 0.05$ and $\delta = 0.3$, the VaR framework outperforms the baseline methods in terms of mean robust performance in most domains. We also see that all the baselines and the VaR framework have similar mean robust performance in the Riverswim domain. We conjecture that this is due to the simplicity of the domain. Furthermore, as expected, we find that the robust performance of BCR Robust MDPs with span-optimized (weighted) ambiguity sets ($WBCR \ell_1$, $WBCR \ell_\infty$) is relatively higher than the robust performance of Robust MDPs with unweighted BCR ambiguity sets ($BCR \ell_1$, $BCR \ell_\infty$). However, we find that even Robust MDPs with span-optimized BCR ambiguity sets are unable to outperform the robust performance of our VaR_α framework.

Figure 4 shows how the robust performance of the VaR framework and the baselines varies for different confidence level values ($1 - \delta$). As expected, the robust performance of the VaR is higher for smaller values of confidence level ($1 - \delta$), but more importantly, we observe that in most domains, the robust performance of VaR and $VaRN$ dominates the robust performance of the baselines for a wide range of values of confidence level δ .

Figure 2 compares the robust performance of the VaR framework and the baselines on both train and test models. The trends in the robust performance of the VaR framework and the baselines are similar on both train and test models.

6 Limitations and Conclusion

The main limitation of the VaR framework is that it does not consider the correlations in the uncertainty of transition probabilities across states and actions [13, 21, 22]. However, due to the non-convex nature of the percentile-criterion [3, 29], constructing a tractable VaR_α Bellman operator that considers these correlations is not feasible. One plausible solution is to use a Conditional Value at Risk Bellman operator [9, 19] which is convex and lower bounds the Value at Risk measure. We leave the analysis of this approach for future work. Empirical analysis of the VaR_α framework in domains with continuous state-action spaces is also an interesting avenue for future work. It would be

valuable to rigorously test the VaR_α framework in large RL domains and compare its performance against other Robust RL methods [11, 13, 22].

In conclusion, we propose a novel dynamic programming algorithm that optimizes a tight lower-bound approximation on the percentile criterion without explicitly constructing ambiguity sets. We theoretically show that our algorithm implicitly constructs tight uncertainty sets that are smaller in size than any optimized Bayesian credible region, and thus computes less conservative policies with the same confidence guarantees on returns. We also derive finite-sample and asymptotic bounds on the performance loss due to our approximation. Finally, our experimental results demonstrate the efficacy of our method in several domains.

References

- [1] Praveen Agarwal, Mohamed Jleli, and Bessem Samet. *Banach contraction principle and applications*, pages 1–23. Springer Singapore, Singapore, 2018.
- [2] Kishan Panaganti Badrinath and Dileep Kalathil. Robust reinforcement learning using least squares policy iteration with provable performance guarantees. In Marina Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 511–520. PMLR, 18–24 Jul 2021.
- [3] Bahram Behzadian, Reazul Hasan Russel, Marek Petrik, and Chin Pang Ho. Optimizing percentile criterion using robust MDPs. In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 1009–1017, 2021.
- [4] Bahram Behzadian, Marek Petrik, and Chin Pang Ho. Fast algorithms for L_∞ -constrained s-rectangular robust MDPs. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, 2021.
- [5] Dimitri P. Bertsekas. *Abstract dynamic programming*. Athena Scientific, 2013. ISBN 978-1-886529-42-7.
- [6] Stéphane Boucheron, Gábor Lugosi, and Pascal Massart. *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 02 2013.
- [7] G. P. Box and G. C. Tiao. Bayesian inference in statistical analysis. *International Statistical Review*, 43:242, 1973.
- [8] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [9] Yinlam Chow and Mohammad Ghavamzadeh. Algorithms for CVaR optimization in MDPs. In *International Conference on Neural Information Processing Systems*, pages 3509–3517. MIT Press, 2014.
- [10] E. Delage and S. Mannor. Percentile optimization for Markov decision processes with parameter uncertainty. *Operations Research*, 58(1):203–213, 2010.
- [11] Esther Derman, Daniel Mankowitz, Timothy A Mann, and Shie Mannor. Soft-robust actor-critic policy-gradient. In *Proceedings of the 34th Annual Conference on Uncertainty in Artificial Intelligence (UAI)*, 2018.
- [12] Esther Derman, Daniel J. Mankowitz, Timothy A. Mann, and Shie Mannor. A Bayesian approach to robust reinforcement learning. In *Proceedings of the 35th Annual Conference on Uncertainty in Artificial Intelligence (UAI)*, 2019.
- [13] Vineet Goyal and Julien Grand-Clement. Robust Markov decision process: Beyond rectangularity. *Mathematics of Operations Research*, 48(1):203–226, 2023.
- [14] Julien Grand-Clement and Christian Kroer. First-order methods for Wasserstein distributionally robust MDP. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pages 2010–2019, 2021.
- [15] Vishal Gupta. Near-optimal Bayesian ambiguity sets for distributionally robust optimization. *Management Science*, 56(9), 2019.

- [16] C. A. R. Hoare. Algorithm 65: Find. *Communications of the ACM*, 4(7):321–322, 1961.
- [17] Marc Kery and Michael Schaub. *Bayesian Population Analysis Using WinBUGS*. Elsevier, 2012.
- [18] Sasche Lange, Thomas Gabel, and Martin Riedmiller. *Batch reinforcement learning*. Springer, 2012.
- [19] Elita A. Lobo, Mohammad Ghavamzadeh, and Marek Petrik. Soft-robust algorithms for batch reinforcement learning, 2021.
- [20] Daniel J. Mankowitz, Nir Levine, Rae Ung Jeong, Abbas Abdolmaleki, Jost Tobias Springenberg, Timothy Mann, Todd Hester, and Martin A. Riedmiller. Robust reinforcement learning for continuous control with model misspecification. In *Proceedings of the 2020 International Conference on Learning Representations (ICLR)*, 2020.
- [21] Shie Mannor, Ofir Mebel, and Huan Xu. Lightning does not strike twice: Robust MDPs with coupled uncertainty. In *Proceedings of the 29th International Conference on Machine Learning (ICML)*, pages 451–458, 2012.
- [22] Shie Mannor, Ofir Mebel, and Huan Xu. Robust MDPs with K-rectangular uncertainty. *Mathematics of Operations Research*, 41(4):1484–1509, 2016.
- [23] P. Massart. The tight constant in the Dvoretzky-Kiefer-Wolfowitz inequality. *The Annals of Probability*, 18:1269 – 1283, 1990.
- [24] Ariel Neufeld and Julian Sester. Robust q -learning algorithm for Markov decision processes under wasserstein uncertainty, 2023.
- [25] Arnab Nilim and Laurent El Ghaoui. Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- [26] Martin L. Puterman. *Markov Decision Processes: Discrete stochastic dynamic programming*. John Wiley and Sons, Inc., 2nd edition, 2005.
- [27] R. Tyrrell Rockafellar. *Convex analysis*. Princeton University Press, 1970.
- [28] R. Tyrrell Rockafellar and Stanislav Uryasev. Optimization of conditional value-at-risk. *Journal of Risk*, 2:21–41, 2000.
- [29] Reazul Hasan Russel and Marek Petrik. Beyond confidence regions: Tight Bayesian ambiguity sets for robust MDPs. In *Proceedings of the 32nd Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- [30] Sergey Sarykalin, Gaia Serraino, and Stan Uryasev. *Value-at-risk vs. conditional value-at-risk in risk management and optimization*, chapter 13, pages 270–294. INFORMS, 2008.
- [31] A. Shapiro, W. Tekaya, J. P. da Costa, and M. P. Soares. Risk-neutral and risk-averse stochastic dual dynamic programming method. *European Journal of Operational Research*, pages 375–391, 2013.
- [32] Alexander L. Strehl and Michael L. Littman. A theoretical analysis of model-based interval estimation. In *Proceedings of the 22nd International Conference on Machine Learning, ICML '05*, page 856–863. Association for Computing Machinery, 2005. ISBN 1595931805.
- [33] Richard S. Sutton and Andrew G. Barto. *Reinforcement Learning: An introduction*. A Bradford Book, 2018.
- [34] A. W. van der Vaart. *Asymptotic statistics*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 1998.
- [35] Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov decision processes. *Mathematics of Operations Research*, 38(1):153–183, 2013.

- [36] Huan Xu and Shie Mannor. The robustness-performance trade-off in Markov decision processes. In *Proceedings of the 19th Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2006.
- [37] Huan Xu and Shie Mannor. Distributionally robust Markov decision processes. *Mathematics of Operations Research*, 37(2):288–300, 2012.
- [38] Paul Herbert Zipkin. *Foundations of inventory management*. McGraw-Hill, Boston, 2000.

A Additional theoretical results

Definition A.1 (Subvariance). For any function $v : \mathcal{S} \rightarrow \mathbb{R}$, a scalar c and $\gamma \in [0, 1)$, the operator \mathfrak{T} satisfies the Translation subvariance property if

$$\mathfrak{T}(v + c\mathbf{1}) = (\mathfrak{T}v) + \gamma c\mathbf{1} .$$

Definition A.2 (Monotonicity). For any function $v : \mathcal{S} \rightarrow \mathbb{R}$ and $u : \mathcal{S} \rightarrow \mathbb{R}$, such that $v \preceq u$ the operator \mathfrak{T} satisfies the Monotonicity property if

$$\mathfrak{T}v \preceq \mathfrak{T}u .$$

Lemma A.3 (Contraction Mapping [5]). For any two bounded functions $u : \mathcal{S} \rightarrow \mathbb{R}$, $v : \mathcal{S} \rightarrow \mathbb{R}$, and $\gamma \in [0, 1)$, the operator \mathfrak{T} is a contraction mapping if it satisfies Monotonicity and Translation subvariance properties. In particular, for all u, v there holds

$$\|\mathfrak{T}u - \mathfrak{T}v\|_\infty \leq \gamma \|u - v\|_\infty .$$

Furthermore, for any two bounded functions $u : \mathcal{S} \rightarrow \mathbb{R}$, and $v : \mathcal{S} \rightarrow \mathbb{R}$, the operator \mathfrak{T} is a non-expansive mapping if it satisfies Monotonicity and Translation invariance properties. In particular, for all u, v there holds

$$\|\mathfrak{T}u - \mathfrak{T}v\|_\infty \leq \|u - v\|_\infty .$$

The proof of Lemma A.3 follows directly from Proposition 2.1.3 in [5]. We re-derive the proof for the sake of completeness.

Proof. Denote

$$c = \max_{s \in \mathcal{S}} |u(s) - v(s)| .$$

Therefore for all $s \in \mathcal{S}$,

$$u(s) - c \leq v(s) \leq u(s) + c .$$

Applying \mathfrak{T} to these inequalities and using the Subvariance (Definition A.1) and Monotonicity (Definition A.2) properties, we obtain that for all $s \in \mathcal{S}$,

$$(\mathfrak{T}u)(s) - \gamma c \leq (\mathfrak{T}v)(s) \leq (\mathfrak{T}u)(s) + \gamma c .$$

It follows that for all $s \in \mathcal{S}$,

$$|(\mathfrak{T}v)(s) - (\mathfrak{T}u)(s)| \leq \gamma c ,$$

and therefore $\|\mathfrak{T}u - \mathfrak{T}v\|_\infty \leq \gamma c$, proving the stated result. \square

Proposition A.4 (Validity). Suppose that the reward for any state s and action a , and next state s' is independent of the next state s' , i.e., $R(s, a, s') = R(s, a)$, where $R(s, a) \in \mathbb{R}$. Then, the following properties hold for all value functions $u, v \in \mathbb{R}^{\mathcal{S}}$.

1. The operator \mathcal{T}_{VaR} is contraction mapping on $\mathbb{R}^{\mathcal{S}}$: $\|\mathcal{T}_{\text{VaR}}u - \mathcal{T}_{\text{VaR}}v\|_\infty \leq \gamma \|u - v\|_\infty$.
2. The operator \mathcal{T}_{VaR} is monotone: $u \succeq v \Rightarrow \mathcal{T}_{\text{VaR}}u \succeq \mathcal{T}_{\text{VaR}}v$.
3. The equality $\mathcal{T}_{\text{VaR}}\hat{v} = \hat{v}$ has a unique solution.

Proof. From Lemma A.3, we know that an operator is a contraction mapping if it satisfies Monotonicity and Subvariance property.

In this proof, we will show that the VaR Bellman operator $\mathcal{T}_{\text{VaR}}^\pi$ satisfies Monotonicity and Subvariance property, and therefore, is a contraction mapping.

We will use shorthand $r_{s,a}$ to denote the reward corresponding to state s and action a , i.e., $r_{s,a} = R(s, a)$.

First, we show that \mathcal{T}_{VaR} satisfies Subvariance property. Consider any $c \in \mathbb{R}$ and state s . Then,

$$\begin{aligned}
(\mathcal{T}_{\text{VaR}}(\mathbf{v} + c\mathbf{1}))(s) &= \max_{a \in \mathcal{A}} \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top(r_{s,a} + \gamma(\mathbf{v} + c\mathbf{1}))] \\
&\stackrel{(a)}{=} \max_{a \in \mathcal{A}} \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top(r_{s,a} + \gamma\mathbf{v} + \gamma c\mathbf{1})] \\
&\stackrel{(b)}{=} \max_{a \in \mathcal{A}} \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top(r_{s,a} + \gamma\mathbf{v}) + \gamma c] \\
&\stackrel{(c)}{=} \max_{a \in \mathcal{A}} \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top(r_{s,a} + \gamma\mathbf{v})] + \gamma c \\
&\stackrel{(d)}{=} \mathcal{T}_{\text{VaR}}\mathbf{v} + \gamma c .
\end{aligned}$$

(a) follows from simple algebraic manipulations, (b) follows from $\gamma c \tilde{\mathbf{p}}_{s,a}^\top \mathbf{1} = \gamma c$, (c) follows from the translational invariance property of VaR measure [30], and (d) follows the definition of the VaR Bellman operator $\mathcal{T}_{\text{VaR}}^\pi$.

Next, we show that \mathcal{T}_{VaR} satisfies Monotonicity property.

Let \mathbf{u} and \mathbf{v} be any two value functions such that $\mathbf{v} \preceq \mathbf{u}$. Consider any state $s \in \mathcal{S}$. Then,

$$\begin{aligned}
(\mathcal{T}_{\text{VaR}}\mathbf{v})(s) - (\mathcal{T}_{\text{VaR}}\mathbf{u})(s) &= \max_{a \in \mathcal{A}} \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top(r_{s,a} + \gamma\mathbf{v})] - \max_{a \in \mathcal{A}} \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top(r_{s,a} + \gamma\mathbf{u})] \\
&\stackrel{(a)}{=} \max_{a \in \mathcal{A}} \text{VaR}_\alpha[r_{s,a} + \gamma\tilde{\mathbf{p}}_{s,a}^\top\mathbf{v}] - \max_{a \in \mathcal{A}} \text{VaR}_\alpha[r_{s,a} + \gamma\tilde{\mathbf{p}}_{s,a}^\top\mathbf{u}] \\
&\stackrel{(b)}{=} \max_{a \in \mathcal{A}}(r_{s,a} + \gamma \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{v}]) - \max_{a \in \mathcal{A}}(r_{s,a} + \gamma \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{u}]) \\
&\stackrel{(c)}{\leq} 0 \\
(\mathcal{T}_{\text{VaR}}\mathbf{v})(s) &\leq (\mathcal{T}_{\text{VaR}}\mathbf{u})(s)
\end{aligned}$$

(a) follows from the fact that $r_{s,a}$ is independent of the next state s' , (b) follows from the translational invariance property of the VaR measure [30], (c) follows from the fact that, for any action a , $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{u}] \geq \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{v}]$ because $\mathbf{u} \succeq \mathbf{v}$ and the VaR measure satisfies monotonicity property [30]. Thus, $\forall a \in \mathcal{A}$, $r_{s,a} + \gamma \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{v}] \leq r_{s,a} + \gamma \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{u}]$ and therefore, $\max_{a \in \mathcal{A}}(r_{s,a} + \gamma \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{v}]) \leq \max_{a \in \mathcal{A}}(r_{s,a} + \gamma \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top\mathbf{u}])$.

Thus, we prove that \mathcal{T}_{VaR} is a γ -contraction mapping and a monotone operator. Since \mathcal{T}_{VaR} is a contraction operator on a Banach space, the Banach fixed point theorem [1] implies that the operator \mathcal{T}_{VaR} has a unique solution $\hat{\mathbf{v}}$, i.e., $\mathcal{T}_{\text{VaR}}\hat{\mathbf{v}} = \hat{\mathbf{v}}$. \square

Proposition A.5. Let $\hat{\mathbf{q}} \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ represent the optimal action-value function corresponding to the optimal VaR policy $\hat{\pi}$. For any $k \in \mathbb{Z}$, let \mathbf{f}^k represent the empirical estimate of the k^{th} action-value function \mathbf{q}^k in Algorithm 3.1 computed with atmost ζ error, i.e., $\|\mathbf{f}^k - \mathbf{q}^k\|_\infty \leq \zeta$. If $k \geq \frac{\ln(\frac{R_{\max}}{(\varepsilon - \zeta)(1 - \gamma)})}{1 - \gamma}$ where $R_{\max} = \max_{s \in \mathcal{S}, a \in \mathcal{A}, s' \in \mathcal{S}} r(s, a, s')$ and $\varepsilon > \zeta$, then,

$$\|\hat{\mathbf{q}} - \mathbf{f}^k\|_\infty \leq \varepsilon .$$

Proof. Let $\hat{\mathbf{f}}$ be the empirical fixed point of the empirical VaR Bellman operator $\hat{\mathcal{T}}_{\text{VaR}}$ defined for any value function \mathbf{v} and state s as

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

To prove the result, we assume that the empirical VaR Bellman operator $\hat{\mathcal{T}}_{\text{VaR}}$ is also a contraction mapping like the VaR Bellman operator \mathcal{T}_{VaR} . Then,

$$\begin{aligned}
\|\hat{\mathbf{q}} - \mathbf{f}^k\|_\infty &\stackrel{(a)}{=} \|\hat{\mathbf{q}} - \hat{\mathbf{f}} + \hat{\mathbf{f}} - \mathbf{f}^k\| \\
&\stackrel{(b)}{\leq} \|\hat{\mathbf{q}} - \hat{\mathbf{f}}\| + \|\hat{\mathbf{f}} - \mathbf{f}^k\|_\infty \\
&\stackrel{(c)}{\leq} \zeta + \|\hat{\mathcal{T}}_{\text{VaR}}\hat{\mathbf{f}} - \hat{\mathcal{T}}_{\text{VaR}}\mathbf{f}^{k-1}\| \\
&\stackrel{(d)}{\leq} \zeta + \gamma\|\hat{\mathbf{f}} - \mathbf{f}^{k-1}\| \\
&\stackrel{(e)}{\leq} \gamma\|\hat{\mathcal{T}}_{\text{VaR}}\hat{\mathbf{f}} - \hat{\mathcal{T}}_{\text{VaR}}\mathbf{f}^{k-2}\| + \zeta \\
&\stackrel{(f)}{\leq} \gamma^k \frac{R_{\max}}{1-\gamma} + \zeta .
\end{aligned}$$

(a) follows from simply adding and subtracting $\hat{\mathbf{f}}$, (b) follows from triangle inequality, the first term of (c) follows from the assumption that $\|\mathbf{q}^k - \mathbf{f}^k\|_\infty \leq \zeta$ for any $k \in \mathbb{Z}$ and the second term follows from the properties of VaR $_\alpha$ Bellman operator, (d) follows from the contraction property of $\widehat{\mathcal{T}}_{\text{VaR}}$, (e) follows from applying the same procedure as in (a) to step (d), and (f) follows from unrolling (e) over k time steps and $\|\hat{\mathbf{f}} - \mathbf{f}^0\|_\infty \leq R_{\max}/1-\gamma$ for $\mathbf{f}^0 = \mathbf{0}$.

We find k such that $\|\hat{\mathbf{q}} - \mathbf{f}^k\|_\infty \leq \varepsilon$,

$$\begin{aligned}
\frac{\gamma^k}{1-\gamma}(R_{\max}) + \zeta &\leq \varepsilon \\
k &\geq \frac{\ln\left(\frac{R_{\max}}{(\varepsilon-\zeta)(1-\gamma)}\right)}{1-\gamma} .
\end{aligned}$$

□

Given a policy π , a state $s \in \mathcal{S}$, and a transition model P , let

$$(\mathcal{T}_P^\pi \mathbf{v})(s) = \sum_{a \in \mathcal{A}} \pi(s, a) \mathbf{p}_{s,a}^\top \mathbf{w}_{s,a} \text{ and } (\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})(s) = \sum_{a \in \mathcal{A}} \pi(s, a) \text{VaR}_\alpha [\hat{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] ,$$

represent the Bellman evaluation operator corresponding to transition model P and the VaR Bellman evaluation operator, respectively.

Lemma A.6. *For any policy π , let $\hat{\mathbf{v}}^\pi$ and \mathbf{v}^π be the fixed point of the VaR policy evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ and Bellman policy evaluation operator \mathcal{T}_P^π . If the VaR Bellman policy evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ dominates the Bellman policy evaluation operator \mathcal{T}_P^π at $\hat{\mathbf{v}}^\pi$, i.e., $\mathcal{T}_{\text{VaR}}^\pi \hat{\mathbf{v}}^\pi \preceq \mathcal{T}_P^\pi \hat{\mathbf{v}}^\pi$, then, the fixed point of the VaR Bellman evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ dominates the fixed point of the Bellman evaluation operator \mathcal{T}_P^π , i.e., $\hat{\mathbf{v}}^\pi \preceq \mathbf{v}^\pi$.*

We note that in contrast to the Bellman policy evaluation operator \mathcal{T}_P^π , the VaR Bellman policy evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ is a function of the random variable \hat{P} and is not dependent on the transition model P assumed in this setting.

Using the assumption $\mathcal{T}_{\text{VaR}}^\pi \hat{\mathbf{v}}^\pi \preceq \mathcal{T}_P^\pi \hat{\mathbf{v}}^\pi$, and from $\hat{\mathbf{v}}^\pi = \mathcal{T}_{\text{VaR}}^\pi \hat{\mathbf{v}}^\pi$ and $\mathbf{v}^\pi = \mathcal{T}_P^\pi \mathbf{v}^\pi$, we get by algebraic manipulations:

Proof.

$$\hat{\mathbf{v}}^\pi - \mathbf{v}^\pi = \mathcal{T}_{\text{VaR}}^\pi \hat{\mathbf{v}}^\pi - \mathcal{T}_P^\pi \mathbf{v}^\pi \preceq \mathcal{T}_P^\pi \hat{\mathbf{v}}^\pi - \mathcal{T}_P^\pi \mathbf{v}^\pi \preceq \gamma P_\pi (\hat{\mathbf{v}}^\pi - \mathbf{v}^\pi) .$$

Here P_π is the transition probability function corresponding to policy π . Subtracting $\gamma P_\pi (\hat{\mathbf{v}}^\pi - \mathbf{v}^\pi)$ from the above inequality gives,

$$(I - \gamma P^\pi)(\hat{\mathbf{v}}^\pi - \mathbf{v}^\pi) \preceq \mathbf{0} .$$

where I is the identity matrix. $(I - \gamma P_\pi)^{-1}$ is monotone as can be seen from its Neumann series.

$$\hat{\mathbf{v}}^\pi - \mathbf{v}^\pi \preceq (I - \gamma P_\pi)^{-1} \mathbf{0} = \mathbf{0} .$$

which proves the result. \square

Proposition A.7. *For any policy π , the VaR Bellman evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ dominates the BCR Bellman evaluation operator $\mathcal{T}_{\text{BCR}}^\pi$ i.e., $\mathcal{T}_{\text{VaR}}^\pi \mathbf{v} \succeq \mathcal{T}_{\text{BCR}}^\pi \mathbf{v}$. Similar results hold for policy optimization operators, i.e., $\mathcal{T}_{\text{VaR}} \mathbf{v} \succeq \mathcal{T}_{\text{BCR}} \mathbf{v}$.*

Proof. Recall the VaR Bellman and BCR Bellman operators defined for any value function \mathbf{v} , policy $\pi \in \Pi_D$, confidence level α , and state s as

$$\begin{aligned} (\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})(s) &= \text{VaR}_\alpha[\tilde{\mathbf{p}}^\top \mathbf{w}_{s,\pi(s)}] \\ (\mathcal{T}_{\text{BCR}}^\pi \mathbf{v})(s) &= \min_{\mathbf{p}_{s,\pi(s)} \in \mathcal{P}_{s,\pi(s)}} \tilde{\mathbf{p}}_{s,\pi(s)}^\top \mathbf{w}_{s,\pi(s)} . \end{aligned}$$

Suppose that for any value function \mathbf{v} , confidence level α , state s and action a , $\text{VaR}_\alpha[\tilde{\mathbf{p}}^\top \mathbf{w}_{s,a}] < \min_{\mathbf{p}_{s,a} \in \mathcal{P}_{s,a}^{\text{BCR}}} \tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}$. We know that,

$$\begin{aligned} \Pr(\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} \leq \text{VaR}_\alpha[\tilde{\mathbf{p}}^\top \mathbf{w}_{s,a}]) &> \alpha \\ \implies \Pr(\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} > \text{VaR}_\alpha[\tilde{\mathbf{p}}^\top \mathbf{w}_{s,a}]) &< 1 - \alpha \\ \implies \Pr(\tilde{\mathbf{p}}_{s,a} \in \mathcal{P}_{s,a}^{\text{BCR}}) &< 1 - \alpha , \end{aligned}$$

which is a contradiction since the BCR ambiguity set \mathcal{P}^{BCR} is constructed such that $\Pr[\tilde{P} \in \mathcal{P}^{\text{BCR}}] \geq 1 - \alpha$.

Since, we proved that $\text{VaR}_\alpha[\tilde{\mathbf{p}}^\top \mathbf{w}_{s,a}] \geq \min_{\mathbf{p}_{s,a} \in \mathcal{P}_{s,a}} \tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}$ for any action a , state s and value function \mathbf{v} , it follows that $(\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})(s) \geq (\mathcal{T}_{\text{BCR}}^\pi \mathbf{v})(s)$ and $(\mathcal{T}_{\text{VaR}} \mathbf{v})(s) \geq (\mathcal{T}_{\text{BCR}} \mathbf{v})(s)$ for all value functions \mathbf{v} and state s . \square

B Proofs

B.1 Proof of Proposition 3.1

Proposition 3.1 (Lower Bound Percentile Criterion). *For any $\delta \in (0, 1)$, if we set the confidence level α in the operator $\mathcal{T}_{\text{VaR}}^\pi$ to δ/s , then for every policy $\pi \in \Pi$: $\Pr_{\tilde{P}}[\hat{\mathbf{v}}^\pi \preceq \tilde{\mathbf{v}}^\pi | \mathcal{D}] \geq 1 - \delta$, where \tilde{P} is a realization of the posterior distribution of the transition model P conditioned on observed transitions \mathcal{D} .*

Proof. Let $\alpha = \delta/s$. Recall that for any policy π , state $s \in \mathcal{S}$, transition model \tilde{P} , the Bellman evaluation operator \mathcal{T}^π and the VaR Bellman evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ are defined as

$$(\mathcal{T}^\pi \mathbf{v})(s) = \sum_{a \in \mathcal{A}} \pi(s, a) \tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} \text{ and } (\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})(s) = \sum_{a \in \mathcal{A}} \pi(s, a) \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] ,$$

respectively.

It is also important to note that the Bellman operator \mathcal{T}^π is defined for a random transition model \tilde{P} .

Let $\hat{\mathbf{v}}^\pi$ be the fixed point of $\mathcal{T}_{\text{VaR}}^\pi$ conditioned on observed transitions \mathcal{D} , and let $\tilde{\mathbf{v}}^\pi$ be a random variable that represents the fixed point of \mathcal{T}^π for a given realization \tilde{P} of the posterior distribution of the transition model P given \mathcal{D} . Then, applying Lemma A.6 to $\mathcal{T}_{\text{VaR}}^\pi$ and \mathcal{T}^π , we have, $\hat{\mathbf{v}}^\pi \preceq \mathbf{v}^\pi$ implies

$$\mathcal{T}_{\text{VaR}}^\pi \hat{\mathbf{v}}^\pi \preceq \mathcal{T}^\pi \hat{\mathbf{v}}^\pi$$

That is for each state s ,

$$\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,\pi(s)}^\top \hat{\mathbf{v}}^\pi] \leq \tilde{\mathbf{p}}_{s,\pi(s)}^\top \hat{\mathbf{v}}^\pi . \quad (11)$$

Using the equation (11), we can bound the probability that the VaR value function lower bounds the true value. \square

$$\Pr_{\tilde{\mathcal{P}}}[\hat{\mathbf{v}}^\pi \preceq \tilde{\mathbf{v}}^\pi | \mathcal{D}] = \Pr_{\tilde{\mathcal{P}}}[\forall s \in \mathcal{S} : \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,\pi(s)}^\top \hat{\mathbf{v}}^\pi] \leq \tilde{\mathbf{p}}_{s,\pi(s)}^\top \hat{\mathbf{v}}^\pi | \mathcal{D}] . \quad (12)$$

From the definition of VaR, we know that for any state s and action a ,

$$\Pr_{\tilde{\mathcal{P}}}[\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{v}}^\pi] \leq \tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{v}}^\pi | \mathcal{D}] \geq 1 - \alpha, \quad (13)$$

Therefore, using union bound and (12) in (13), we can write

$$\Pr_{\tilde{\mathcal{P}}}[\hat{\mathbf{v}}^\pi \succ \tilde{\mathbf{v}}^\pi | \mathcal{D}] \leq \sum_{s \in \mathcal{S}} \Pr_{\tilde{\mathcal{P}}}[\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,\pi(s)}^\top \hat{\mathbf{v}}^\pi] > \tilde{\mathbf{p}}_{s,\pi(s)}^\top \hat{\mathbf{v}}^\pi | \mathcal{D}] .$$

Thus,

$$\Pr[\hat{\mathbf{v}}^\pi \succ \tilde{\mathbf{v}}^\pi | \mathcal{D}] = \sum_{s \in \mathcal{S}} \frac{\delta}{S} = S \frac{\delta}{S} = \delta .$$

B.2 Proof of Proposition 3.2

Proposition 3.2. *Suppose that $\tilde{\mathbf{p}}_{s,a}$ for any state s and action a , is a multivariate sub-Gaussian with mean $\bar{\mathbf{p}}_{s,a}$ and covariance factor $\Sigma_{s,a}$, i.e., $\mathbb{E}[\exp(\lambda(\tilde{\mathbf{p}}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \mathbf{w})] \leq \exp(\lambda^2 \mathbf{w}^\top \Sigma_{s,a} \mathbf{w}/2)$, $\forall \lambda \in \mathbb{R}, \forall \mathbf{w} \in \mathbb{R}^S$. Then, for any state $s \in \mathcal{S}$, \mathcal{T}_{VaR} satisfies*

$$(\mathcal{T}_{\text{VaR}} \mathbf{v})(s) \geq \max_{a \in \mathcal{A}} \left(\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \sqrt{2 \ln(1/\alpha)} \sqrt{\mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a}} \right) .$$

Proof.

$$\begin{aligned} (\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})_s &= \max_{a \in \mathcal{A}} \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] \\ &\stackrel{(a)}{=} \max_{a \in \mathcal{A}} \inf \{ t \mid \Pr(\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} \leq t) > \alpha \} \\ &\stackrel{(b)}{=} \max_{a \in \mathcal{A}} \inf \{ t \mid \Pr((\tilde{\mathbf{p}}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \mathbf{w}_{s,a} \leq (t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a})) > \alpha \} \\ &\stackrel{(c)}{=} \max_{a \in \mathcal{A}} \inf \{ t \mid \Pr(\exp((\tilde{\mathbf{p}}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \mathbf{w}_{s,a}) \leq \exp(t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a})) > \alpha \} \\ &\stackrel{(d)}{\geq} \max_{a \in \mathcal{A}} \inf \left\{ t \mid \inf_{\lambda > 0} \frac{\exp(\lambda^2 \mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a}/2)}{\exp(\lambda(t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}))} > \alpha \right\} \\ &\stackrel{(e)}{=} \max_{a \in \mathcal{A}} \inf \left\{ t \mid \exp\left(\frac{-(t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a})^2}{2 \mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a}}\right) > \alpha \right\} \\ &\stackrel{(f)}{=} \max_{a \in \mathcal{A}} \inf \left\{ t \mid (t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a})^2 < -2 \ln(\alpha) \mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a} \right\} \\ &\stackrel{(g)}{=} \max_{a \in \mathcal{A}} \inf \left\{ t \mid (t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}) \in \left(-\sqrt{2 \ln(1/\alpha)} \sqrt{\mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a}}, \sqrt{2 \ln(1/\alpha)} \sqrt{\mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a}} \right) \right\} \\ &\stackrel{(h)}{=} \max_{a \in \mathcal{A}} \left(\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \sqrt{2 \ln(1/\alpha)} \sqrt{\mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a}} \right) . \end{aligned}$$

Equality (a) follows from the definition of VaR, (b) follows from subtracting $\bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}$ on both sides, (c) follows from taking exponential on both sides, (d) follows from using the upper-bound given by Chernoff bound for a sub-Gaussian distribution [6] i.e.,

$\Pr(\exp((\tilde{\mathbf{p}}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \mathbf{w}_{s,a}) \leq \exp(t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a})) \leq \inf_{\lambda > 0} \frac{\exp(\lambda^2 \mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a} / 2)}{\exp(\lambda(t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}))}$, (e) follows from solving for λ and getting $\lambda = \frac{(t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a})}{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}$, (f) follows from taking \ln on both sides, (g) follows from simple algebraic manipulations and (h) follows from taking the infimum of the solution interval of t .

Solving for t in step (g), we get $t = \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] = \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \sqrt{2 \ln(1/\alpha)} \sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}$ which proves the stated result. \square

B.3 Proof of Proposition 3.3

Proposition 3.3. *Suppose that \tilde{P} is normally distributed, i.e., for any state s and action a , $\tilde{\mathbf{p}}_{s,a} \sim \mathcal{N}(\bar{\mathbf{p}}_{s,a}, \boldsymbol{\Sigma}_{s,a})$.*

Then, \mathcal{T}_{VaR} for any state $s \in \mathcal{S}$ takes the form

$$(\mathcal{T}_{\text{VaR}} \mathbf{v})(s) = \max_{a \in \mathcal{A}} \left(\bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}} \right) .$$

Proof. Consider the VaR Bellman optimality operator defined for any state s and value function \mathbf{v} as

$$(\mathcal{T}_{\text{VaR}} \mathbf{v})(s) = \max_{a \in \mathcal{A}} \text{VaR}_\alpha [\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] . \quad (14)$$

From the theory of multivariate normal distributions [6], we know that, for any state s and action a , since $\tilde{\mathbf{p}}_{s,a}$ is normally distributed $\mathcal{N}(\bar{\mathbf{p}}_{s,a}, \boldsymbol{\Sigma}_{s,a})$, $\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}$ is also normally distributed $\mathcal{N}(\bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}, \mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a})$. To find the $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ for any state s and action a , it is sufficient to find t such that $\Pr(\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} \geq t) = 1 - \alpha$.

$$\begin{aligned} \Pr \left(\frac{(\tilde{\mathbf{p}} - \bar{\mathbf{p}}_{s,a})^\top \mathbf{w}_{s,a}}{\sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}} > \frac{t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}}{\sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}} \right) &= 1 - \alpha \\ 1 - \Phi \left(\frac{t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}}{\sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}} \right) &= 1 - \alpha \\ \left(\frac{t - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}}{\sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}} \right) &= \Phi^{-1}(\alpha) \\ t &= \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} + \Phi^{-1}(\alpha) \sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}} \\ t &= \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}} , \end{aligned}$$

The first equation follows from subtracting $\bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}$ and dividing by $\sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}$ on both sides. The second equality follows from the definition of CDF of $\mathcal{N}(0, 1)$ and the third equality follows from simple algebraic manipulations.

Substituting the value of $t = \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] = \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top \boldsymbol{\Sigma}_{s,a} \mathbf{w}_{s,a}}$ in (14), we obtain the stated results. \square

Proposition 3.3 is useful when the Bernstein von Mises theorem holds for transition probability \tilde{P} , i.e., as the number of posterior samples N goes to ∞ , the posterior distribution of $\tilde{\mathbf{p}}_{s,a}$ for any state s and action a converges to a multivariate normal centered at the maximum likelihood true transition model $\bar{\mathbf{p}}_{s,a}$ with covariance matrix given by the Fisher information matrix $I(\mathbf{p}_{s,a}^*)^{-1}/N$ of the true transition probabilities $\mathbf{p}_{s,a}^*$ [34].

B.4 Proof of Theorem 3.4

Theorem 3.4 (Performance). *Let \hat{v} be the fixed point of the VaR Bellman optimality operator \mathcal{T}_{VaR} , and π^* be the optimal policy in (1). Let $\rho^* = \text{VaR}_\alpha [\rho(\pi^*, \tilde{P})]$ denote the optimal percentile returns and $\hat{\rho} = \mathbf{p}_0^\top \hat{v}$ denote the lower bound on the percentile returns computed using the Bellman operator \mathcal{T}_{VaR} . For any $\delta \in (0, 1)$, we set the confidence level $\alpha = \delta/(2SA)$ in \mathcal{T}_{VaR} . Then, with probability at least $1 - \delta$, the performance loss with respect to ρ^* is*

$$\rho^* - \hat{\rho} \leq \frac{1}{1 - \gamma} \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} (\text{VaR}_{1-\alpha} [\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}] - \text{VaR}_\alpha [\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}]) . \quad (6)$$

Proof. We prove this theorem in three parts.

Let $\hat{v} \in \mathbb{R}^S$ be the fixed point of the VaR_α Bellman operator \mathcal{T}_{VaR} , i.e., $\hat{v} = \mathcal{T}_{\text{VaR}} \hat{v}$. Recall that the VaR_α returns is given by $\hat{\rho} = \mathbf{p}_0^\top \hat{v}$. Furthermore, let \mathcal{T}_P^π represent the Bellman evaluation operator for a given policy $\pi \in \Pi$ and a transition probability model P . Then, \mathcal{T}_P^π is defined for each $s \in \mathcal{S}$ as

$$(\mathcal{T}_P^\pi \mathbf{v})_s = \mathbf{p}_{s,\pi(s)}^\top \mathbf{w}_{s,\pi(s)} ,$$

where $\mathbf{w}_{s,\pi(s)} = \mathbf{r}_{s,\pi(s)} + \gamma \cdot \mathbf{v}$.

It is well known that the Bellman operator \mathcal{T}_P^π is a contraction mapping, monotone, and has a unique fixed point. Let $\tilde{\pi} \in \arg \max_{\pi \in \Pi} \rho(\pi, \tilde{P})$. Let \tilde{v} be the unique fixed point of $\mathcal{T}_{\tilde{P}}^{\tilde{\pi}}$:

$$\tilde{v} = \mathcal{T}_{\tilde{P}}^{\tilde{\pi}} \tilde{v} .$$

Note that it is well-known that:

$$\mathbf{p}_0^\top \tilde{v} = \rho(\tilde{\pi}, \tilde{P}) .$$

First, we show that the lower bound on the percentile returns computed by the VaR Bellman operator \mathcal{T}_{VaR} , i.e., $\mathbf{p}_0^\top \hat{v}$ is less than than the returns corresponding to policy $\tilde{\pi}$, i.e., $\mathbf{p}_0^\top \tilde{v}$ with high confidence $1 - \delta$.

We can write $\hat{\rho} = \mathbf{p}_0^\top \hat{v}$ as

$$\begin{aligned} \mathbf{p}_0^\top \hat{v} &\stackrel{(a)}{\leq} \text{VaR}_\alpha [\rho(\hat{\pi}, \tilde{P})] \\ &\stackrel{(b)}{\leq} \rho(\hat{\pi}, \tilde{P}) \\ &\stackrel{(c)}{\leq} \rho(\tilde{\pi}, \tilde{P}) = \mathbf{p}_0^\top \tilde{v} . \end{aligned}$$

Inequality (a) follows because the VaR Bellman optimality operator optimizes a lower bound on the percentile criterion with high confidence $1 - \delta$. Inequality (b) follows from the definition of VaR. Inequality (c) follows because $\tilde{\pi}$ is optimal for \tilde{P} . The above equation implies

$$0 \leq \mathbf{p}_0^\top (\hat{v} - \tilde{v}) \leq \|\hat{v} - \tilde{v}\|_\infty$$

Now, we establish the probabilistic bound based on bounding the Bellman residual. We will use the following result to establish an upper bound on $\|\hat{v} - \tilde{v}\|_\infty$.

$$\begin{aligned} (\mathcal{T}_{\tilde{P}}^{\tilde{\pi}} \hat{v} - \hat{v})_s &\stackrel{(a)}{=} (\mathcal{T}_{\tilde{P}}^{\tilde{\pi}} \hat{v} - \mathcal{T}_{\text{VaR}} \hat{v})_s \\ &\stackrel{(\text{def})}{=} \tilde{\mathbf{p}}_{s,\tilde{\pi}(s)}^\top \hat{\mathbf{z}}_{s,\tilde{\pi}(s)} - \text{VaR}_\alpha \left[\tilde{\mathbf{p}}_{s,\tilde{\pi}(s)}^\top \hat{\mathbf{z}}_{s,\tilde{\pi}(s)} \right] \\ &\stackrel{(b)}{\leq} \tilde{\mathbf{p}}_{s,\tilde{\pi}(s)}^\top \hat{\mathbf{z}}_{s,\tilde{\pi}(s)} - \text{VaR}_\alpha \left[\tilde{\mathbf{p}}_{s,\tilde{\pi}(s)}^\top \hat{\mathbf{z}}_{s,\tilde{\pi}(s)} \right] \\ &\stackrel{(c)}{\leq} \max_{a \in \mathcal{A}} (\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a} - \text{VaR}_\alpha [\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a}]) \\ &\stackrel{(d)}{\leq} \max_{a \in \mathcal{A}} (\text{VaR}_{1-\alpha} [\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a}] - \text{VaR}_\alpha [\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a}]) \quad \text{with probability } 1 - 2\alpha A . \end{aligned} \quad (15)$$

(a) follows from $\hat{\mathbf{v}}$ being the fixed point of \mathcal{T}_{VaR} , (b) follows from the optimality of $\hat{\pi}$, (c) follows from simple algebraic manipulations, and (d) follows from $\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a} \geq \text{VaR}_{1-\alpha}[\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a}]$ with probability α , $\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a} \leq \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{z}}_{s,a}]$ with probability α , and finally, taking a union bound over the set of actions \mathcal{A} yields the given probabilistic bound.

For any state s , let $\hat{c}_s = \max_{a \in \mathcal{A}} (\text{VaR}_{1-\alpha}[\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}] - \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}])$. Applying the inequality in (15) to all states, we get

$$(\mathcal{T}_{\tilde{\mathbf{P}}} \tilde{\mathbf{v}} - \hat{\mathbf{v}}) \preceq \hat{\mathbf{c}} .$$

We can now use the standard dynamic programming bounding technique to bound $\|\tilde{\mathbf{v}} - \hat{\mathbf{v}}\|_\infty$ as follows:

$$\begin{aligned} 0 &\stackrel{(a)}{\preceq} \tilde{\mathbf{v}} - \hat{\mathbf{v}} \\ &\stackrel{(b)}{\preceq} \tilde{\mathbf{v}} - \mathcal{T}_{\tilde{\mathbf{P}}} \hat{\mathbf{v}} + \mathcal{T}_{\tilde{\mathbf{P}}} \hat{\mathbf{v}} - \hat{\mathbf{v}} \\ &\stackrel{(c)}{\preceq} \tilde{\mathbf{v}} - \mathcal{T}_{\tilde{\mathbf{P}}} \hat{\mathbf{v}} + \hat{\mathbf{c}} \\ \tilde{\mathbf{v}} - \hat{\mathbf{v}} &\preceq \mathcal{T}_{\tilde{\mathbf{P}}} \tilde{\mathbf{v}} - \mathcal{T}_{\tilde{\mathbf{P}}} \hat{\mathbf{v}} + \hat{\mathbf{c}} . \end{aligned}$$

We have (a) because $\hat{\mathbf{v}} \preceq \tilde{\mathbf{v}}$ because $\mathcal{T}_{\text{VaR}} \tilde{\mathbf{v}} \preceq \tilde{\mathbf{v}}$ and thus, $\tilde{\mathbf{v}} \succeq \mathcal{T}_{\text{VaR}} \mathcal{T}_{\text{VaR}} \tilde{\mathbf{v}} \succeq \dots \succeq \mathcal{T}_{\text{VaR}} \dots \mathcal{T}_{\text{VaR}} \tilde{\mathbf{v}} \succeq \hat{\mathbf{v}}$ because $\hat{\mathbf{v}}$ is the fixed point of \mathcal{T}_{VaR} and \mathcal{T}_{VaR} is monotone. We have (b) from simply adding and subtracting $\mathcal{T}_{\tilde{\mathbf{P}}} \hat{\mathbf{v}}$ and (c) follows from (15).

Applying the ℓ_∞ norm on both sides (i.e., taking the max over all states $s \in \mathcal{S}$), we get

$$\begin{aligned} \|\tilde{\mathbf{v}} - \hat{\mathbf{v}}\|_\infty &\leq \|\mathcal{T}_{\tilde{\mathbf{P}}} \tilde{\mathbf{v}} - \mathcal{T}_{\tilde{\mathbf{P}}} \hat{\mathbf{v}} + \hat{\mathbf{c}}\|_\infty \\ &\stackrel{(a)}{\leq} \gamma \|\tilde{\mathbf{v}} - \hat{\mathbf{v}}\|_\infty + \|\hat{\mathbf{c}}\|_\infty \\ \|\tilde{\mathbf{v}} - \hat{\mathbf{v}}\|_\infty &\stackrel{(b)}{\leq} \frac{\max_{s \in \mathcal{S}} \hat{c}_s}{1 - \gamma} \quad \text{with probability } (1 - 2\alpha SA) . \end{aligned}$$

Inequality (a) follows by the triangle inequality, and inequality (b) follows from applying the results in equation (15) and taking a union bound over all states in \mathcal{S} .

Thus, setting $\alpha = \delta/2SA$, we get $\Pr \left[\|\tilde{\mathbf{v}} - \hat{\mathbf{v}}\|_\infty \leq \frac{\max_{s \in \mathcal{S}} \hat{c}_s}{1 - \gamma} \right] \leq 1 - \delta$.

Finally, to prove the bound on ρ^* and $\hat{\rho} = \mathbf{p}_0^\top \hat{\mathbf{v}}$, we need to show $\rho^* \leq \eta$, where $\eta = \hat{\rho} + \frac{\max_{s \in \mathcal{S}} \hat{c}_s}{1 - \gamma}$. Suppose that the contradiction $\rho^* > \eta$ holds true. Realize that ρ^* is optimal in (1), and therefore, must satisfy

$$\Pr \left[\max_{\pi \in \Pi} \rho(\pi, \tilde{P}) \geq \rho^* \right] \geq 1 - \delta .$$

Recall from the statement of the theorem that

$$\Pr \left[\max_{\pi \in \Pi} \rho(\pi, \tilde{P}) \geq \eta \right] \leq \delta .$$

We can now derive a contradiction.

$$\delta \geq \Pr \left[\max_{\pi \in \Pi} \rho(\pi, \tilde{P}) \geq \eta \right] \stackrel{(a)}{\geq} \Pr \left[\max_{\pi \in \Pi} \rho(\pi, \tilde{P}) \geq \rho^* \right] \geq 1 - \delta .$$

where inequality (a) follows from the assumption $\rho^* \geq \eta$. Thus, we get $\delta \geq 1 - \delta$, which is a contradiction for $\delta < 0.5$. The lower bound $0 \leq \rho^* - \hat{\rho}$ follows from the optimality of ρ^* and Proposition 3.1, which proves the proposition. \square

B.5 Proof of Theorem 3.5

Theorem 3.5 (Asymptotic Performance). *For any $\delta \in (0, 1)$, set $\alpha = \delta/(2SA)$ in \mathcal{T}_{VaR} . Let $I(\mathbf{p}_{s,a}^*)^{-1}$ for any state s and action a , be the Fisher Information matrix corresponding to the true transition*

probabilities $\mathbf{p}_{s,a}^*$. Furthermore, let $\sigma_{\max} = \max_{s \in \mathcal{S}, a \in \mathcal{A}} \sqrt{\hat{\mathbf{w}}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \hat{\mathbf{w}}_{s,a}}$ represent the maximum asymptotic standard deviation of the returns estimate $\tilde{\mathbf{p}}_{s,a}^\top \hat{\mathbf{w}}_{s,a}$ for any state-action pair (s, a) . Then, with probability at least $1 - \delta$, the asymptotic performance of the VaR framework $\hat{\rho}$ w.r.t. the optimal percentile returns ρ^* satisfies

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

Proof. To prove this theorem, we assume that Bernstein von Mises theorem [34] holds for \tilde{P} , i.e., the posterior distribution of transition probabilities $\tilde{\mathbf{p}}_{s,a}$ for any state s and action a converges in the limit to a multivariate normal $\mathcal{N}(\tilde{\mathbf{p}}_{s,a}, I(\mathbf{p}_{s,a}^*))$ centered at the maximum likelihood estimator of the true transition model $\mathbf{p}_{s,a}^*$ with covariance matrix given by the Fisher information matrix $I(\mathbf{p}_{s,a}^*)/N$ of the true transition probabilities $\mathbf{p}_{s,a}^*$.

Therefore, for any state s and action a , $\lim_{N \rightarrow \infty} \sqrt{N}(\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}) \rightsquigarrow \mathcal{N}(0, \mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*) \mathbf{w}_{s,a})$.

We know from Proposition 3.3, that Value at risk of a univariate normal random variable $X \sim \mathcal{N}(\mu, \sigma)$ can be written as $\text{VaR}_\alpha[X] = \mu + \Phi^{-1}(\alpha)\sigma$. Therefore, applying this result to the R.H.S of Equation (6) gives

$$\begin{aligned} \lim_{N \rightarrow \infty} \sqrt{N}(\rho^* - \hat{\rho}) &\leq \lim_{N \rightarrow \infty} \frac{1}{1 - \gamma} \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \left(\sqrt{N} \tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} + \Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right. \\ &\quad \left. - \left(\sqrt{N} \tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} + \Phi^{-1}(\alpha) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right) \right) \\ &= \frac{1}{1 - \gamma} \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \left(\Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right. \\ &\quad \left. - \Phi^{-1}(\alpha) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right) \\ &= \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \left(\Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right. \\ &\quad \left. - \Phi^{-1}(\alpha) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right) \\ &= \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \left((\Phi^{-1}(1 - \alpha) - \Phi^{-1}(\alpha)) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right) \\ &= \max_{s \in \mathcal{S}} \max_{a \in \mathcal{A}} \left(2\Phi^{-1}(1 - \alpha) \sqrt{\mathbf{w}_{s,a}^\top I(\mathbf{p}_{s,a}^*)^{-1} \mathbf{w}_{s,a}} \right) . \end{aligned} \tag{16}$$

We prove the second inequality in Theorem 3.5, by leveraging the sub-Gaussian bounds for a standard normal distribution $\mathcal{N}(0, \sigma)$ to show that $\Phi^{-1}(1 - \alpha) \leq \sqrt{2 \ln(1/\alpha)}$.

We know that for a standard normal distribution $X \sim \mathcal{N}(0, \sigma)$ where $\sigma \in \mathbb{R}$, the follow sub-Gaussian bounds holds true [6].

$$\Pr \left(-\sqrt{2 \ln(2/\alpha)} \sigma \leq X \leq \sqrt{2 \ln(2/\alpha)} \sigma \right) \geq 1 - \alpha . \tag{17}$$

It is well known that for a standard normal distribution $\mathcal{N}(0, \sigma)$, the following equation holds.

$$\Pr \left(-\Phi^{-1}(1 - \alpha/2) \sigma \leq X \leq \Phi^{-1}(1 - \alpha/2) \sigma \right) = 1 - \alpha \tag{18}$$

Comparing equation (17) and (18), we get

$$\begin{aligned} \Phi^{-1}(1 - \alpha/2) &\leq \sqrt{2 \ln(2/\alpha)} \\ &\equiv \Phi^{-1}(1 - \alpha) \leq \sqrt{2 \ln(1/\alpha)} . \end{aligned} \tag{19}$$

Substituting equation (19) in equation (16), proves the second inequality of the theorem. \square

B.6 Proof of Proposition 3.6

Proposition 3.6 (Time Complexity). *The time complexity of a single iteration of the loop in line 3 of the VaR Value Iteration (Algorithm 3.1) is in $\mathcal{O}(SAN)$, where N is the number of samples of the posterior samples of \tilde{P} .*

Proof. The proposition follows from the fact that any quantile of an array of real values can be computed in linear time using the Quick Select algorithm [16] and a single iteration of the loop in line 3 of Algorithm 3.1 computes quantile of returns SA times. \square

B.7 Proof of Proposition 3.7

Proposition 3.7 (Empirical Error Bound). *For any state s , action a and value function v , let $\widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ represent the empirical estimate of α -percentile of returns $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ and Φ_f represent the cumulative density function (CDF) of the random estimate of returns $\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}$. Suppose that Φ_f is differentiable at the point $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ and let $m = \Phi'_f(\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}])$ represents the density of estimate of returns at point $\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$. Let N^* be the number of posterior samples required to obtain empirical error $\varepsilon \in \mathbb{R}$, with confidence $1 - \zeta$, i.e., $\Pr \left[\left| \widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] - \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] \right| > \varepsilon \right] \leq \zeta$. Then, $\lim_{\varepsilon \rightarrow 0} N^* \varepsilon^2 = \ln(2/\zeta)/2m^2$.*

Proof. To prove this theorem, we first compute the derivative of the inverse of the CDF $\partial \Phi_f^{-1} / \partial \alpha$ as follows. From the definition of the cdf Φ_f and VaR, we know that, for any $\alpha \in (0, 0.5)$, $\Phi_f^{-1}(\alpha) = \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$.

From the inverse-function theorem, we get,

$$\begin{aligned} (\Phi_f^{-1}(\alpha))' &= \frac{1}{\Phi_f'(\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}])} \\ &= \frac{1}{\Phi_f'(\Phi_f^{-1}(\alpha))} \\ &= \frac{1}{m}. \end{aligned}$$

Equipped with the above result, we can now proceed to prove the main result.

To prove the result, we need to find N^* such that

$$\begin{aligned} \Pr \left[\text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] - \varepsilon \leq \widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] \leq \text{VaR}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] + \varepsilon \right] &\geq 1 - \zeta \\ \stackrel{(a)}{=} \Pr \left[\Phi_f^{-1}(\alpha - \varepsilon m) \leq \widehat{\text{VaR}}_\alpha[\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] \leq \Phi_f^{-1}(\alpha + \varepsilon m) \right] &\geq 1 - \zeta. \end{aligned} \quad (20)$$

Equation (a) follows from applying a first order Taylor expansion to Φ_f^{-1} around the point α . We apply the following results to obtain a bound on N^* .

Let \hat{F} and F represent the empirical CDF and the true CDF of a random variable \tilde{Z} . Suppose that the empirical estimate of the CDF \hat{F} is estimated using N^* samples from the true distribution of \tilde{Z} and $0 < \zeta < 1$ represents the desired level of confidence guarantees, Then, from DWK inequality [23], we know that

$$\Pr \left(\|\hat{F} - F\|_\infty \geq \sqrt{\ln(2/\zeta)/2N^*} \right) \leq \zeta$$

The above equation implies that

$$\Pr \left(\exists p \in (0, 1) : F^{-1}(p) < \hat{F}^{-1}(p - l_t) \text{ or } F^{-1}(p) > \hat{F}^{-1}(p + u_t) \right) \leq \zeta \quad (21)$$

where $l_t = u_t = \sqrt{\ln(2/\zeta)/2N^*}$.

Thus, applying equation (21) to (20), i.e., $l_t = \sqrt{\ln(2/\zeta)/2N^*} = \varepsilon m$ gives $N^* = \ln(2/\zeta)/2\varepsilon^2 m^2$.

\square

B.8 Proof of Proposition 4.1

Proposition 4.1 (Equivalence). *Let \hat{v}^π be the fixed point of the VaR Bellman evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ for each $\pi \in \Pi^D$, i.e. $\hat{v}^\pi = (\mathcal{T}_{\text{VaR}}^\pi \hat{v}^\pi)$, where Π^D is the set of all deterministic policies. Then, the optimal VaR policy $\hat{\pi}$ solves*

$$\max_{\pi \in \Pi^D} \min_{P \in \mathcal{P}^{\text{VaR}, \hat{v}^\pi}} \rho(\pi, P) . \quad (8)$$

Proof. Consider the VaR Bellman optimality operator.

$$\forall s \in \mathcal{S}, \quad \mathbf{v} \in \mathbb{R}^S, \quad (\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})_s = \max_{a \in \mathcal{A}} \text{VaR}_\alpha [\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}] . \quad (22)$$

Notice that $\text{VaR}_\alpha [\tilde{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a}]$ is convex in $\tilde{\mathbf{p}}_{s,a}$. Thus, using the definition of support functions of convex functions [8], we can write Equation (22) as

$$\forall s \in \mathcal{S}, \quad \mathbf{v} \in \mathbb{R}^S, \quad (\mathcal{T}_{\text{VaR}}^\pi \mathbf{v})_s = \max_{a \in \mathcal{A}} \min_{P \in \mathcal{P}_{s,a}^{\text{VaR}, \mathbf{v}}} \mathbf{p}_{s,a}^\top \mathbf{w}_{s,a} , \quad (23)$$

$$\text{where } \mathcal{P}_{s,a}^{\text{VaR}, \mathbf{v}} = \{ \mathbf{p}_{s,a} \in \Delta^S \mid \mathbf{p}_{s,a}^\top \mathbf{w}_{s,a} \geq \text{VaR}_\alpha [\mathbf{p}_{s,a}^\top \mathbf{w}_{s,a}] \} .$$

Equation (23) represents the VaR Bellman operator $\mathcal{T}_{\text{VaR}}^\pi$ (policy evaluation) as a Robust Bellman operator with a special ambiguity set that is dependent on the value function \mathbf{v} . Thus, from Equation 23 and theory of value iteration in RL, we get that the fixed point of VaR Bellman optimality operator $\mathcal{T}_{\text{VaR}}^\pi$ (policy optimization) is the value function of the optimal policy that solves the following optimization problem.

$$\begin{aligned} & \arg \max_{\pi \in \Pi^D} \min_{P \in \mathcal{P}^{\text{VaR}, \mathbf{v}^\pi}} \rho(\pi, P) \\ \mathcal{P}^{\text{VaR}, \mathbf{v}^\pi} = & \times \mathcal{P}_{s,a}^{\text{VaR}, \mathbf{v}^\pi}, \quad \mathcal{P}_{s,a}^{\text{VaR}, \mathbf{v}^\pi} = \{ \mathbf{p}_{s,a} \in \Delta^S \mid \mathbf{p}_{s,a}^\top \mathbf{v}^\pi \geq \text{VaR}_\alpha [\tilde{\mathbf{p}}_{s,a}^\top \mathbf{v}^\pi] \} . \end{aligned}$$

□

B.9 Proof of Proposition 4.2

Proposition 4.2. *For any policy π , the fixed point of the VaR policy evaluation operator $\mathcal{T}_{\text{VaR}}^\pi$ dominates the fixed point of the Bellman evaluation operator $\mathcal{T}_{\text{BCR}}^\pi$, i.e., $\mathcal{T}_{\text{BCR}}^\pi \cdots \mathcal{T}_{\text{BCR}}^\pi \mathbf{v} \preceq \mathcal{T}_{\text{VaR}}^\pi \cdots \mathcal{T}_{\text{VaR}}^\pi \mathbf{v}$ for any \mathbf{v} . Similar results hold for policy optimization operators \mathcal{T}_{VaR} and \mathcal{T}_{BCR} , i.e., $\mathcal{T}_{\text{BCR}} \cdots \mathcal{T}_{\text{BCR}} \mathbf{v} \preceq \mathcal{T}_{\text{VaR}} \cdots \mathcal{T}_{\text{VaR}} \mathbf{v}$ for any \mathbf{v} .*

We prove the first part of the proposition using induction. First, we verify that $(\mathcal{T}_{\text{VaR}}^\pi \mathbf{v}) \succeq (\mathcal{T}_{\text{BCR}}^\pi \mathbf{v})$ for any \mathbf{v} .

Proof.

$$\mathcal{T}_{\text{VaR}}^\pi \mathbf{v} \stackrel{(a)}{\succeq} \mathcal{T}_{\text{BCR}}^\pi \mathbf{v} \quad \forall \mathbf{v} \in \mathbb{R}^S$$

(a) follows from the definition of $\mathcal{T}_{\text{VaR}}^\pi$ and $\mathcal{T}_{\text{BCR}}^\pi$ since the constraints of BCR ambiguity sets encompass the constraints of the ambiguity sets implicitly constructed by $\mathcal{T}_{\text{VaR}}^\pi$ (Proposition A.7).

Suppose that for any $\mathbf{v} \in \mathbb{R}^S$ and $k \in \mathcal{Z}$, applying $\mathcal{T}_{\text{VaR}}^\pi$ and $\mathcal{T}_{\text{BCR}}^\pi$ k times to \mathbf{v} yields, $(\mathcal{T}_{\text{VaR}}^\pi)^k \mathbf{v} \succeq (\mathcal{T}_{\text{BCR}}^\pi)^k \mathbf{v}$. Then,

$$\begin{aligned} (\mathcal{T}_{\text{VaR}}^\pi)^{k+1} \mathbf{v} & \stackrel{(a)}{\succeq} \mathcal{T}_{\text{VaR}}^\pi (\mathcal{T}_{\text{BCR}}^\pi)^k \mathbf{v} \\ (\mathcal{T}_{\text{VaR}}^\pi)^{k+1} \mathbf{v} & \stackrel{(b)}{\succeq} \mathcal{T}_{\text{BCR}}^\pi (\mathcal{T}_{\text{BCR}}^\pi)^k \mathbf{v} \\ (\mathcal{T}_{\text{VaR}}^\pi)^{k+1} \mathbf{v} & \succeq (\mathcal{T}_{\text{BCR}}^\pi)^{k+1} \mathbf{v} . \end{aligned}$$

(a) follows from the monotonicity property of $\mathcal{T}_{\text{VaR}}^\pi$, (b) follows from the fact that $\mathcal{T}_{\text{VaR}}^\pi \mathbf{w} \succeq \mathcal{T}_{\text{BCR}}^\pi \mathbf{w}$ for any \mathbf{w} which in turn, follows from the definition of $\mathcal{T}_{\text{VaR}}^\pi$ and $\mathcal{T}_{\text{BCR}}^\pi$. Therefore, $\mathcal{T}_{\text{VaR}}^\pi (\mathcal{T}_{\text{BCR}}^\pi)^k \mathbf{v} \succeq \mathcal{T}_{\text{BCR}}^\pi (\mathcal{T}_{\text{BCR}}^\pi)^k \mathbf{v}$. Thus, by induction, it follows that $\mathcal{T}_{\text{VaR}}^\pi \cdots \mathcal{T}_{\text{VaR}}^\pi \mathbf{v} \succeq \mathcal{T}_{\text{BCR}}^\pi \cdots \mathcal{T}_{\text{BCR}}^\pi \mathbf{v}$ for any \mathbf{v} .

Now, we prove the second part of the proposition. Let $\hat{\mathbf{v}}$ and \mathbf{v}^{BCR} be the fixed points of \mathcal{T}_{VaR} and \mathcal{T}_{BCR} respectively. Supposed that $\mathbf{v}^{\text{BCR}} \succ \hat{\mathbf{v}}$, then

$$\begin{aligned} \mathcal{T}_{\text{VaR}} \mathbf{v}^{\text{BCR}} &\stackrel{\text{(a)}}{\succ} \mathcal{T}_{\text{VaR}} \hat{\mathbf{v}} \\ \mathcal{T}_{\text{VaR}} \mathcal{T}_{\text{VaR}} \mathbf{v}^{\text{BCR}} &\stackrel{\text{(b)}}{\succ} \mathcal{T}_{\text{VaR}} \mathcal{T}_{\text{VaR}} \hat{\mathbf{v}} \\ \mathcal{T}_{\text{VaR}} \dots \mathcal{T}_{\text{VaR}} \mathcal{T}_{\text{VaR}} \mathbf{v}^{\text{BCR}} &\stackrel{\text{(c)}}{\succ} \mathcal{T}_{\text{VaR}} \dots \mathcal{T}_{\text{VaR}} \mathcal{T}_{\text{VaR}} \hat{\mathbf{v}} \\ \hat{\mathbf{v}} &\stackrel{\text{(d)}}{\succ} \hat{\mathbf{v}} . \end{aligned}$$

(a), (b), (c) follows from the monotonicity property of \mathcal{T}_{VaR} and (d) follows from the fact that $\hat{\mathbf{v}}$ is the unique fixed point of \mathcal{T}_{VaR} . The last equation is a contradiction. therefore, it must be that $\mathbf{v}^{\text{BCR}} \preceq \hat{\mathbf{v}}$. \square

B.10 Proof of Theorem 4.3

Theorem 4.3 (Asymptotic Radii of VaR Ambiguity Sets). *Let $\bar{P} = (\bar{\mathbf{p}}_{s,a})_{s \in \mathcal{S}, a \in \mathcal{A}}$ be the Maximum Likelihood estimate of transition probabilities computed from data \mathcal{D} and $\Sigma = (I(\mathbf{p}_{s,a}^*)^{-1})_{s \in \mathcal{S}, a \in \mathcal{A}}$ be the corresponding covariance matrix. Then, $\forall s \in \mathcal{S}, a \in \mathcal{A}$,*

$$\lim_{N \rightarrow \infty} \sqrt{N}(\mathcal{P}_{s,a}^{\text{VaR}} - \bar{\mathbf{p}}_{s,a}) = \left\{ \mathbf{p}_{s,a} \in \Delta^S \mid \|\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a}\|_{\Sigma_{s,a}^{-1}} \leq \Phi^{-1}(1 - \alpha) \right\} - \bar{\mathbf{p}}_{s,a} . \quad (9)$$

To prove this theorem, we first establish some primary results.

We note that to prove this theorem, we assume that Bernstein von Mises theorem holds for transition probabilities \tilde{P} .

For any state s , action a and value function \mathbf{v} , consider the 1-step VaR $_{\alpha}$ Bellman update

$$f_{s,a}(\mathbf{v}) = \text{VaR}_{\alpha} \left[\tilde{\mathbf{p}}_{s,a}^{\top} (\mathbf{r}_{s,a} + \gamma \mathbf{v}) \right] , \quad (24)$$

Furthermore, let $\mathbf{w}_{s,a} = \mathbf{r}_{s,a} + \gamma \mathbf{v}$. If the posterior of transition probability $\tilde{\mathbf{p}}_{s,a}$ has mean $\bar{\mathbf{p}}_{s,a} \in \mathbb{R}^S$ and covariance matrix $\Sigma_{s,a} = I(\mathbf{p}_{s,a}^*)^{-1} \in \mathbb{R}^{S \times S}$ and Bernstein von Mises theorem [34] holds for posterior of transition probability $\tilde{\mathbf{p}}_{s,a}$, then as $N \rightarrow \infty$, the returns $\tilde{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a} \sim \mathcal{N} \left(\bar{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a}, \frac{\mathbf{w}_{s,a}^{\top} I(\mathbf{p}_{s,a}^*) \mathbf{w}_{s,a}}{N} \right)$. Then, we can write equation (24) as

$$f_{s,a}(\mathbf{v}) = \left(\bar{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a} - \Phi^{-1}(1 - \alpha) \|\mathbf{w}_{s,a}\|_{\Sigma_{s,a}^{-1}} \right) , \quad (25)$$

where $\Phi^{-1}(1 - \alpha)$ represents the $1 - \alpha$ percentile of standard normal distribution. (25) follows from the analytical form of VaR of a normal random variable, i.e., for any normal random variable \tilde{Y} with mean μ , variance matrix σ and confidence level α , $\text{VaR}_{\alpha}[\tilde{Y}] = \mu - \Phi^{-1}(1 - \alpha)\sigma$.

Since $f_{s,a}(\mathbf{v})$ is convex in \mathbf{v} when $\tilde{\mathbf{p}}_{s,a}$ is normally distributed, we can use the definition of support functions of a closed convex set [8] to construct a unique ambiguity set $\mathcal{P}_{s,a}^{\text{VaR}}$ of the form

$$\mathcal{P}_{s,a}^{\text{VaR}} = \left\{ \mathbf{p}_{s,a} \mid \bar{\mathbf{p}}_{s,a} \leq \mathbf{p}_{s,a}^{\top} \mathbf{w}_{s,a} - \Phi^{-1}(1 - \alpha) \|\mathbf{w}_{s,a}\|_{\Sigma_{s,a}^{-1}}, \forall \mathbf{v} \in \mathbb{R}^S \right\} . \quad (26)$$

The above equation implies that

$$\forall \mathbf{v}, \quad \min_{\mathbf{p}_{s,a} \in \mathcal{P}_{s,a}^{\text{VaR}}} \mathbf{p}_{s,a}^{\top} \mathbf{w}_{s,a} = \text{VaR}_{\alpha} \left[\tilde{\mathbf{p}}_{s,a}^{\top} \mathbf{w}_{s,a} \right] .$$

Using basic algebraic manipulations as shown in Proposition B.1, we can alternatively express the ambiguity set $\mathcal{P}_{s,a}^{\text{VaR}}$ in the form of a semi-ellipsoidal ball with radius $\Phi^{-1}(1 - \alpha)$.

$$\mathcal{P}_{s,a}^{\text{VaR}} = \left\{ \mathbf{p}_{s,a} \mid \frac{1}{\sqrt{N}} \|\bar{\mathbf{p}}_{s,a} - \mathbf{p}_{s,a}\|_{\Sigma_{s,a}^{-1}} \leq \frac{\Phi^{-1}(1 - \alpha)}{\sqrt{N}} \right\} . \quad (27)$$

We note that although $\Sigma_{s,a}$ are invertible, we can proceed as described in Section 3.1 in [15] to find an approximation of their inverses.

Proposition B.1. Consider the two ambiguity sets given in equations (26) and (27). These two representations are equivalent.

Proof. We begin with equation (26),

$$\begin{aligned} \mathbf{p}_{s,a}^\top \mathbf{w}_{s,a} &\stackrel{(a)}{\leq} \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} - \Phi^{-1}(1-\alpha) \|\mathbf{w}_{s,a}\|_{\Sigma_{s,a}^{-1}} \\ \mathbf{p}_{s,a}^\top \mathbf{w}_{s,a} - \bar{\mathbf{p}}_{s,a}^\top \mathbf{w}_{s,a} &\stackrel{(b)}{\leq} -\Phi^{-1}(1-\alpha) \|\mathbf{w}_{s,a}\|_{\Sigma_{s,a}^{-1}} \\ \mathbf{w}_{s,a}^\top (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a}) (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \mathbf{w}_{s,a} &\stackrel{(c)}{\leq} \Phi^{-1}(1-\alpha)^2 \mathbf{w}_{s,a}^\top \Sigma_{s,a} \mathbf{w}_{s,a} . \end{aligned}$$

The first and second equation follow from the definition of $\mathcal{P}_{s,a}^{\text{VaR}}$ and simple algebraic manipulations. The third equation follows by squaring on both sides.

Equation (c) is obtained by simply squaring Equation (b) on both sides. Using the basic properties of semi-positive definite matrices, we can write the above equation as

$$\begin{aligned} &(\Phi^{-1}(1-\alpha)^2 \Sigma_{s,a} - (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})(\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top) \succeq 0 \\ &(\Sigma_{s,a})^\top (\Phi^{-1}(1-\alpha)^2 \Sigma_{s,a} - (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})(\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top) \Sigma_{s,a} \succeq 0 \\ &(\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})(\Sigma_{s,a}^{-1})^\top ((\Phi^{-1}(1-\alpha)^2 \Sigma_{s,a} - (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})(\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top) \Sigma_{s,a}^{-1} (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})) \stackrel{(1)}{\succeq} 0 \\ &\Phi^{-1}(1-\alpha)^2 (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \Sigma_{s,a}^{-1} (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a}) - ((\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \Sigma_{s,a} (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a}))^2 \succeq 0 \end{aligned}$$

Equation (1) holds since $\mathbf{U}^\top \mathbf{M} \mathbf{U} \succeq 0 \forall \mathbf{U}, \mathbf{M} \succeq 0$.

$$\begin{aligned} (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})^\top \Sigma_{s,a} (\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a}) &\leq \Phi^{-1}(1-\alpha)^2 \\ \|(\mathbf{p}_{s,a} - \bar{\mathbf{p}}_{s,a})\|_{\Sigma_{s,a}^{-1}} &\leq \Phi^{-1}(1-\alpha) . \end{aligned}$$

The proof for the other direction is simply the reverse of this proof and hence we omit it. \square

B.11 Proof of Theorem 4.4

Theorem 4.4 (Asymptotic Radius of Bayesian Credible Regions). For any state s and action a , let $\mathcal{P}_{s,a}^{\text{BCR}}$ be any Bayesian credible region and $\bar{\mathbf{p}}_{s,a}$ be the maximum likelihood estimator based on data \mathcal{D} . Furthermore, $\xi < \sqrt{\chi_{S,1-\alpha}^2} / \Phi^{-1}(1-\alpha)$. Then, $\forall s \in \mathcal{S}, a \in \mathcal{A}$,

$$\lim_{N \rightarrow \infty} \sqrt{N} (\mathcal{P}_{s,a}^{\text{BCR}} - \bar{\mathbf{p}}_{s,a}) \not\subseteq \lim_{N \rightarrow \infty} \sqrt{N} \xi (\mathcal{P}_{s,a}^{\text{VaR}} - \bar{\mathbf{p}}_{s,a}) . \quad (10)$$

This theorem follows directly from Theorem 10 in [15]. For the sake of clarity, we re-derive the theorem below.

We note that to prove this theorem, we assume that Bernstein von Mises theorem holds for transition probabilities \tilde{P} .

Lemma B.2. For any positive semidefinite matrices \mathbf{A}, \mathbf{B} , and any \mathbf{v} , the following inequality holds true.

$$\| \|\mathbf{v}\|_{\mathbf{A}^{-1}} - \|\mathbf{v}\|_{\mathbf{B}^{-1}} \| \leq \sqrt{|\mathbf{v}^\top (\mathbf{A} - \mathbf{B}) \mathbf{v}|} \leq \|\mathbf{v}\|_{\mathbf{A}^{-1}} \sqrt{\|\mathbf{A} - \mathbf{B}\|_F} .$$

This proof of this lemma can be found in Lemma C.2 in [15]. We will use this lemma to prove the theorem.

For any state-action pair (s, a) , we will use the shorthand $\boldsymbol{\mu}_N$ and Σ_N to represent the mean $\bar{\mathbf{p}}_{s,a}$ and the covariance matrix $\Sigma_{s,a}$ of transition probabilities $\tilde{\mathbf{p}}_{s,a}$ for any state s and action a , such that $\bar{\mathbf{p}}_{s,a}$ and $\Sigma_{s,a}$ are estimated using N samples from the posterior distribution of transition probabilities f . We will also use the shorthand $\tilde{\mathbf{p}}$ to represent $\tilde{\mathbf{p}}_{s,a}$ and \mathbf{p}^* to represent $\mathbf{p}_{s,a}^*$.

Define $\mathcal{P}(\mathcal{D}, \tau)$ for any $\tau \in \mathbb{R}$ and dataset \mathcal{D} , as

$$\mathcal{P}(\mathcal{D}, \tau) = \left\{ \mathbf{p} \in \Delta^S \mid \frac{1}{N} \|\mathbf{p} - \boldsymbol{\mu}_N\|_{\Sigma_N} \leq \tau \right\} .$$

Notice that $\mathcal{P}(\mathcal{D}, \Phi^{-1}(1-\alpha)/\sqrt{N})$ is the asymptotic ambiguity set \mathcal{P}^{VaR} in Theorem 4.3. Let $\mathcal{P}(\mathcal{D})$ represent a Bayesian credible region for any dataset \mathcal{D} . Consider any \mathcal{D} such that $\mathcal{P}(\mathcal{D}) - \boldsymbol{\mu}_N \subseteq \xi(\mathcal{P}(\mathcal{D}, \Phi^{-1}(1-\alpha)/\sqrt{N}) - \boldsymbol{\mu}_N)$. Since $\mathcal{P}(\mathcal{D})$ is a credible region,

$$\begin{aligned} 1 - \varepsilon &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N \in \mathcal{P}(\mathcal{D}) - \boldsymbol{\mu}_N) \\ &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N \in \xi(\mathcal{P}(\mathcal{D}, \Phi^{-1}(1-\alpha)/\sqrt{N}) - \boldsymbol{\mu}_N)) \\ &= \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)^\top \boldsymbol{\Sigma}_N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N) \leq \xi^2 \Phi^{-1}(1-\alpha)^2/N) . \end{aligned}$$

Since $\xi \leq \frac{\sqrt{\chi_{S,1-\alpha}^2}}{\Phi^{-1}(1-\alpha)}$, there exists $\delta \geq 0$, such that $\xi^2 \Phi^{-1}(1-\alpha)^2 \leq \chi_{S,1-\delta-\alpha}^2 - \delta$. Fix such δ . Then,

$$\begin{aligned} &\Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)^\top \boldsymbol{\Sigma}_N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N) \leq \xi^2 \Phi^{-1}(1-\alpha)^2/N) \\ &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)^\top \boldsymbol{\Sigma}_N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N) \leq (\chi_{S,1-\delta-\alpha}^2 - \delta)/N) \\ &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)^\top I(\boldsymbol{p}_{s,a}^*)(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N) \leq (\chi_{S,1-\delta-\alpha}^2)/N) \\ &\quad + \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(N^{-1}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)^\top (N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}^*))(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N) \leq -\delta/N) \\ &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\|\sqrt{N}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)\|_{I(\boldsymbol{p}^*)^{-1}}^2 \leq \chi_{S,1-\alpha-\delta}^2) + Z(\mathcal{D}) \\ &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\|\zeta\|_{I(\boldsymbol{p}^*)^{-1}}^2 \leq \chi_{S,1-\alpha-\delta}^2) + R_N(\mathcal{D}) + Z(\mathcal{D}) . \end{aligned}$$

where $\zeta \sim \mathcal{N}(0, I(\boldsymbol{p}^*)^{-1})$, $R_N(\mathcal{D}) = \sup_{\mathcal{A}} |\Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\sqrt{N}(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)) - \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\zeta \in \mathcal{A})|$ denotes the total variational distance for the realization \mathcal{D} , where \mathcal{A} is a measurable subset of Δ^S (see Bernstein von Mises theorem in Theorem 6 in [15]), $I(\boldsymbol{p}^*)$ is the Fisher information matrix of $\Pr_{\mathcal{D}|\boldsymbol{p}^*}$, and $Z(\mathcal{D}) = \Pr((\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N)^\top (N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}_{s,a}^*))(\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N) > \frac{\delta}{N})$. The first probability is almost $1 - \alpha - \delta$. Thus, for any \mathcal{D} such that the theorem statement does not hold, $R_N(\mathcal{D}) + S(\mathcal{D}) > \delta$.

Fix $t > 0$ such that $\Pr(\|\zeta\|_{I(\boldsymbol{p}^*)^{-1}} > t) \leq 0.5\delta$. From, the second inequality in Lemma B.2,

$$\begin{aligned} Z(\mathcal{D}) &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\|\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N\|_{I(\boldsymbol{p}^*)^{-1}} \sqrt{\|N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}^*)\|_F} > \sqrt{\delta/N}) \\ &\leq \Pr_{\tilde{\boldsymbol{p}}|\mathcal{D}}(\|\tilde{\boldsymbol{p}} - \boldsymbol{\mu}_N\|_{I(\boldsymbol{p}^*)^{-1}} > t/\sqrt{N}) + \mathbb{I}(\sqrt{\|N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}^*)\|_F} > \sqrt{\delta}/t) \\ &\leq \Pr(\|\zeta\|_{I(\boldsymbol{p}^*)^{-1}} > t) + R(\mathcal{D}) + \mathbb{I}(\sqrt{\|N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}^*)\|_F} > \sqrt{\delta}/t) \\ &0.5\delta + R_N(\mathcal{D}) + \mathbb{I}(\sqrt{\|N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}^*)\|_F} > \sqrt{\delta}/t) . \end{aligned}$$

Thus, for any randomly drawn $\tilde{\mathcal{D}}$,

$$\begin{aligned} \Pr_{\tilde{\mathcal{D}}}(\text{Theorem does not hold}) &\leq \Pr_{\tilde{\mathcal{D}}}(R_N(\tilde{\mathcal{D}}) + Z(\tilde{\mathcal{D}}) > \delta) \\ &\leq \Pr_{\tilde{\mathcal{D}}}(2R_N(\tilde{\mathcal{D}}) + \mathbb{I}(\sqrt{\|N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}^*)\|_F} > \sqrt{\delta}/t) > 0.5\delta) \\ &\leq \Pr_{\tilde{\mathcal{D}}}(2R_N(\tilde{\mathcal{D}}) > 0.25\delta) + \Pr(\mathbb{I}(\sqrt{\|N^{-1} \boldsymbol{\Sigma}_N^{-1} - I(\boldsymbol{p}^*)\|_F} > \sqrt{\delta}/t) > 0.25\delta) . \end{aligned}$$

From the assumption that the Bernstein-Von-Mises theorem holds for $\tilde{\mathcal{P}}$, the first probability tends to zero, and since $N \boldsymbol{\Sigma}_N \rightarrow I(\boldsymbol{p}^*)^{-1}$ for a large number of samples, the second probability tends to zero as well. This proves the theorem.

Methods	Riverswim	Inventory	Population-Small	Population
VaR	100.27 ± 8.62	483.08 ± 0.2	-1117.57 ± 120.01	-1856.06 ± 83.74
BCR l_1	108.9 ± 10.56	391.39 ± 17.14	-2578.25 ± 52.02	-2956.65 ± 432.84
BCR l_∞	95.96 ± 0.0	254.43 ± 22.41	-5437.67 ± 23.13	-6422.09 ± 13.44
WBCR l_1	108.9 ± 10.56	481.07 ± 2.29	-2251.96 ± 342.54	-2468.16 ± 90.22
WBCR l_∞	95.96 ± 0.0	239.76 ± 0.0	-5133.85 ± 42.72	-5917.99 ± 69.29
VaRN	123.78 ± 7.67	482.92 ± 0.59	-1514.44 ± 12.31	-1806.12 ± 1.77

Table 2: Shows the mean and standard deviation of the robust (percentile) returns at $\delta = 0.30$ achieved by different robust methods in Riverswim, Inventory, Population-Small and Population domains.

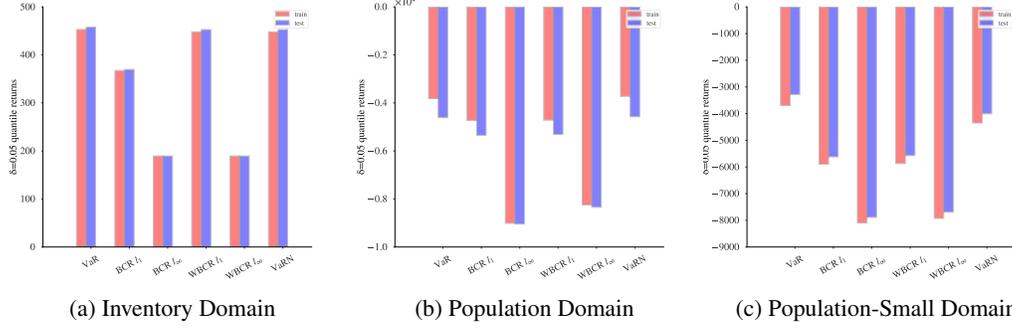


Figure 2: Comparison of test and train robust returns achieved by VaR , $VaRN$, $BCR l_1$, $BCR l_\infty$, $WBCR l_1$ and $WBCR l_\infty$ agents at confidence level $\delta = 0.05$ in Inventory, Population-Small and Population domain. VaR framework achieves the highest robust returns in all the domains on test and train datasets. All the RL agents are trained on the original train dataset.

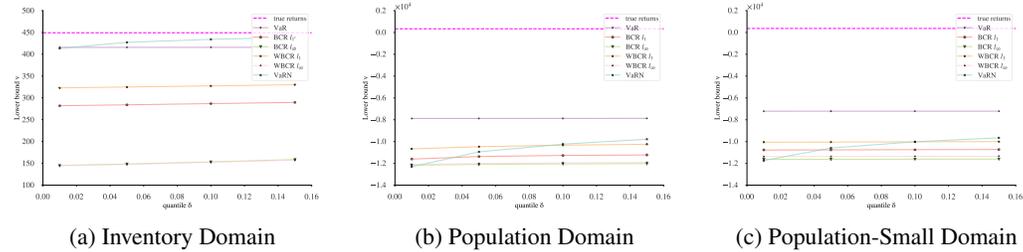


Figure 3: Comparison of robust lower bound value achieved by VaR , $VaRN$, $BCR l_1$, $BCR l_\infty$, $WBCR l_1$ and $WBCR l_\infty$ ambiguity sets for different confidence levels δ in Inventory, Population-Small and Population domain. $VaRN_\alpha$ achieves the highest robust returns in all the domains.

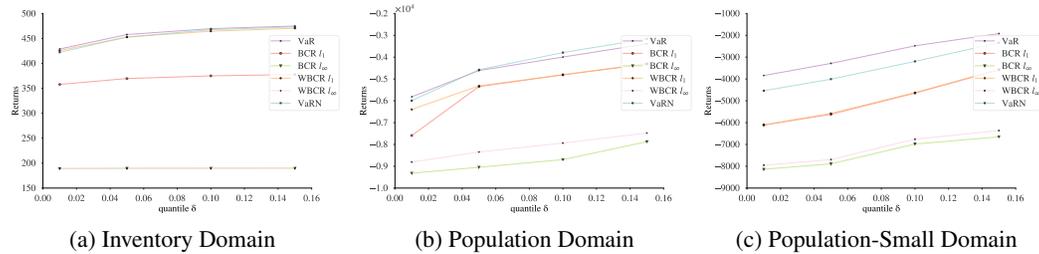


Figure 4: Comparison of robust returns on test dataset achieved by VaR , $VaRN$, $BCR l_1$, $BCR l_\infty$, $WBCR l_1$ and $WBCR l_\infty$ agents (trained on the original train dataset) for different confidence levels δ in Inventory and Population and Population-Small domain.

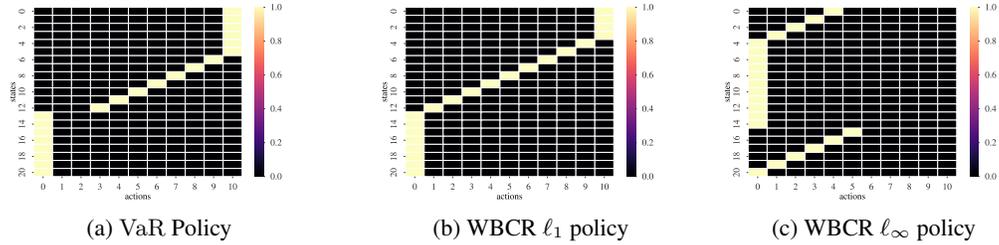


Figure 5: Comparison of robust policies corresponding to VaR, WBCR ℓ_1 , WBCR ℓ_∞ ambiguity sets at confidence level $\delta = 0.1$ in Inventory domain.

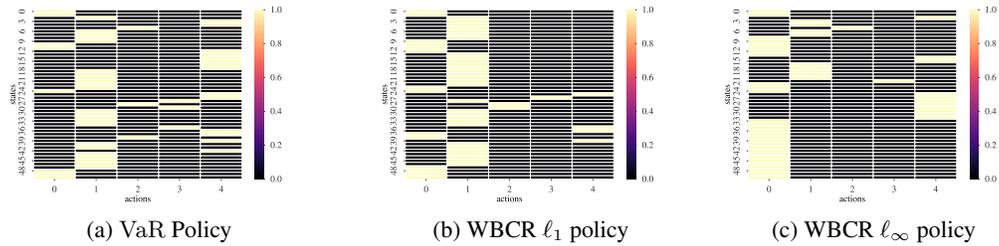


Figure 6: Comparison of robust policies corresponding to VaR, WBCR ℓ_1 , WBCR ℓ_∞ ambiguity sets at confidence level $\delta = 0.1$ in Population-Small domain.

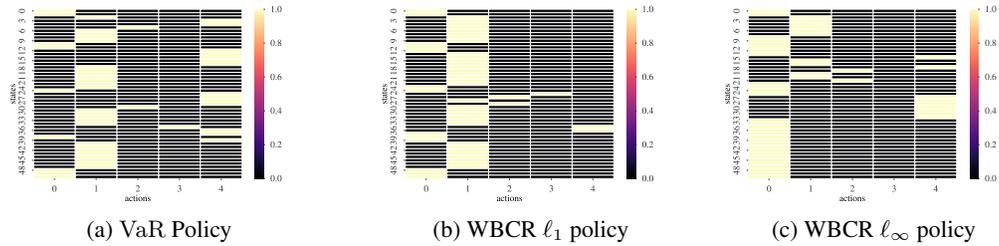


Figure 7: Comparison of robust policies corresponding to VaR, WBCR ℓ_1 , WBCR ℓ_∞ ambiguity sets at confidence level $\delta = 0.1$ in Population domain.

C Experiments

D Implementation Details

Hyperparameters for Riverswim Domain	
Number of train models	200
Number of test models	200
Hyperparameters for Inventory Domain	
Number of train models	100
Number of test models	100
Hyperparameters for Population-Small Domain	
Number of train models	50
Number of test models	50
Hyperparameters for Population Domain	
Number of train models	1000
Number of test models	1000

D.1 Code

We have provided the code in the supplementary materials. Since the dataset for the population domain is very large, we were unable to add it to the supplementary materials. We will make the dataset publicly available after the paper is published.

D.2 Machine Specifications

We ran all the experiments on MacBook Air (M2 2022) with 16GB Memory and 8 cores. The total computational time ~ 3 hours.