

BAYESIAN RISK-SENSITIVE POLICY GRADIENT FOR MDPs WITH GENERAL LOSS FUNCTIONS

Anonymous authors

Paper under double-blind review

ABSTRACT

Motivated by many application problems, we consider Markov decision processes (MDPs) with a general convex loss function and unknown parameters. To mitigate the epistemic uncertainty associated with unknown parameters, we take a Bayesian approach to estimate the parameters from data and impose a coherent risk functional (with respect to the Bayesian posterior distribution) on the loss. Since this formulation usually does not satisfy the interchangeability principle, it does not admit Bellman equations and cannot be solved by approaches based on dynamic programming. Therefore, We propose a policy gradient optimization method, leveraging the dual representation of coherent risk measures and extending the envelope theorem to continuous cases. We then show the algorithm converges to a stationary point with the rate of $\mathcal{O}(T^{-1/2} + r^{-1/2})$, where T is the number of policy gradient iterations and r is the sample size of the gradient estimator. We further extend our algorithm to an episodic setting, and show the extended algorithm converges to a globally optimal policy and provide bounds on the number of iterations needed to achieve an error bound $\mathcal{O}(\epsilon)$ in each episode.

1 INTRODUCTION

Markov decision process (MDP) is a paradigm for modeling sequential decision making under uncertainty, with a primary focus on identifying an optimal policy that minimizes the (discounted) expected total cost. However, the standard form of MDP is not sufficient for modeling some practical problems. For example, consider a self-driving car in a dynamic urban environment. It must reach its destination efficiently yet safely amid unpredictable events, calling for a general optimization objective within the MDP framework. At the same time, it faces incomplete knowledge of its environment, such as road conditions. Thus, the decision maker confronts two challenges: defining a general performance measure for intrinsic uncertainty and handling epistemic uncertainty about the environment. This paper addresses both challenges simultaneously within the MDP framework.

There is extensive literature addressing general loss functions and epistemic uncertainty separately. For instance, risk-sensitive objectives have been explored in the contexts of MDPs (Howard & Matheson, 1972; Ruszczyński, 2010; Mannor & Tsitsiklis, 2011; Petrik & Subramanian, 2012), stochastic optimal control (Borkar & Meyn, 2002; Moon, 2020), and stochastic programming (Shapiro, 2012; Pichler et al., 2022) literature. These objectives cannot be simply represented as the total expected cost. Epistemic uncertainty arises when some MDP parameters, such as transition probabilities, are unknown and must be estimated from available data. This discrepancy between the estimated and true MDP is referred to as epistemic uncertainty. Numerous approaches have been proposed to address epistemic uncertainty in MDPs, with robust MDP (Nilim & Ghaoui, 2004; Iyengar, 2005; Delage & Mannor, 2010; Wiesemann et al., 2013; Petrik & Russel, 2019) being one of the most widely adopted formulations. A more flexible and less conservative formulation, coined as Bayesian Risk MDP, was recently proposed by (Lin et al., 2022).

However, no existing literature addresses both a general loss function and epistemic uncertainty simultaneously. To the best of our knowledge, this paper is the first to consider this problem. In this work, we study MDPs with a general loss function, particularly focusing on convex functions in terms of the occupancy measure. Additionally, to handle both epistemic uncertainty and intrinsic uncertainty, we take a Bayesian approach to estimate the unknown parameters (such as transition probabilities) with a fixed batch of data, and impose a coherent risk functional (with respect to the

054 Bayesian posterior distribution) on the general convex loss function. Therefore, the problem is framed
055 as an offline optimization task with a two-layer composite objective: the outer general coherent
056 risk measure and the inner general convex loss function. To determine the optimal policy for this
057 composite problem, we use a policy optimization approach, which directly optimizes policies and
058 accommodates complex, high-dimensional representations such as neural networks. This method
059 typically employs parameterized policies and utilizes a policy gradient approach, introduced by
060 (Sutton et al., 1999), to search for the optimal solution. For the outer layer, the coherent risk measure
061 admits a dual representation as demonstrated by (Shapiro et al., 2021), which can be expressed as
062 the supremum of the expectation over a risk envelope set. We extended the envelope theorem in
063 (Milgrom & Segal, 2002) to obtain the policy gradient. A similar approach was taken by (Tamar et al.,
064 2015), but their consideration of a discrete parameter space limits the applicability of their method to
065 our problem. Our extension from the discrete case to the continuous case for the envelope theorem
066 may be of independent interest. The derived policy gradient involves the gradient of the loss function
067 with respect to the policy parameter, which can then be estimated by different methods. In particular,
068 we adapt the recent variational approach proposed by (Zhang et al., 2020) to construct the gradient
069 estimator. Other methods, such as zeroth-order estimation method proposed by (Balasubramanian &
070 Ghadimi, 2022), could also be used to estimate the inner gradient. By incorporating the inner gradient
071 estimator into the policy gradient, we derive the gradient estimator for the composed objective and
072 use policy gradient descent to optimize the problem. To make our approach more applicable with
073 new observed data, we further extend our approach to the episodic setting, where the agent iteratively
074 applies the current policy to gather more data and updates the policy based on new environment
estimates informed by the additional data.

075 Our choice of policy gradient method for this problem is not only due to its popularity but also
076 because algorithms based on dynamic programming are not applicable to general loss function that is
077 not in the standard form of cumulative sum. Therefore, our approach is completely different from
078 most robust MDPs or Bayesian risk MDPs which relies on dynamic programming. However, for
079 feasibility of the policy gradient method, we assume the general loss function is convex. The convex
080 loss functions are widely used, as discussed by (Pennings & Smidts, 2003), and are sufficiently
081 general to encompass many of the previously mentioned examples (e.g., risk-sensitive MDPs and
082 constrained MDPs) as special cases. More discussions about convex MDP is offered in Appendix B.
083 The standard expected total cost can be viewed as such a special case, where the loss function is a
084 linear function of the occupancy measure. The dynamic programming approach to solving MDPs
085 involves the use of Bellman equations. However, the derivation of Bellman equations relies on the
086 interchangeability principle, which may not hold for general convex loss functions. For a more
087 detailed discussion on why the interchangeability principle fails for general convex loss functions,
088 we refer readers to (Rockafellar & Wets, 2009) for the expectation operator and (Shapiro, 2017) for
089 general risk functionals. It is also worth noting that the Bayesian approach has been considered by
090 (Duff, 2002; Poupard et al., 2006; Abbasi-Yadkori & Szepesvári, 2015; Imani et al., 2018; Derman
091 et al., 2020; Lin et al., 2022; Wang & Zhou, 2023), where the Bayesian update accounts for future
data realization and enables the use of dynamic programming algorithms.

092 For our composite problem, there have been dedicated efforts to solve MDPs with some special
093 objectives using the policy gradient algorithm. For example, (Chow & Ghavamzadeh, 2014) applied
094 Conditional Value-at-Risk(CVaR) to the total cost and developed policy gradient and actor-critic
095 algorithms, each utilizing a distinct method to estimate the gradient and update policy parameters
096 in the descent direction. In contrast, we consider a broader composition of a general coherent
097 risk measure and a general loss function, allowing more flexible objectives. Note that although
098 the composition of a coherent risk measure and a convex loss function is convex in the occupancy
099 measure, it is generally non-convex in the policy parameters, which introduces additional challenges
100 for our convergence analysis. Finally, the work most relevant to ours is perhaps (Zhang et al., 2020),
101 which addresses a reinforcement learning (RL) problem with a general convex loss function and
102 derives the variational policy gradient theorem with a global convergence guarantee. However, our
103 work differs in that we consider an offline planning problem in an MDP with unknown transition
104 probabilities, which are estimated from data. Therefore, we address not only a general convex loss
105 function but also epistemic uncertainty. This introduces additional challenges related to risk measures,
and the robustness of the proposed formulation is a key consideration.

106 Our contributions are summarized as follows: (1) We propose a Bayesian risk formulation for MDPs
107 with a general convex loss function and develop a policy gradient algorithm to solve for the optimal

108 policy. The proposed formulation jointly mitigates both epistemic and intrinsic uncertainty; (2) We
 109 extend the envelope theorem to the dual representation of the coherent risk measure, and then apply the
 110 envelope theorem to derive the policy gradient. Our extension from the discrete case to the continuous
 111 case for the envelope theorem may be of independent interest; (3) We prove the convergence of the
 112 proposed algorithm and establish its convergence rate as $\mathcal{O}(T^{-1/2} + r^{-1/2})$, where T is the number
 113 of policy gradient iterations and r is the sample number of the gradient estimator; (4) We extend
 114 our policy gradient algorithm to the episodic setting, and prove the asymptotic convergence of the
 115 episodic minimizer of our Bayesian formulation to a global minimizer of the MDP problem under the
 116 true environment. Moreover, we show the number of iterations required in any episode to maintain
 117 an optimality gap $\mathcal{O}(\epsilon)$ under our Bayesian formulation.

118 2 PROBLEM FORMULATION

119 Consider an infinite-horizon Markov Decision Process (MDP) over a finite state space \mathcal{S} and a finite
 120 action space \mathcal{A} . For each state $s \in \mathcal{S}$ and action $a \in \mathcal{A}$, a transition to the next state s' follows the
 121 transition kernel P^* , i.e. $s' \sim P^*(\cdot|s, a)$. A stationary policy π is defined as a function mapping from
 122 the state space to a probability simplex $\Delta(\cdot)$ over the action space. Given any transition probability
 123 P , define $\lambda^{\pi, P}$ to be the discounted state-action occupancy measure under policy π :

$$124 \lambda_{sa}^{\pi, P} = \sum_{t=0}^{\infty} \gamma^t \cdot \mathbb{P}(s_t = s, a_t = a \mid \pi, s_0 \sim \tau, P) \quad (1)$$

125 for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, where τ is the initial distribution, $\gamma \in (0, 1)$ is the discount factor.

126 As mentioned in introduction, in many application problems such as a self-driving car in a dynamic
 127 urban environment, the decision maker faces two kinds of challenges: the epistemic uncertainty
 128 about the environment and a general performance measure for the intrinsic uncertainty. In this paper,
 129 we aim to address both challenges together. We consider a general loss function $F(\lambda, P)$ defined
 130 over the occupancy measure λ and transition kernel P , which is assumed to be convex in λ . In
 131 practice, the true distribution P^* is usually unknown and needs to be estimated. In this work, we take
 132 a Bayesian approach to estimate the environment. We assume that the transition kernel $P^* \equiv P_{\theta^*}$
 133 is parameterized by θ^* , where $\theta^* \in \Theta$ is the true but unknown parameter value, $\Theta \subseteq \mathbb{R}^p$ is the
 134 parameter space, and p is the dimension of Θ . Many real-world problems exhibit the characteristic
 135 of relying on a parametric assumption. In the example of a self-driving car, the noise in sensor
 136 measurements may be assumed to follow an unknown Gaussian distribution.

137 Under the parametric assumption, we assume we have access to some data which are state transitions
 138 $\zeta = (s, a, s')$, where s' follows the distribution $P_{\theta^*}(\cdot|s, a)$ and define $P_{\theta^*}(\zeta) := P_{\theta^*}(s'|s, a)$. Now
 139 given a fixed batch of data $\zeta^{(N)}$ of N samples, we can update the posterior distribution (denoted
 140 by μ_N) on the parameter θ using the Bayes rule: $\mu_N(\theta) = \frac{P_{\theta}(\zeta^{(N)})\mu_0(\theta)}{\int_{\Theta} P_{\theta'}(\zeta^{(N)})\mu_0(\theta')d\theta'}$, where μ_0 is a prior
 141 distribution of θ we assume. Furthermore, as discussed before, model mis-specification caused by the
 142 lack of data could lead to sub-optimality of the learned policy when it is implemented in a real-world
 143 setting. Hence, we further impose a risk functional on the objective with respect to (w.r.t.) the
 144 Bayesian posterior to account for the epistemic uncertainty, which results in the following composed
 145 formulation:

$$146 \min_{\pi} \rho_{\theta \sim \mu_N}(F(\lambda^{\pi, P_{\theta}}, P_{\theta})) \quad (2)$$

147 where ρ is a general coherent risk measure ¹ w.r.t. the posterior μ_N . We aim to solve problem
 148 equation 2 in this paper. Detailed introduction about coherent risk measures can be found in (Artzner
 149 et al., 1999). By this formulation, we look for a policy that minimizes a performance measure taking
 150 into account to the epistemic uncertainty caused by lack of data for a general convex loss function. If
 151 F is a linear function of λ , i.e. $F(\lambda, P) = \langle \lambda, c \rangle$ for a cost vector $c \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{A}|}$, and the posterior
 152 μ_N is a singleton on the true parameter θ^* , then the risk measure just considers the performance on
 153 this singleton and equation 2 reduces to the classical MDP problem. In Appendix H, we give two
 154 examples that are not in the classical form of MDP (but fall into our framework).

155 ¹Let $(\Omega, \mathcal{F}, \mathbb{P})$ w.r.t. the posterior μ_N be a probability space and \mathcal{X} be a linear space of \mathcal{F} -measurable
 156 functions $X : \Omega \rightarrow \mathbb{R}$. A risk measure is a function $\rho : \mathcal{X} \rightarrow \mathbb{R}$ which assigns to a random variable X a
 157 real number representing its risk. A coherent risk measure satisfies properties of monotonicity, sub-additivity,
 158 homogeneity, and translational invariance.

3 POLICY GRADIENT ALGORITHM: DERIVATION AND ESTIMATION

As discussed in the introduction, the dynamic programming type of algorithms may not be applicable to a general convex loss function $F(\cdot)$. So we adopt the policy gradient algorithm, which directly optimizes parameterized policies. Consider a stochastic parameterized policy $\pi_\alpha : \mathcal{S} \rightarrow \Delta(\mathcal{A})$, parameterized by $\alpha \in W \subset \mathbb{R}^d$. To directly work on the parameterized policy, we denote $F(\lambda^{\pi_\alpha, P_\theta}, P_\theta)$ by $C(\alpha, \theta)$. The policy optimization problem equation 2 then becomes:

$$\min_{\alpha} G(\alpha) := \rho_{\theta \sim \mu_N}(C(\alpha, \theta)). \quad (3)$$

It is worth mentioning that $G(\alpha)$ is not necessarily a convex function though F is convex w.r.t. λ . So we can only reach a stationary point of $G(\alpha)$ by the policy gradient descent method (see more detailed discussion in Section 4.2). In the rest of the section, we derive the policy gradient to the proposed formulation equation 3 using the envelope theorem, and construct the policy gradient estimator. It should be noted that our proposed formulation allows for flexible methods to estimate the policy gradient, including the variational approach such as in (Zhang et al., 2020), and the zeroth-order method such as in (Balasubramanian & Ghadimi, 2022).

3.1 PRELIMINARIES

To ensure the objective $G(\alpha)$ is well defined, we first assume that $C(\alpha, \theta) \in \mathcal{Z} := L_p(\Theta, \mu_N)$.

Assumption 3.1. $C(\alpha, \theta) \in \mathcal{Z} = \{f : \|f\|_p := (\int_{\Theta} |f(\theta)|^p d\mu_N(\theta))^{1/p} < \infty\}, \forall \alpha \in W$, for some $p \geq 1$.

The choice of p depends on the specific coherent risk measure. For example, p can be chosen as 1 for CVaR introduced in Example 1. Let $\mathcal{B} := \{\xi \in \mathcal{Z}^* : \int_{\Theta} \xi(\theta) \mu_N(\theta) d\theta = 1, \xi \succeq 0\}$, where $\mathcal{Z}^* := L_q(\Theta, \mu_N)$ is the dual space of \mathcal{Z} with $1/p + 1/q = 1$. As shown in (Shapiro et al., 2021), a coherent risk measure has a well-known dual representation.

Theorem 1. (Theorem 6.6 in (Shapiro et al., 2021).) *A risk measure $\rho : \mathcal{Z} \rightarrow \mathbb{R}$ is coherent if and only if there exists a convex bounded and closed set (also known as risk envelope) $\mathcal{U} = \mathcal{U}(\mu_N) \subset \mathcal{B}$ such that $\rho(Z) = \max_{\xi \in \mathcal{U}(\mu_N)} \mathbb{E}_{\xi}[Z]$, where $\mathbb{E}_{\xi}[Z] := \int_{\Theta \in \Theta} Z(\theta) \xi(\theta) \mu_N(\theta) d\theta$.*

Note ξ could be viewed as perturbation on the posterior μ_N that satisfies certain conditions, and the risk measure can be understood as the extreme performance for these perturbations. Theorem 1 implies that a functional ρ defined by $\rho(Z) = \max_{\xi \in \mathcal{U}} \mathbb{E}_{\xi}[Z]$ is a coherent risk measure if $\mathcal{U} \subset \mathcal{B}$ is convex, bounded and closed. In this paper we only focus on a class of coherent risk measures ρ following Definition 3.1 throughout the paper.

Definition 3.1. *For each given policy parameter $\theta \in \mathbb{R}^K$, there exists an expression for the risk envelope \mathcal{U} of the coherent risk measure ρ in the following form:*

$$\mathcal{U}(\mu_N) = \{\xi \in \mathcal{Z}^* : g_e(\xi, \mu_N) = 0, \forall e \in \mathcal{E}, f_i(\xi, \mu_N) \leq 0, \\ \forall i \in \mathcal{I}, \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) d\theta = 1, \xi(\theta) \geq 0, \|\xi\|_q \leq B_q\},$$

where $g_e(\xi, \mu_N)$ is an affine function in ξ , $f_i(\xi, \mu_N)$ is a convex function in ξ , $\|\cdot\|_q$ is the L_q norm in \mathcal{Z}^* , and there exists a strictly feasible point $\tilde{\xi}$. \mathcal{E} and \mathcal{I} here denote the sets of equality and inequality constraints, respectively. Furthermore, for any given $\xi \in \mathcal{B}$, $f_i(\xi, \mu_N)$ and $g_e(\xi, \mu_N)$ are twice differentiable in μ_N , and there exists a $M > 0$ such that $\forall \omega \in \Omega$:

$$\max \left\{ \max_{i \in \mathcal{I}} \left| \frac{df_i(\xi, \mu_N)}{d\mu_N(\theta)} \right|, \max_{e \in \mathcal{E}} \left| \frac{dg_e(\xi, \mu_N)}{d\mu_N(\theta)} \right| \right\} \leq M.$$

The conditions on g_e and f_i ensure that risk envelope $\mathcal{U}(\mu_N)$ is a convex closed set, and the condition $\|\xi\|_q \leq B_q$ makes $\mathcal{U}(\mu_N)$ bounded. A similar assumption is considered in Assumption 2.2 (Tamar et al., 2015). The assumption about bounded derivatives can be easily satisfied if Θ is compact. While (Tamar et al., 2015) only consider the case where Θ is finite, we extend it to the continuous case, leading to a functional problem over an infinite dimensional space instead of a finite-dimensional case. Therefore, we extend the result in (Tamar et al., 2015) to the infinite dimensional space, which

is shown in Theorem 2. The more technical discussion about differences between our approach and that of Tamar et al. (2015) is provided in Appendix E. Notably, the function forms of $g_e(\cdot)$ and $f_i(\cdot)$ can be exactly specified for a given coherent risk measure. We refer the readers to Appendix I and Section 6.3.2 (Shapiro et al., 2021) for some examples of the envelope set for coherent risk measures, which cover most common coherent risk measures.

3.2 DERIVATION OF POLICY GRADIENT

According to Theorem 1, we can write the coherent risk measure as a maximization problem equation 4, where the decision variable is ξ and the objective is a linear functional of ξ , and define the Lagrangian function equation 5 for problem equation 4:

$$\begin{aligned} \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) &= \max_{\xi: \xi \in \mathcal{U}(\mu_N)} \mathbb{E}_\xi[C(\alpha, \theta)] = \max_{\xi: \xi \in \mathcal{U}(\mu_N)} \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) C(\alpha, \theta) d\theta. \quad (4) \\ L_\alpha(\xi, \lambda^P, \lambda^\mathcal{E}, \lambda^\mathcal{I}) &= \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) C(\alpha, \theta) d\theta - \lambda^P \left(\int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) d\theta - 1 \right) \\ &\quad - \sum_{e \in \mathcal{E}} \lambda^\mathcal{E}(e) g_e(\xi, \mu_N) - \sum_{i \in \mathcal{I}} \lambda^\mathcal{I}(i) f_i(\xi, \mu_N). \quad (5) \end{aligned}$$

Using the Lagrangian relaxation equation 5, we derive the policy gradient for equation 4 in Theorem 2. For this purpose, we need some mild assumptions on the objective function.

Assumption 3.2. (1) $\nabla_\lambda F(\lambda, P)$ is uniformly bounded by $L_{F, \infty}$ for any λ and P w.r.t. $\|\cdot\|_\infty$; (2) $\nabla_\alpha C(\alpha, \theta)$ is $L_{\theta, 2}$ -Lipschitz continuous w.r.t. $\theta \in \Theta$ and $\|\cdot\|_2$ for any $\alpha \in W$; (3) $\nabla C(\alpha, \theta)$ is uniformly bounded by B for any $\alpha \in W$ and $\theta \in \Theta$ w.r.t. $\|\cdot\|_2$; (4) $\Theta \subseteq \mathbb{R}^p$ is compact and convex; (5) W , the domain of α , is bounded by B_W .

Assumption 3.2 requires the uniform boundedness and Lipschitz continuity of ∇C and ∇F , where $C(\alpha, \theta) = F(\lambda^{\pi_{\alpha, P_\theta}}, P_\theta)$. One sufficient condition easy to verify for Assumption 3.2 to hold is: each component in the composed function $F(\lambda^{\pi_{\alpha, P_\theta}}, P_\theta)$ is (somewhere) twice continuously differentiable w.r.t parameters α, θ , and the domains of two parameters are compact convex sets. Standard assumptions on the policy parameterization in RL are closely related to Assumption 3.2; we provide a detailed discussion and demonstrate the assumption on the classic Linear Quadratic Control (LQC) problem in Appendix D.

Theorem 2. Assume that Assumption 3.1 and 3.2 hold, and ρ satisfies Definition 3.1. Assume that μ_N is a Radon measure (see Appendix G.1 for definition). Define $\xi^* \in \arg \max_{\xi \geq 0, \|\xi\|_q \leq B_q} \min_{\lambda^P \geq 0, \lambda^\mathcal{E}, \lambda^\mathcal{I}} L_\alpha(\xi, \lambda^P, \lambda^\mathcal{E}, \lambda^\mathcal{I})$. Then we have the policy gradient

$$g(\alpha) := \nabla_\alpha \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) = \int_{\theta \in \Theta} \xi_\alpha^*(\theta) \mu_N(\theta) \nabla_\alpha C(\alpha, \theta) d\theta. \quad (6)$$

Proof details of Theorem 2 can be found in Appendix G.2. Theorem 2 implies that we can plug in a saddle point of Lagrangian equation 5 into equation 6 to get the policy gradient. However, equation 6 still involves ∇C , the gradient of the loss function, and the integration w.r.t. the posterior μ_N . In the next subsection, we show how to estimate the policy gradient in equation 6.

3.3 CONSTRUCTION OF THE POLICY GRADIENT ESTIMATOR

In this section, we focus on how to estimate the policy gradient $g(\alpha)$ and denote its estimator by $\hat{g}(\alpha)$. We first need to find ξ^* in Theorem 2. For some coherent risk measures, the closed-form expression of ξ^* is known. For CVaR with risk level $\beta \in (0, 1)$, $\xi^*(\theta) = \frac{1}{1-\beta}$ if $C(\alpha, \theta) \geq v_\beta$ and 0 otherwise, where v_β is the β -quantile of $C(\alpha, \theta)$. For a general coherent risk measure, we can use the approach sample average approximation (SAA). We first sample $\theta_k, k = 1, \dots, r$, from μ_N , and then solve the following SAA problem for the solution $\xi^*(\theta_k)$ for each k :

$$\begin{aligned} \max_{\substack{\xi \geq 0, \\ (\sum_{k=1}^r |\xi(\theta_k)|^q)/r \leq B_q}} \min_{\lambda^P \geq 0, \lambda^\mathcal{E}, \lambda^\mathcal{I}} &\frac{1}{r} \sum_{k=1}^r \xi(\theta_k) C(\alpha, \theta_k) - \lambda^P \left(\frac{1}{r} \sum_{k=1}^r \xi(\theta_k) - 1 \right) \\ &- \sum_{e \in \mathcal{E}} \lambda^\mathcal{E}(e) g_e(\xi^{(r)}, \mu_N(r)) - \sum_{i \in \mathcal{I}} \lambda^\mathcal{I}(i) f_i(\xi^{(r)}, \mu_N(r)) \quad (7) \end{aligned}$$

Notice the objective in equation 7 is linear w.r.t. $\lambda^{\mathcal{I}}, \lambda^{\mathcal{E}}$ and concave w.r.t ξ , and the domain of ξ is a convex bounded set in \mathbb{R}^r . Thus, equation 7 can be solved by any max-min optimization algorithm for a concave-convex function, such as alternating gradient descent ascent. Here we assume that we can solve equation 7 to derive $\xi^*(\theta_k)$ accurately for each k . Apart from ξ^* , we need to estimate $\nabla_{\alpha} C(\alpha, \theta)$ and the integral $\int_{\theta \in \Theta} \xi_{\alpha}^*(\theta) \mu_N(\theta) \nabla_{\alpha} C(\alpha, \theta)$. To estimate $\nabla_{\alpha} C(\alpha, \theta)$, any plug-in estimation method satisfies our demand. Here, we adopt the variational policy gradient theorem in (Zhang et al., 2020), which considers the policy gradient for a concave function defined on the occupancy measure for a RL problem. Different from our Bayesian-risk problem with a general loss function, (Zhang et al., 2020) only considers the inner-layer F of our objective equation 2 in the online setting. It should also be noticed that their method can be replaced by other methods such as the zeroth-order estimation method in (Balasubramanian & Ghadimi, 2022). While the variational policy gradient theorem require access to the conjugate function F^* , which may be difficult to calculate in some cases, zeroth-order method only requires function evaluation of F but leads to higher computational cost in general cases.

Lemma 3.1. (Theorem 3.1 in (Zhang et al., 2020)) Suppose F is convex and continuously differentiable in an open neighborhood of $\lambda^{\pi_{\alpha}, P_{\theta}}$. Fix the transition kernel P_{θ} and denote $V(\alpha; z)$ to be the expected cumulative cost of policy π_{α} when the cost function is z , and assume $\nabla_{\alpha} V(\alpha; z)$ always exists. Then we have

$$\nabla_{\alpha} C(\alpha, \theta) = - \lim_{\delta \rightarrow 0^+} \operatorname{argmin}_{x \in \mathbb{R}^{SA}} \sup_{z \in \mathbb{R}^{SA}} \{V(\alpha; z) + \delta \nabla_{\alpha} V(\alpha; z)^{\top} x - F^*(z) + \frac{\delta}{2} \|x\|^2\}, \quad (8)$$

where $V(\alpha; z) = \langle z, \lambda(\alpha, \theta) \rangle$, $\nabla_{\alpha} V(\alpha; z)^{\top} x = \langle z, \nabla_{\alpha} \lambda(\alpha, \theta) x \rangle$, $F^*(z) := \sup_{x \in \mathbb{R}^{SA}} x^{\top} z - F(x)$ is the Fenchel conjugate of F .

It may have a high computational cost if we directly estimate each part at a specific α in $\nabla_{\alpha} C = \nabla_{\lambda} F \cdot \nabla_{\alpha} \lambda$. The variational policy gradient method bypasses this issue by changing this problem into a problem of calculating some linear functions and the conjugate function at any z , shown in equation 8. (Zhang et al., 2020) considers an online setting and thus they need to interact with the environment to estimate $\nabla_{\alpha} C$. In our offline setting, we can directly solve equation 8 to get $\nabla_{\alpha} C$. An example algorithm to solve equation 8 is given in Appendix G.3. To evaluate the integral $\int_{\theta \in \Theta} \xi_{\alpha}^*(\theta) \mu_N(\theta) \nabla_{\alpha} C(\alpha, \theta) d\theta$, we use samples θ_k to construct the policy gradient estimator

$$\nabla_{\alpha} \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) \approx \widehat{g}(\alpha) := \frac{1}{r} \sum_{k=1}^r \xi^*(\theta_k) \nabla_{\alpha} C(\alpha, \theta_k). \quad (9)$$

In this paper, we assume the access to samples from the posterior distribution μ_N . While computing the posterior often requires costly methods like Markov Chain Monte Carlo (MCMC), utilizing a conjugate prior yields closed-form updates for the posterior parameters, improving efficiency. Computing the posterior typically requires expensive methods such as Markov Chain Monte Carlo (MCMC). However, by utilizing a conjugate prior, we obtain a closed-form expression for the posterior parameters, making the calculation more efficient.

3.4 FULL ALGORITHM AND EPISODIC SETTING

To perform policy gradient optimization, we iteratively use the gradient descent step equation 10, where η_t is the step size, and $\operatorname{Proj}_W(x) = \arg \min_{y \in W} \|y - x\|_2^2$ is the projection into the parameter space W . We summarize the full algorithm in Algorithm 1 in Appendix F.

$$\alpha_{t+1} = \arg \min_{\alpha \in W} \left(\widehat{g}(\alpha_t), \alpha - \alpha_t \right) + \frac{\eta_t}{2} \|\alpha - \alpha_t\|^2 = \operatorname{Proj}_W \left(\alpha_t - \frac{1}{\eta_t} \widehat{g}(\alpha_t) \right). \quad (10)$$

So far we have considered the offline setting with a fixed batch of data, but in many application problems data can be collected periodically. Again, consider a self-driving car as an example: the car is trained in an offline setting and then deployed to a real environment for a test drive while collecting more data from the environment. The collected data can be then used to learn about the environment and update the policy. This process can be repeated iteratively. Thus, we extend our approach to an episodic setting as described above. A potential approach is to use Algorithm 1 to make policy updates during each episode, as detailed in Algorithm 2 in Appendix F.

4 CONVERGENCE ANALYSIS

In this section, we analyze the convergence properties of Algorithm 1 and Algorithm 2. We begin by establishing the error bound for the policy gradient estimator. Next, we demonstrate the finite-time first-order convergence rate is $\mathcal{O}(T^{-1/2} + r^{-1/2})$, where T is the number of policy gradient iterations and r is the sample number of the gradient estimator. Furthermore, we prove the consistency of the proposed Bayesian risk formulation, meaning that the optimal policy obtained through this formulation converges to the one obtained by solving the true problem as the number of initial data points N approaches infinity. Lastly, for the episodic setting we show the number of iterations required in any episode to maintain an $\mathcal{O}(\epsilon)$ -error bound over all episodes, which implies the global convergence of Algorithm 2 as N goes to infinity.

4.1 ESTIMATION ERROR OF THE POLICY GRADIENT

Assumption 4.1. Assume ξ^* in Theorem 2 satisfies $\sup_{\alpha \in W} \text{Var}_{\theta \sim \mu_N} [\xi^*(\theta) \nabla C(\alpha, \theta)] = \sigma_\xi < \infty$.

Assumption 4.1 requires the uniformly bounded variance of $\xi^* \nabla C$. It is hard to show some property of ξ^* in a general case as the envelope set is given in a general form. One sufficient condition for Assumption 4.1 to hold is that ξ^* is bounded on Θ . As Θ is a compact and convex set, it is not a strong condition.

Theorem 3. Assume Assumption 3.1, 3.2 and 4.1 hold. By using r samples for gradient estimator in equation 9, the gradient estimation error is $\mathbb{E} [\|\hat{g}(\alpha) - g(\alpha)\|_2^2] \leq \frac{\sigma_\xi}{r}, \forall \alpha \in W$.

Theorem 3 implies that the sample complexity of $\Theta(1/\epsilon)$ is required to achieve the estimation error $\mathcal{O}(\epsilon)$. Please refer to Appendix G.4 for the detailed proof.

4.2 CONVERGENCE ANALYSIS

First we make an assumption about the Lipschitz continuity of $g(\alpha)$ in Assumption G.1.

Theorem 4. (Stationary convergence) Suppose that Assumption 3.1, 3.2, 4.1 and G.1 hold. By choosing $\eta_t = 2L_G$ in Algorithm 1, it holds that $\mathbb{E} \|\nabla G(\alpha_{out})\| \leq \mathcal{O}(T^{-1/2} + r^{-1/2})$.

Theorem 4 shows that the gradient bound of the output policy consists of two parts: an asymptotically diminishing error bound $T^{-1/2}$ in the exact setting and an estimation error bound $r^{-1/2}$ of the policy gradient. The total sample complexity from the posterior μ_N to achieve accuracy $\mathcal{O}(\epsilon)$ is $\mathcal{O}(\epsilon^{-4})$ by choosing $T = \epsilon^{-2}$ and $r = \epsilon^{-2}$. The proof and assumptions are shown in Appendix G.5. Because of the intrinsic non-convex structure under a fixed posterior, only convergence to a stationary point can be achieved under a fixed posterior (or in other words, in a fixed episode). Detailed discussion about the non-convex structure is provided in Appendix C. However, global convergence can be achieved in the episodic setting as the posterior updates and converges. This is shown in Corollary 1 later.

Theorem 5. (Consistency of episodic optimal policy) Suppose that Assumption 3.1, 3.2, 4.1, G.1 and G.2 hold. Define $G_i(\alpha) := \rho_{\theta \sim \mu_i}(C(\alpha, \theta))$, which is the objective for the posterior μ_i with data size i . Then we have $D_i := \sup_{\alpha \in W} |G_i(\alpha) - C(\alpha, \theta^*)|$ and $E_i := \sup_{\alpha \in W} \|\nabla_\alpha G_i(\alpha) - \nabla_\alpha C(\alpha, \theta^*)\|_2$ tend to 0 with probability 1 as $i \rightarrow \infty$, where the probability is w.r.t. the data-generating distribution. Moreover, $C(\alpha_i^*, \theta^*) - C(\alpha^*, \theta^*) \rightarrow 0$ with probability 1 as $i \rightarrow \infty$, where α_i^* is a global minimizer of $G_i(\alpha)$ and α^* is a global minimizer of $C(\alpha, \theta^*)$.

As the data size N increases, the posterior distribution converges to a Dirac measure, which is a point mass at the true parameter θ^* . Consequently, the performance of the optimal policy for the posterior μ_N converges to the optimal policy under the true environment θ^* , as demonstrated in Theorem 5. Additional assumptions are required to ensure the convergence of a series of posteriors. Broadly speaking, it is necessary that all parameters and all data points have positive probabilities of being sampled under both the prior and posterior distributions, and that the interchangeability of limits and integrals is satisfied. Detailed proof and assumptions for Theorem 5 are provided in Appendix G.6. In the episodic setting, we iteratively use the current policy for data collection and posterior updates, and perform policy updates based on the updated posterior, as described in Algorithm 2. A natural question arises: how many iterations are required within a given episode to achieve a certain level of accuracy? This is addressed in Theorem 6. Notably, the inner map $\lambda(\alpha, \theta^*)$

from policy parameter to occupancy measure is not necessarily convex in α , though $F(\lambda, \theta^*)$ is convex in λ . However, the hidden convexity can be utilized to get the global optimality under the true environment, regardless of the gradient estimation method. By utilizing the local bijection assumption of $\lambda(\cdot, \theta^*)$, a stationary point is still globally optimal, shown by Theorem 5.13 in (Zhang et al., 2021), which requires Assumption G.3. Assumption G.3 can be satisfied when λ is a locally differentiable bijection on a compact convex set W .

Theorem 6. (*Finite-episode error bound*) Suppose that Assumption 3.1, 3.2, 4.1, G.1, G.2 and G.3 hold. Assume that $G_i(\alpha)$ defined in Theorem 5 has $L_{G,i}$ -Lipschitz continuous gradient. Let $\{\alpha_{i,j}\}_{i=1}^N \}_{j=0}^{t_i}$ be the generated policy parameter sequences for N episodes by Algorithm 2. For any $\epsilon > 0$, if we choose $t_i = \Theta(L_{G,i}(E_{i-1} + D_i)\epsilon^{-2})$, $r = \Theta(\epsilon^{-2})$, then we can keep a constant gradient bound $\mathbb{E} \left[\left(\sum_{j=0}^{t_i-1} \|\nabla G_i(\alpha_{i,j})\|_2 \right) / t_i \right] \leq \epsilon$ for each i . Furthermore, $\mathbb{E}C(\alpha_{out}, \theta^*) - C(\alpha^*, \theta^*) \leq \mathcal{O}(\epsilon + E_N)$, where D_i, E_i are defined in Theorem 5.

Corollary 1. (*Global convergence*) Using the same assumptions and notations in Theorem 6, $\mathbb{E}C(\alpha_{out}, \theta^*) - C(\alpha^*, \theta^*) \leq \mathcal{O}(\epsilon)$ for any $\epsilon > 0$ as $N \rightarrow \infty$.

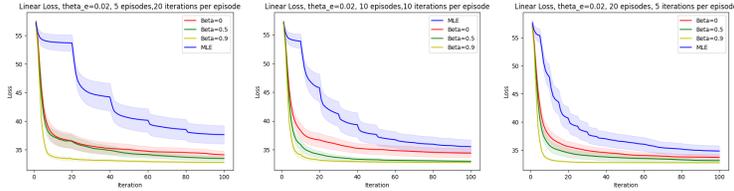
Theorem 6 offers theoretical advice on how to choose the iteration number in each episode. Generally speaking, we need fewer iterations to keep the gradient bound when i grows since D_i, E_i approaches 0. Corollary 1 is a direct result of Theorem 5 and Theorem 6, and Corollary 1 implies global convergence to the true optimal policy in the episodic setting (when the data size N increases to infinity). Detailed proof can be found in Appendix G.7.

5 NUMERICAL EXPERIMENTS

We evaluate our proposed formulation and algorithm on the offline Frozen Lake problem (Ravichandiran, 2018), an OpenAI benchmark. We refer readers to Appendix K for a detailed description of the problem. We consider different convex loss functions, including the mean and Kullback-Leibler (KL) divergence, for various tasks. We compare the Bayesian Risk Policy Gradient (BR-PG) algorithm with CVaR risk measure under different risk levels $\beta = 0, 0.5, 0.9$, respectively, with two other methods. The first is the empirical approach, which fits a maximum likelihood estimator (MLE) to the data and solves the MDP using the estimated parameters. The second is a modified offline version of distributionally robust Q-learning (DRQL) (Liu et al., 2022), which uses Q-learning to optimize worst-case performance over a KL divergence ball centered at the MLE kernel. When the risk level β approaches 1, Bayesian-risk performance is similar as the worst-case performance. Since we are considering an offline planning problem, we modify the DRQL to interact with an offline simulator that uses the transition kernel with the MLE parameters derived from the data. For a fair comparison, we conduct DRQL experiments with different radii of the KL divergence ball. **Linear Loss.** We consider the linear loss function, which corresponds to the total discounted cost in a classical MDP problem. This is referred to as one replication, and we repeat for 50 replications using different independent data sets. **Episodic Case.** We consider the episodic setting with 50 replications where the data collection and policy update are alternatively conducted. **Mimicking a policy.** Here we consider a different problem of mimicking an expert policy still using Frozen Lake environment and 50 replications. The loss function to minimize is defined as the KL divergence between state occupancy measure under the current policy and the expert state distribution. More implement details for three tasks can be found in Appendix K.

Table 1: Results for frozen lake problem. Linear loss and positive-sided variance at different risk levels α are reported for different algorithms and different data sizes with linear loss function. Standard errors are reported in parentheses. Escape probability $\theta_e = 0.02$ and number of data points is $N = 5$ and 50.

Approach	N=5		N=50	
	linear loss	positive-sided variance	linear loss	positive-sided variance
BR-PG ($\beta = 0$)	33.886(0.347)	5.212	32.784 (0.00825)	0.0026
BR-PG ($\beta = 0.5$)	33.104 (0.127)	0.710	32.757 (0.00516)	0.00119
BR-PG ($\beta = 0.9$)	32.854 (0.0641)	0.193	32.741 (0.00283)	0.000376
Empirical Method	37.057(0.927)	34.387	33.340 (0.0936)	0.380
DRQL(radius=0.05)	37.936(0.389)	26.354	34.365(0.366)	5.139
DRQL(radius=1)	35.216(0.732)	22.213	32.924(0.105)	0.519
DRQL(radius=20)	36.255(0.813)	24.622	32.855(0.063)	0.179
Optimal Policy under True Model	32.499		32.499	



(a) 5×20 total iterations (b) 10×10 total iterations (c) 20×5 total iterations

Figure 1: Results for episodic case with different episode numbers and iterations per episode under the same escape probability $\theta_e = 0.02$ and 50 replications. Here the loss function is still chosen to be the linear loss. 95% confidence intervals are reported by the shaded bands.

Conclusions. In each replication, data points are randomly sampled from the true distribution. While facing the epistemic uncertainty, BR-PG algorithm provides robustness across different loss functions. Table 1 shows that our BR-PG algorithm has lower linear loss, standard error and positive-sided variance (psv), demonstrating more robustness in the sense of balancing the mean and variability of the actual cost. In contrast, the empirical approach performs badly when the data size is small, e.g. $N = 5$, indicating that it is not robust against the epistemic uncertainty and suffers from the scarcity of data. DRQL also performs better than empirical method but worse than our algorithm in the sense of having larger mean and variance of the loss. Figure 1 shows that the loss of our algorithm decreases quickly in spite of few data. In the episodic case, the loss function decreases faster with more episodes (but the same total number of iterations), due to more collected data with more episodes. The loss function of our BR-PG method decreases more quickly in early episodes, which is shown by two differences between Figure 2a and Figure 2b. First, the 95% confidence interval, shown in the shaded band around each curve, is narrower for $N = 50$. Second, the absolute loss of $N = 50$ decreases by about 20% compared with $N = 5$. Figure 2 demonstrates the better performance of our proposed BR-PG algorithm compared to the empirical approach, where we achieve smaller loss and lower variability, for the policy mimicking task. From Table 1 and Figure 1, we can see when there are more data, the posterior distribution used in BR-PG algorithm and the MLE estimator used in the empirical approach converges to the true parameter as data size increases, which reduces to solving an MDP with known transition probability, and therefore, the optimal policies and the actual costs tend to be similar.

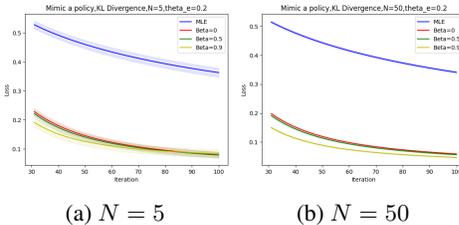


Figure 2: Results for loss function "KL Divergence" with data sizes $N = 5$ and 50 under $\theta_e = 0.02$. 95% confidence intervals are reported in the shaded area.

6 CONCLUSIONS

In this paper, we develop a Bayesian risk approach to jointly address both epistemic and intrinsic uncertainty in the infinite-horizon MDP. For a general coherent risk measure and a general convex loss function, we design a policy gradient algorithm for the proposed formulation and demonstrate the algorithm’s convergence at a rate of $\mathcal{O}(T^{-1/2} + r^{-1/2})$. Furthermore, we establish the consistency of an online episodic extension and provide bounds on the number of iterations required to maintain a constant gradient bound $\mathcal{O}(\epsilon)$ for each episode. The numerical experiments confirm the stationary analysis of the proposed algorithm and demonstrate the robustness of the formulation under various loss functions.

REPRODUCIBILITY STATEMENT

Numerical details appear in the appendices; all code is included in the supplementary material.

REFERENCES

- 486
487
488 Abbasi-Yadkori, Y. and Szepesvári, C. Bayesian optimal control of smoothly parameterized systems.
489 In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pp. 2–11,
490 2015.
- 491 Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. On the theory of policy gradient methods:
492 Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):
493 1–76, 2021.
- 494
495 Altman, E. *Constrained Markov decision processes*. Routledge, 2021.
- 496 Artzner, P., Delbaen, F., Eber, J.-M., and Heath, D. Coherent measures of risk. *Mathematical Finance*,
497 9(3):203–228, 1999.
- 498
499 Bai, Q., Bedi, A. S., Agarwal, M., Koppel, A., and Aggarwal, V. Achieving zero constraint violation
500 for concave utility constrained reinforcement learning via primal-dual approach. *Journal of*
501 *Artificial Intelligence Research*, 78:975–1016, 2023.
- 502
503 Balasubramanian, K. and Ghadimi, S. Zeroth-order nonconvex stochastic optimization: Handling
504 constraints, high dimensionality, and saddle points. *Foundations of Computational Mathematics*,
505 pp. 1–42, 2022.
- 506 Barakat, A., Fatkhullin, I., and He, N. Reinforcement learning with general utilities: Simpler variance
507 reduction and large state-action space. In *International Conference on Machine Learning*, pp.
508 1753–1800. PMLR, 2023.
- 509
510 Borkar, V. S. and Meyn, S. P. Risk-sensitive optimal control for Markov decision processes with
511 monotone cost. *Mathematics of Operations Research*, 27(1):192–209, 2002.
- 512 Brezis, H. and Brézis, H. *Functional analysis, Sobolev spaces and partial differential equations*,
513 volume 2. Springer, 2011.
- 514
515 Chow, Y. and Ghavamzadeh, M. Algorithms for CVaR optimization in MDPs. In Ghahramani,
516 Z., Welling, M., Cortes, C., Lawrence, N. D., and Weinberger, K. Q. (eds.), *Advances in Neural*
517 *Information Processing Systems*, 2014.
- 518
519 Delage, E. and Mannor, S. Percentile optimization for Markov decision processes with parameter
520 uncertainty. *Operations research*, 58(1):203–213, 2010.
- 521
522 Derman, E., Mankowitz, D., Mann, T., and Mannor, S. A Bayesian approach to robust reinforcement
523 learning. In *Proceedings of The 35th Uncertainty in Artificial Intelligence Conference*, pp. 648–658,
524 2020.
- 524
525 Duff, M. *Optimal Learning: Computational procedures for Bayes-adaptive Markov decision pro-*
526 *cesses*. Ph.D. dissertation, University of Massachusetts Amherst, Amherst, MA, 2002.
- 527
528 Hong, L. J. and Liu, G. Simulating sensitivities of conditional value at risk. *Management Science*, 55
529 (2):281–293, 2009.
- 529
530 Howard, R. A. and Matheson, J. E. Risk-sensitive Markov decision processes. *Management science*,
531 18(7):356–369, 1972.
- 532
533 Imani, M., Ghoreishi, S. F., and Braga-Neto, U. M. Bayesian control of large MDPs with unknown
534 dynamics in data-poor environments. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K.,
535 Cesa-Bianchi, N., and Garnett, R. (eds.), *Proceedings of the 31th International Conference on*
Neural Information Processing Systems, 2018.
- 536
537 Iyengar, G. N. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):257–280,
538 2005.
- 539
Lin, Y., Ren, Y., and Zhou, E. Bayesian risk Markov decision processes. In *Advances in Neural*
Information Processing Systems, volume 35, pp. 17430–17442, 2022.

- 540 Liu, Z., Bai, Q., Blanchet, J., Dong, P., Xu, W., Zhou, Z., and Zhou, Z. Distributionally robust
541 q -learning. In *International Conference on Machine Learning*, pp. 13623–13643. PMLR, 2022.
- 542
- 543 Mannor, S. and Tsitsiklis, J. N. Mean-variance optimization in Markov decision processes. In *Pro-*
544 *ceedings of the 28th International Conference on International Conference on Machine Learning*,
545 pp. 177–184, 2011.
- 546 Milgrom, P. and Segal, I. Envelope theorems for arbitrary choice sets. *Econometrica*, 70(2):583–601,
547 2002.
- 548
- 549 Moon, J. Generalized risk-sensitive optimal control and hamilton–jacobi–bellman equation. *IEEE*
550 *Transactions on Automatic Control*, 66(5):2319–2325, 2020.
- 551 Nesterov, Y. and Spokoiny, V. Random gradient-free minimization of convex functions. *Foundations*
552 *of Computational Mathematics*, 17:527–566, 2017.
- 553
- 554 Nilim, A. and Ghaoui, L. Robustness in Markov decision problems with uncertain transition matrices.
555 In Thrun, S., Saul, L., and Schölkopf, B. (eds.), *Advances in Neural Information Processing*
556 *Systems*, 2004.
- 557 Papini, M., Binaghi, D., Canonaco, G., Pirota, M., and Restelli, M. Stochastic variance-reduced
558 policy gradient. In *International conference on machine learning*, pp. 4026–4035. PMLR, 2018.
- 559
- 560 Pennings, J. M. and Smidts, A. The shape of utility functions and organizational behavior. *Manage-*
561 *ment Science*, 49(9):1251–1263, 2003.
- 562 Petrik, M. and Russel, R. H. Beyond confidence regions: Tight Bayesian ambiguity sets for robust
563 mdps. In Wallach, H., Larochelle, H., Beygelzimer, A., d'Alché-Buc, F., Fox, E., and Garnett, R.
564 (eds.), *Advances in Neural Information Processing Systems*, volume 32, 2019.
- 565
- 566 Petrik, M. and Subramanian, D. An approximate solution method for large risk-averse Markov
567 decision processes. In *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial*
568 *Intelligence*, pp. 805–814, 2012.
- 569
- 570 Pichler, A., Liu, R. P., and Shapiro, A. Risk-averse stochastic programming: Time consistency and
571 optimal stopping. *Operations Research*, 70(4):2439–2455, 2022.
- 572
- 573 Poupart, P., Vlassis, N., Hoey, J., and Regan, K. An analytic solution to discrete Bayesian reinforce-
574 ment learning. In *Proceedings of the 23rd International Conference on Machine Learning*, pp.
575 697–704, 2006.
- 576
- 577 Ravichandiran, S. *Hands-on reinforcement learning with Python: master reinforcement and deep*
578 *reinforcement learning using OpenAI gym and tensorflow*. Packt Publishing Ltd, 2018.
- 579
- 580 Rockafellar, R. T. and Wets, R. J.-B. *Variational analysis*, volume 317. Springer Science & Business
581 Media, 2009.
- 582
- 583 Ruszczyński, A. Risk-averse dynamic programming for Markov decision processes. *Mathematical*
584 *programming*, 125:235–261, 2010.
- 585
- 586 Shapiro, A. Minimax and risk averse multistage stochastic programming. *European Journal of*
587 *Operational Research*, 219(3):719–726, 2012.
- 588
- 589 Shapiro, A. Interchangeability principle and dynamic equations in risk averse stochastic programming.
590 *Operations Research Letters*, 45(4):377–381, 2017.
- 591
- 592 Shapiro, A., Dentcheva, D., and Ruszczyński, A. *Lectures on stochastic programming: modeling and*
593 *theory*. SIAM, 2021.
- 594
- 595 Shapiro, A., Zhou, E., and Lin, Y. Bayesian distributionally robust optimization. *SIAM Journal on*
596 *Optimization*, 33(2):1279–1304, 2023.
- 597
- 598 Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy gradient methods for reinforcement
599 learning with function approximation. In Solla, S. A., Leen, T. K., and Müller, K.-R. (eds.),
600 *Advances in Neural Information Processing Systems*, pp. 1057–1063, 1999.

594 Tamar, A., Chow, Y., Ghavamzadeh, M., and Mannor, S. Policy gradient for coherent risk measures.
595 *Advances in neural information processing systems*, 28, 2015.
596

597 Wang, Y. and Zhou, E. Bayesian risk-averse q-learning with streaming observations. In *Advances in*
598 *Neural Information Processing Systems*, 2023.

599 Wiesemann, W., Kuhn, D., and Rustem, B. Robust Markov decision processes. *Mathematics of*
600 *Operations Research*, 38(1):153–183, 2013.
601

602 Ying, D., Guo, M. A., Ding, Y., Lavaei, J., and Shen, Z.-J. Policy-based primal-dual methods
603 for convex constrained markov decision processes. In *Proceedings of the AAAI Conference on*
604 *Artificial Intelligence*, volume 37, pp. 10963–10971, 2023.

605 Ying, D., Zhang, Y., Ding, Y., Koppel, A., and Lavaei, J. Scalable primal-dual actor-critic method for
606 safe multi-agent rl with general utilities. *Advances in Neural Information Processing Systems*, 36,
607 2024.
608

609 Zhang, J., Koppel, A., Bedi, A. S., Szepesvari, C., and Wang, M. Variational policy gradient method
610 for reinforcement learning with general utilities. *Advances in Neural Information Processing*
611 *Systems*, pp. 4572–4583, 2020.

612 Zhang, J., Ni, C., Szepesvari, C., Wang, M., et al. On the convergence and sample efficiency of
613 variance-reduced policy gradient method. *Advances in Neural Information Processing Systems*, 34:
614 2228–2240, 2021.

615 Zhang, J., Bedi, A. S., Wang, M., and Koppel, A. Multi-agent reinforcement learning with general
616 utilities via decentralized shadow reward actor-critic. In *Proceedings of the AAAI Conference on*
617 *Artificial Intelligence*, volume 36, pp. 9031–9039, 2022.
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647

A THE USE OF LARGE LANGUAGE MODELS (LLMS)

In this work, the large language model (LLM) was used solely for text grammar refinement and code debugging.

B REVIEW ON CONVEX RL

Our problem is highly relevant to convex RL, which generalizes cumulative reward on a convex general-utility objective instead of cumulative reward. Specifically, our problem is closely tied to convex RL, which extends the traditional cumulative reward framework to a convex general-utility objective. Prior research has explored policy gradient methods to address convex RL. For instance, (Zhang et al., 2020) demonstrates that the policy gradient of convex RL can be formulated as a min-max optimization problem. To reduce estimator variance, (Zhang et al., 2021) introduces an off-policy policy gradient estimator that leverages mini-batch techniques and truncation mechanisms, while (Barakat et al., 2023) employs a recursive approach to handle large state-action spaces. In the domain of multi-agent convex RL, (Zhang et al., 2022) assumes global state observability and proposes a trajectory-based actor-critic method. Recent studies have also focused on safe convex RL, where the objective is to maximize a convex utility function under convex safety constraints. For example, (Ying et al., 2023) develops a primal-dual algorithm with strong guarantees on the optimality gap and constraint violations, achieving an $\mathcal{O}(1/\epsilon^3)$ bound in the convex-concave case with zero constraint violation. Building on this, (Bai et al., 2023) improves the bound to $\mathcal{O}(1/\epsilon^2)$. Furthermore, (Ying et al., 2024) extends the primal-dual framework to multi-agent convex safe RL.

C DISCUSSION ON NON-CONVEX STRUCTURE

For any fixed environment parameter θ , the set of occupancy measure is convex, which makes global convergence achievable. But the global convergence cannot be achieved under the Bayesian setting because the set of occupancy measure is nonconvex! Consider the simple case that $\Theta = \{\theta_1, \theta_2\}$ and the state space and action space are finite. Then for any fixed policy parameter α and environment parameter θ_i , the occupancy measure $\lambda(\alpha, \theta_i)$ is a $|S||A|$ -dimensional vector. Under any fixed θ_i , for any $t \in [0, 1]$ and policy parameters α_1, α_2 , the convex combination $t\lambda(\alpha_1, \theta_i) + (1-t)\lambda(\alpha_2, \theta_i)$ is also an occupancy measure corresponding to some policy α_3 under θ_i . This hidden convexity can be utilized to achieve global convergence for a fixed θ . However, the occupancy measure is a $2|S||A|$ -dimensional vector under the Bayesian setting, and we need an additional constraint to guarantee that the occupancy measures under different environments correspond to the same policy α :

$$\frac{\lambda(\alpha, \theta_1)_{s,a}}{\sum_{a'} \lambda(\alpha, \theta_1)_{s,a'}} = \frac{\lambda(\alpha, \theta_2)_{s,a}}{\sum_{a'} \lambda(\alpha, \theta_2)_{s,a'}}, \forall s, a.$$

This constraint means that the agent chooses the action a at state s with the same probability under two environments, which makes the set of occupancy measure nonconvex. In other words, for any $t \in [0, 1], \alpha_1, \alpha_2$, there may not exist a policy α_3 such that $t(\lambda(\alpha_1, \theta_1), \lambda(\alpha_1, \theta_2)) + (1-t)(\lambda(\alpha_2, \theta_1), \lambda(\alpha_2, \theta_2)) = (\lambda(\alpha_3, \theta_1), \lambda(\alpha_3, \theta_2))$. Thus, the global convergence cannot be achieved due to the lack of intrinsic convexity under the Bayesian setting with a fixed posterior. However, as the data size N increases, the posterior distribution converges to a Dirac measure, which is a point mass at the true parameter θ^* . Consequently, the performance of the optimal policy for the posterior μ_N converges to the optimal policy under the true environment, as demonstrated in Theorem 5. What's more, the global optimality gap will converge to any accuracy ϵ when the data size N increases, as shown in Theorem 6. In other words, global convergence can be achieved in our episodic setting (when the data size N increases to infinity).

D DISCUSSION ON ASSUMPTION 3.2

Recall that $C(\alpha, \theta) := F(\lambda^{\pi_\alpha, P_\theta}, P_\theta)$. By chain rule, $\nabla_\alpha C(\alpha, \theta) = \nabla_\lambda F \cdot \nabla_\alpha \lambda^{\pi_\alpha}$. Similar to classical Policy Gradient Theorem, it holds that $\nabla_\alpha \lambda^{\pi_\alpha}(s, a) = \lambda^{\pi_\alpha}(s, a) \nabla_\alpha \log \pi_\alpha(a | s)$. Thus, the behavior of $\nabla_\alpha C(\alpha, \theta)$ depends on the regularity of both $\nabla_\lambda F$ and $\nabla_\alpha \log \pi_\alpha(a | s)$. In classical RL analysis like Agarwal et al. (2021); Papini et al. (2018), one typically assumes that $\nabla_\alpha \log \pi_\alpha(a | s)$ is either Lipschitz continuous or uniformly bounded, and per-step rewards $r(s, a)$ are assumed to be bounded for all (s, a) pairs, which together implies the smoothness of the expected return with respect to the policy parameter. In our setting, we replace the usual expected return with a general convex loss function F . Consequently, we must impose smoothness conditions on F . Specifically, if $F(\lambda)$ is Lipschitz continuous with Lipschitz-continuous gradient, and $\nabla_\alpha \log \pi_\alpha(a | s)$ is Lipschitz continuous and bounded, then Assumption 3.2 is satisfied.

We can demonstrate that the classic Linear-Quadratic Regulator (LQR), perhaps the simplest continuous-action benchmark, satisfies all of these assumptions.

The dynamics is

$$s_{t+1} = As_t + Ba_t + w_t, \quad w_t \sim \mathcal{N}(0, \Sigma_w)$$

and the policy is Gaussian

$$a_t \sim \pi_\alpha(\cdot | s_t) = \mathcal{N}(Ks_t, \Sigma_a)$$

with parameter $\alpha = \text{vec}(K)$. Define the loss function F to be any convex function in λ . Since both the transition and the policy are affine and Gaussian, the joint law (s_t, a_t) is Gaussian at every t . Then the occupancy measure is a discounted summation of Gaussian distributions. Recall that the behavior of $\nabla_\alpha C(\alpha, \theta)$ depends on the regularity of both $\nabla_\lambda F$ and $\nabla_\alpha \log \pi_\alpha(a | s)$. In the example of LQR, the first and second derivatives of $\log \pi_\alpha(a | s)$ are bounded since the policy is Gaussian. Specifically, if $F(\lambda)$ is Lipschitz continuous with Lipschitz-continuous gradient, and since $\nabla_\alpha \log \pi_\alpha(a | s)$ is Lipschitz continuous and bounded in this LQR problem, then Assumptions 3.2 (2)(3) are satisfied.

In the example of LQR, Θ is a set containing possible A, B, Σ_w . Assumptions 3.2 (4) is about the compactness and convexity of Θ , which is not strong. One sufficient condition for Assumption 4.1 to hold is that ζ^* is bounded on Θ . As Θ is a compact and convex set, it is not a strong condition.

E DISCUSSION ON DIFFERENCES BETWEEN OUR METHOD AND TAMAR ET AL. (2015)

We want to obtain $\nabla_\alpha [\rho_{\theta \sim \mu_N}(C(\alpha, \theta))]$, where the derivative is taken with respect to the policy parameter α in the general loss function $C(\alpha, \theta)$. On the other hand, the problem in Tamar et al. (2015) is how to get $\nabla_\alpha [\rho_{\theta \sim \mu_N(\alpha)}(D(\theta))]$ for a random variable $D(\theta)$, where the derivative is taken with respect to α in the distribution μ_N . The difference in settings essentially leads to different forms of Lagrangians and causes the failure to apply their results to our setting. What’s more, Θ is a finite set in their setting, but Θ can be an uncountable continuous subset of some \mathbb{R}^d in our setting. As a result, when we use the Envelope Theorem to prove the result, we are facing an infinite-dimensional optimization problem over functions, which is a much harder problem than their finite-dimensional optimization problem over vectors. Briefly speaking, in our proof we ensure differentiability and integrability conditions hold uniformly, and construct a separable set of disturbance functions as the domain for Lagrangians.

The original envelope theorem, as presented by Milgrom & Segal (2002), primarily addresses finite-dimensional parameter spaces. Tamar et al. (2015) utilize a discrete-parameter space framework for policy gradients in risk-sensitive Markov Decision Processes (MDPs), restricting their applicability to problems with finite and discrete parameter settings. Our extension generalizes the envelope theorem to handle continuous parameter spaces, significantly broadening the applicability of policy gradient methods to a wider class of problems, such as those involving continuous uncertainty sets or continuous Bayesian posterior distributions over model parameters. Specifically, we address the additional complexities introduced by functional optimization in infinite-dimensional spaces, which involves ensuring differentiability and integrability conditions hold uniformly.

This generalization is nontrivial as it requires overcoming challenges associated with infinite-dimensional optimization, such as ensuring the boundedness and continuity of gradients and validating

strong duality under more complex integrative constraints. As such, our contribution facilitates the development of theoretically sound policy gradient methods capable of addressing a broader and more practical range of MDP formulations where uncertainty is represented continuously. This makes our methodology independently interesting to researchers in stochastic control, reinforcement learning, and risk-sensitive optimization.

F ALGORITHMS

Algorithm 1 Bayesian Risk Policy Gradient (BR-PG)

input: Initial α_0 , data $\zeta^{(N)}$ of size N , prior distribution $\mu_0(\theta)$, iteration number T ;

Calculate the posterior $\mu_N(\theta) = \frac{P_\theta(\zeta^{(N)})\mu_0(\theta)}{\int_{\Theta} P_{\theta'}(\zeta^{(N)})\mu_0(\theta')d\theta'}$;

for $t = 0$ to $T - 1$ **do**

Sample $\{\theta_k^t\}_{k=1}^r$ from $\mu_N(\theta)$;

Use the closed-form expression or solve equation 7 to get $\xi^*(\theta_k^t)$;

Solve equation 8 to get $\nabla_\alpha C(\alpha_t, \theta_k^t)$ for $k = 1, \dots, r$;

$\hat{g}_t := \frac{1}{r} \sum_{k=1}^r \xi^*(\theta_k^t) \nabla_\alpha C(\alpha_t, \theta_k^t)$;

$\alpha_{t+1} = \text{Proj}_W \left(\alpha_t - \frac{1}{\eta_t} \hat{g}_t \right)$;

end for

output: Choose α_{out} uniformly from $\alpha_0, \dots, \alpha_{T-1}$.

Algorithm 2 Episodic BR-PG

input: Initial α_0 , prior distribution $\mu_0(\theta)$, total episode number N .

Deploy policy $\pi(\alpha_0)$ to gain the initial data set $\zeta^{(1)}$.

for $i = 1$ to N **do**

if $i = 1$ **then**

$\alpha_{i,0} = \alpha_0$

else

Let $\alpha_{i,0}$ to be uniformly chosen from $\alpha_{i-1,0}, \dots, \alpha_{i-1,t_{i-1}-1}$;

end if

Deploy policy $\pi(\alpha_{i,0})$ to gain a data set $\zeta^{(i)}$.

Calculate the posterior $\mu_i(\theta) = \frac{P_\theta(\zeta^{(i)})\mu_{i-1}(\theta)}{\int_{\Theta} P_{\theta'}(\zeta^{(i)})\mu_{i-1}(\theta')}$;

Use Algorithm 1 with t_i iterations and initial point $\alpha_{i,0}$ to generate the policy parameter sequence $\alpha_{i,1}, \dots, \alpha_{i,t_i}$.

end for

output: Let α_{out} to be randomly and uniformly chosen from $\alpha_{N,0}, \dots, \alpha_{N,t_N-1}$.

810 G PROOF DETAILS

811 G.1 DEFINITION OF RADON MEASURE

812 **Definition G.1.** μ_N is a Radon measure on Θ if (i) $\mu_N(\Theta)$ is finite, (ii) for all Borel set $E \subseteq$
 815 Θ , we have $\mu_N(E) = \inf\{\mu_N(U) : E \subseteq U, U \text{ is open}\}$ and $\mu_N(E) = \sup\{\mu_N(K) : K \subseteq$
 816 $E, K \text{ is compact}\}$.

817 For a continuous parameter space Θ , if the prior is a continuous distribution and the likelihood
 818 function is continuous in θ , then the posterior is Radon. And it always holds for discrete cases. Thus
 819 it hold in most cases that we may care about, and most common probability distributions are Radon
 820 Measures.

821 G.2 PROOF OF THEOREM 2

822 *Proof.*

$$823 \mathcal{U}(\mu_N) = \{\xi : g_e(\xi, \mu_N) = 0, \forall e \in \mathcal{E}, f_i(\xi, \mu_N) \leq 0, \forall i \in \mathcal{I},$$

$$824 \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) = 1, \xi(\theta) \geq 0, \|\xi\|_q \leq B_q\}.$$

825 Define the Lagrangian:

$$826 L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E}) = \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) C(\alpha, \theta) - \sum_{i \in \mathcal{I}} \lambda^\mathcal{I}(i) f_i(\xi, \mu_N) - \sum_{e \in \mathcal{E}} \lambda^\mathcal{E}(e) g_e(\xi, \mu_N), \quad (11)$$

827 and a subtly relaxed envelope

$$828 \mathcal{U}'(\mu_N) = \{\xi : \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) = 1, \xi(\theta) \geq 0, \|\xi\|_q \leq B_q\}.$$

829 As mentioned before, we can rewrite the objective as the value of a max-min problem in equation 4

$$830 \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) = \max_{\xi \in \mathcal{U}'(\mu_N)} \min_{\lambda^\mathcal{I} \geq 0, \lambda^\mathcal{E}} L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E}).$$

831 Two things deserved to be noticed: (i) Slater's condition holds in the primal optimization problem
 832 equation 4 by Definition 3.1. (ii) $L_\theta(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$ is concave in ξ and convex in $(\lambda^\mathcal{I}, \lambda^\mathcal{E})$. Then strong
 833 duality holds for equation 4.

$$834 \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) = \max_{\xi \in \mathcal{U}'(\mu_N)} \min_{\lambda^\mathcal{I} \geq 0, \lambda^\mathcal{E}} L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$$

$$835 = \min_{\lambda^\mathcal{I} \geq 0, \lambda^\mathcal{E}} \max_{\xi \in \mathcal{U}'(\mu_N)} L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E}) \quad (12)$$

836 As $\nabla_\alpha C(\alpha, \theta)$ is uniformly bounded for all θ and α , we have $\nabla_\alpha L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$ is uniformly
 837 bounded w.r.t α and continuous at all $(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$. Then we have $L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$ is absolutely contin-
 838 uous w.r.t α for all $(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$. Since $\nabla_\alpha^2 C(\alpha, \theta)$ is uniformly bounded for all θ and α , we have
 839 $\{L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})\}_{(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})}$ is equi-differentiable in α .

840 As Θ is compact and convex, Θ is a separable metric space with Euclidean metric and its Borel sigma
 841 algebra. Then (Θ, μ_N) is a separable metric measure space. By Theorem 4.13 (Brezis & Brézis,
 842 2011), $L^q(\Theta, \mu_N)$ is separable. Then $\mathcal{U}'(\mu_N) = \{\xi \in L^q(\Theta, \mu_N) : \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) = 1, \xi(\theta) \geq$
 843 $0, \|\xi\|_q \leq B_q\}$ is separable.

844 Define the set of saddle point for equation 12 by $X^* = \arg \max_{\xi \in \mathcal{U}'(\mu_N)} \min_{\lambda^\mathcal{I} \geq 0, \lambda^\mathcal{E}} L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$
 845 and $Y^* = \arg \min_{\lambda^\mathcal{I} \geq 0, \lambda^\mathcal{E}} \max_{\xi \in \mathcal{U}'(\mu_N)} L_\alpha(\xi, \lambda^\mathcal{I}, \lambda^\mathcal{E})$.

846 Then for every selection of saddle point $(\xi_\alpha^*, \lambda_\alpha^{*\mathcal{E}}, \lambda_\alpha^{*\mathcal{I}}) \in X^* \times Y^*$, the Envelope theorem for
 847 saddle-point problems (Theorem 4(Milgrom & Segal, 2002)) shows that

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

$$\begin{aligned}
\nabla_{\alpha} \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) &= \nabla_{\alpha} \max_{\xi \in \mathcal{U}'(\mu_N)} \min_{\lambda^{\mathcal{I}} \geq 0, \lambda^{\mathcal{E}}} L_{\alpha}(\xi, \lambda^{\mathcal{I}}, \lambda^{\mathcal{E}}) \\
&= \nabla_{\alpha} L_{\alpha}(\xi, \lambda^{\mathcal{I}}, \lambda^{\mathcal{E}}) \Big|_{(\xi_{\alpha}^*, \lambda_{\alpha}^{*\mathcal{I}}, \lambda_{\alpha}^{*\mathcal{E}})} \\
&= \int_{\theta \in \Theta} \xi_{\alpha}^*(\theta) \mu_N(\theta) \nabla_{\alpha} C(\alpha, \theta)
\end{aligned} \tag{13}$$

□

G.3 PROOF OF LEMMA 3.1

Proof. Here is a brief proof sketch, and the full proof can be found in the proof of Theorem 3.1 (Zhang et al., 2020). For a convex function, the conjugate of the conjugate is itself. Notice that $V(\alpha; z) + \delta \nabla_{\alpha} V(\alpha; z)^{\top} x - F^*(z) = \langle z, \lambda(\alpha, \theta) + \delta \nabla_{\alpha} \lambda(\alpha, \theta) x \rangle - F^*(z)$. Then we have $\sup_{z \in \mathbb{R}^{SA}} V(\alpha; z) + \delta \nabla_{\alpha} V(\alpha; z)^{\top} x - F^*(z) = F(\lambda(\alpha, \theta) + \delta \nabla_{\alpha} \lambda(\alpha, \theta) x)$. By the first-order condition, we have

$$\operatorname{argmin}_{x \in \mathbb{R}^{SA}} F(\lambda(\alpha, \theta) + \delta \nabla_{\alpha} \lambda(\alpha, \theta) x) + \frac{\delta}{2} \|x\|_2^2 = -\nabla F(\lambda(\alpha, \theta) + \delta \nabla_{\alpha} \lambda(\alpha, \theta) x) \nabla_{\alpha} \lambda(\alpha, \theta) x.$$

By letting $\delta \rightarrow 0^+$ and using the chain rule, we get the result equation 8. □

G.3.1 ALGORITHM FOR SOLVING THEOREM 3.1

Estimate $V(\alpha, z)$: Recall that we consider an offline setting where the transition kernel P_{θ} is assumed to be known for any given θ . For any fixed transition kernel P_{θ} and policy π_{α} , we can estimate the occupancy measure by making a truncation K in the definition of occupancy measure in equation 1:

$$\widehat{\lambda}_{sa}^{\pi, P} = \sum_{t=0}^K \gamma^t \cdot \mathbb{P}(s_t = s, a_t = a \mid \pi, s_0 \sim \tau, P)$$

with the error $\|\widehat{\lambda} - \lambda\|_1 \leq \epsilon_{\lambda} := \gamma^K / (1 - \gamma)$. This error can be made arbitrarily small by increasing K , thus we assume that we can exactly compute occupancy measure. After computing the occupancy measure, $V(\alpha; z) = \langle z, \lambda \rangle$.

Estimate $\nabla_{\alpha} V(\alpha, z)$: The policy gradient theorem (Sutton et al., 1999) shows that

$$\nabla_{\alpha} V(\alpha; z) = \mathbb{E}^{\pi_{\alpha}} \left[\sum_{t=0}^{\infty} \gamma^t Q^{\pi_{\alpha}}(s_t, a_t; z) \cdot \nabla_{\alpha} \log \pi_{\alpha}(a_t \mid s_t) \right]$$

where $Q^{\pi}(s, a; z) := \mathbb{E}^{\pi} [\sum_{t=0}^{\infty} \gamma^t z(s_t, a_t) \mid s_0 = s, a_0 = a, a_t \sim \pi(\cdot \mid s_t)]$ satisfying the Bellman equation

$$Q^{\pi}(s, a; z) = \mathbb{E}[z(s, a)] + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} P_{\theta}(s' \mid s, a) \pi(a' \mid s') Q^{\pi}(s', a'; z). \tag{14}$$

For any given θ , policy π and cost function z , we can solve the Bellman equation equation 14 exactly to get $Q(\cdot, \cdot)$. It can be seen that $\nabla_{\alpha} V(\alpha; z)$ is a linear function of λ :

$$\nabla_{\alpha} V(\alpha; z) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} Q(s, a) \dot{\nabla}_{\alpha} \log \pi_{\alpha}(a \mid s) \lambda(s, a).$$

For any θ , policy π and cost function z , we can calculate the Q value function by solving the Bellman equation:

$$Q^{\pi}(s, a; z) = \mathbb{E}[z(s, a)] + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} P_{\theta}(s' \mid s, a) \pi(a' \mid s') Q^{\pi}(s', a'; z)$$

Then we can use Algorithm 3 to solve equation 8 in Lemma 3.1. It should be noticed that $\delta \nabla_{\alpha} V(\alpha; z)^{\top} x = \mathcal{O}(\delta)$ is omitted when calculating the gradient for z as $\delta \rightarrow 0$. Thus we omit this

term when calculating the gradient for z . To evaluate the integral $\int_{\theta \in \Theta} \xi_\alpha^*(\theta) \mu_N(\theta) \nabla_\alpha C(\alpha, \theta)$, we sample i.i.d θ_k from μ_N for $k = 1, \dots, r$, then we can construct the policy gradient estimator

$$\nabla_\alpha \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) \approx \widehat{g}(\alpha) := \frac{1}{r} \sum_{k=1}^r \xi^*(\theta_k) \nabla_\alpha C(\alpha, \theta_k).$$

Algorithm 3 Alternative Gradient Descent Method

input: initial z_0, x_0 , step sizes a_t, b_t , iteration number T , transition kernel parameter θ , policy parameter α ;

for $t = 0$ to $T - 1$ **do**

$$z_{t+1} = z_t + a_t [\lambda(\alpha, \theta) - \nabla F^*(z_t)]$$

$$x_{t+1} = x_t - b_t [\nabla_\alpha V(\alpha; z) + x_t], \text{ where } \nabla_\alpha V(\alpha; z) = \sum_{s,a} Q(s, a) \dot{\nabla}_\alpha \log \pi_\alpha(a | s) \lambda(s, a)$$

end for

output: $-x_T$.

G.4 PROOF OF THEOREM 3

Proof. By Theorem 2, the true gradient is

$$g(\alpha) = \int_{\theta \in \Theta} \xi_\alpha^*(\theta) \mu_N(\theta) \nabla_\alpha C(\alpha, \theta).$$

And our gradient estimator is

$$\widehat{g}(\alpha) := \frac{1}{r} \sum_{k=1}^r \xi^*(\theta_k) \nabla_\alpha C(\alpha, \theta_k).$$

Then we have

$$\mathbb{E} \|\widehat{g} - g\|_2^2 \leq \frac{1}{r} \mathbb{E} \|\xi^*(\theta_1) \nabla_\alpha C(\alpha, \theta_1) - \int_{\Theta} \xi^*(\theta) \mu_N(\theta) \nabla_\alpha C(\alpha, \theta) d\theta\|_2^2 \leq \frac{\sigma_\xi}{r}.$$

□

G.5 PROOF OF THEOREM 4

First, we make an assumption about G .

Assumption G.1. *There exists some $L_G > 0$ s.t. $g(\alpha)$ is L_G -Lipschitz continuous in α .*

Proof. For ease of notation, denote $g(\alpha_t)$ as g_t and $\widehat{g}(\alpha_t)$ as \widehat{g}_t . By Assumption G.1, we have

$$\begin{aligned} G(\alpha) &\leq G(\alpha_t) + \langle g_t, \alpha - \alpha_t \rangle + \frac{L_G}{2} \|\alpha - \alpha_t\|_2^2 \\ &\leq G(\alpha) + L_G \|\alpha - \alpha_t\|_2^2. \end{aligned} \tag{15}$$

Then we have

$$\begin{aligned}
972 & \\
973 & \\
974 & G(\alpha_{t+1}) \leq G(\alpha_t) + \langle \hat{g}_t, \alpha_{t+1} - \alpha_t \rangle + \langle g_t - \hat{g}_t, \alpha_{t+1} - \alpha_t \rangle + \frac{L_G}{2} \|\alpha_{t+1} - \alpha_t\|_2^2 \\
975 & \\
976 & \leq G(\alpha_t) + \langle \hat{g}_t, \alpha_{t+1} - \alpha_t \rangle + \frac{1}{2L_G} \|g_t - \hat{g}_t\|_2^2 + \frac{L_G}{2} \|\alpha_{t+1} - \alpha_t\|_2^2 + \frac{L_G}{2} \|\alpha_{t+1} - \alpha_t\|_2^2 \\
977 & \\
978 & = G(\alpha_t) + \langle \hat{g}_t, \alpha_{t+1} - \alpha_t \rangle + \frac{1}{2L_G} \|g_t - \hat{g}_t\|_2^2 + L_G \|\alpha_{t+1} - \alpha_t\|_2^2 \\
979 & \\
980 & = \min_{\alpha \in W} G(\alpha_t) + \langle \hat{g}_t, \alpha - \alpha_t \rangle + L_G \|\alpha - \alpha_t\|_2^2 + \frac{1}{2L_G} \|g_t - \hat{g}_t\|_2^2 \\
981 & \\
982 & = \min_{\alpha \in W} G(\alpha_t) + \langle g_t, \alpha - \alpha_t \rangle + L_G \|\alpha - \alpha_t\|_2^2 + \langle \hat{g}_t - g_t, \alpha - \alpha_t \rangle + \frac{1}{2L_G} \|g_t - \hat{g}_t\|_2^2 \\
983 & \\
984 & \leq \min_{\alpha \in W} G(\alpha_t) + \langle g_t, \alpha - \alpha_t \rangle + L_G \|\alpha - \alpha_t\|_2^2 + \frac{L_G}{2} \|\alpha - \alpha_t\|_2^2 + \frac{1}{2L_G} \|g_t - \hat{g}_t\|_2^2 + \frac{1}{2L_G} \|g_t - \hat{g}_t\|_2^2 \\
985 & \\
986 & = \min_{\alpha \in W} G(\alpha_t) + \langle g_t, \alpha - \alpha_t \rangle + \frac{3L_G}{2} \|\alpha - \alpha_t\|_2^2 + \frac{1}{L_G} \|g_t - \hat{g}_t\|_2^2 \\
987 & \\
988 & = G(\alpha_t) - \frac{1}{6L_G} \|g_t\|_2^2 + \frac{1}{L_G} \|g_t - \hat{g}_t\|_2^2 \\
989 & \\
990 &
\end{aligned}$$

991 where the first inequality comes from equation 15, the second inequality comes from Cauchy–Schwarz
992 inequality, the second equality holds because the definition of α_{t+1} , the third inequality holds because
993 of Cauchy–Schwarz inequality again, and the fifth equality holds by taking $\alpha = \alpha_t - \frac{1}{3L_G} g_t$.
994 Telescoping over t , we have

$$995 \\ 996 \\ 997 \frac{\sum_{t=0}^{T-1} \|g_t\|_2^2}{T} \leq \frac{6L_G}{T} (G(\alpha_0) - G(\alpha_T)) + 6 \sum_{t=0}^{T-1} \frac{\|g_t - \hat{g}_t\|_2^2}{T} \\ 998$$

999 Since

$$1000 \mathbb{E} \left[\|g_t - \hat{g}_t\|_2^2 \right] \leq \frac{\sigma_\xi}{r} \\ 1001$$

1002 Then we have

$$1003 \mathbb{E} \|g_{\text{out}}\|_2^2 \leq \frac{6L_G}{T} (G(\alpha_0) - \mathbb{E} G(\alpha_T)) + 6 \frac{\sigma_\xi}{r} \\ 1004$$

1005 Let $T = \epsilon^{-2}$, $r = \epsilon^{-2}$, then

$$1006 \mathbb{E} \|g_{\text{out}}\|_2^2 = \mathcal{O}(\epsilon^2) \\ 1007$$

1008 and then

$$1009 \mathbb{E} \|g_{\text{out}}\|_2 = \mathcal{O}(\epsilon) \\ 1010$$

1011 \square

1012 G.6 PROOF OF THEOREM 5

1013 **Assumption G.2.** (Assumption 3.1 in (Shapiro et al., 2023))

- 1014
- 1015 (1) The set Θ is convex compact with nonempty interior.
 - 1016
 - 1017 (2) $\ln \mu_0(\theta)$ is bounded on Θ , i.e., there are constants $c_1 > c_2 > 0$ such that $c_1 \geq \mu_0(\theta) \geq c_2$
1018 for all $\theta \in \Theta$.
 - 1019
 - 1020 (3) $P^*(\zeta) > 0$ for any $\zeta \in \Xi$.
 - 1021
 - 1022 (4) $P_\theta(\zeta) > 0$, and hence $\mu_N(\theta) > 0$, for all $\xi \in \Xi$ and $\theta \in \Theta$.
 - 1023
 - 1024 (5) $P_\zeta(\xi)$ is continuous in $\theta \in \Theta$.
 - 1025
 - (6) $\ln P_\theta(\zeta)$, $\theta \in \Theta$, is dominated by an integrable (w.r.t. P_*) function.

Assumption G.2 (1), (2) are used to guarantee the uniform convergence of posterior. Assumption G.2 (3), (4) require that all data points has positive probability to be sampled under the prior and posterior. Assumption G.2 (5), (6) are used to exchange the order of limit and integral.

With Assumption G.2, we are now ready to prove Theorem 5. Define a function $\psi(\theta) = \mathbb{E}_{P^*}[\ln P_\theta(\xi)]$ and let $\Theta^* := \{\theta' \in \Theta : \psi(\theta') = \inf_{\theta \in \Theta} \psi(\theta)\}$. For $\epsilon > 0$, define sets

$$V_\epsilon := \{\theta \in \Theta : \psi(\theta^*) - \psi(\theta) \geq \epsilon\}, U_\epsilon := \Theta \setminus V_\epsilon = \{\theta \in \Theta : \psi(\theta^*) - \psi(\theta) < \epsilon\}.$$

First we need to show two intermediate lemmas.

Lemma G.1. (Lemma 3.1. (Shapiro et al., 2023)) Suppose that Assumption G.2 holds. Then for $0 < \epsilon_2 < \epsilon_1 < \epsilon_0$, it follows that w.p. 1 for N large enough

$$\sup_{\theta \in V_{\epsilon_0}} \mu_N(\theta) \leq \kappa(\epsilon_2)^{-1} e^{-N(\epsilon_1 - \epsilon_2)},$$

where V_{ϵ_0} and U_{ϵ_0} are defined in (3.2), and $\kappa(\epsilon_2) := \int_{U_{\epsilon_2}} d\theta$.

Lemma G.2. Suppose that Assumption G.2 holds. $\forall \delta > 0, \exists \epsilon > 0$ such that $d(\theta, \Theta^*) < \delta$ for all $\theta \in U_\epsilon$.

Proof. We prove this lemma by contradiction. Suppose that $\exists \delta_0 > 0$ such that $\forall \epsilon > 0$, there exists $\theta \in \Theta$ satisfying $\psi(\theta^*) - \psi(\theta) < \epsilon$ and $d(\theta, \Theta^*) \geq \delta_0$.

Choose $\epsilon = \frac{1}{n}$ and then get a sequence $\{\theta_n\}_{n=1}^\infty$. As Θ is compact, there exists a subsequence of $\{\theta_n\}_{n=1}^\infty$ that converge to a point $\theta' \in \Theta$ satisfying $d(\theta', \Theta^*) \geq \delta_0$. As ψ is continuous, $\psi(\theta') = \psi(\theta^*)$. Contradiction! \square

Then we can prove Theorem 5

Proof. For any $\delta > 0$, we can choose ϵ_0 such that $d(\theta, \Theta^*) \leq \delta$ for $\theta \in U_{\epsilon_0}$. Then we have

$$\begin{aligned} & |\rho_{\theta \sim \mu_N}(C(\alpha, \theta)) - C(\alpha, \theta^*)| \\ &= \left| \max_{\xi: \xi \in \mathcal{U}(\mu_N)} \int_{\theta \in \Theta} \xi(\theta) \mu_N(\theta) [C(\alpha, \theta) - C(\alpha, \theta^*)] d\theta \right| \\ &\leq \max_{\xi: \xi \in \mathcal{U}(\mu_N)} \int_{U_{\epsilon_0}} \xi(\theta) \mu_N(\theta) |C(\alpha, \theta) - C(\alpha, \theta^*)| d\theta + \max_{\xi: \xi \in \mathcal{U}(\mu_N)} \int_{V_{\epsilon_0}} \xi(\theta) \mu_N(\theta) |C(\alpha, \theta) - C(\alpha, \theta^*)| d\theta \\ &\leq \sup_{\|\theta - \theta^*\| \leq \delta} |C(\alpha, \theta) - C(\alpha, \theta^*)| + 2 \sup_{\alpha \in W, \theta \in \Theta} |C(\alpha, \theta)| \max_{\xi: \xi \in \mathcal{U}(\mu_N)} \int_{V_\epsilon} \xi(\theta) \mu_N(\theta) d\theta \end{aligned}$$

By Holder's Inequality, we have

$$\begin{aligned} \int_{V_\epsilon} \xi(\theta) \mu_N(\theta) d\theta &= \int_{V_\epsilon} \xi(\theta) \mu_N(\theta)^{1/q} \mu_N(\theta)^{1/p} d\theta \\ &\leq \left[\int_{V_\epsilon} \xi(\theta)^q \mu_N(\theta) d\theta \right]^{1/q} \left[\int_{V_\epsilon} \mu_N(\theta) d\theta \right]^{1/p} \\ &\leq B_q \kappa(\epsilon_2)^{-1/p} e^{-N(\epsilon_1 - \epsilon_2)/p} \cdot \text{Vol}(\Theta)^{1/p} \end{aligned}$$

Thus

$$|\rho_{\theta \sim \mu_N}(C(\alpha, \theta)) - C(\alpha, \theta^*)| \leq \delta L_\theta + 2B_q \kappa(\epsilon_2)^{-1/p} e^{-N(\epsilon_1 - \epsilon_2)/p} \text{Vol}(\Theta)^{1/p} \sup_{\alpha \in W, \theta \in \Theta} |C(\alpha, \theta)|,$$

where the inequality holds because $C(\alpha, \theta)$ is B -Lipschitz continuous w.r.t. θ . This implies $D_N \rightarrow 0$ as $N \rightarrow \infty$ since δ is arbitrary. Then we have

$$\begin{aligned} & C(\alpha_N^*, \theta^*) - C(\alpha^*, \theta^*) \\ &= C(\alpha_N^*, \theta^*) - \rho_{\theta \sim \mu_N}(C(\alpha_N^*, \theta)) + \rho_{\theta \sim \mu_N}(C(\alpha_N^*, \theta)) - \rho_{\theta \sim \mu_N}(C(\alpha^*, \theta)) + \rho_{\theta \sim \mu_N}(C(\alpha^*, \theta)) - C(\alpha^*, \theta^*) \\ &\leq 2\delta B + 4B_q \kappa(\epsilon_2)^{-1/p} e^{-N(\epsilon_1 - \epsilon_2)/p} \text{Vol}(\Theta)^{1/p} \sup_{\alpha \in W, \theta \in \Theta} |C(\alpha, \theta)|, \end{aligned}$$

Let $N \rightarrow \infty$ and recall that δ is arbitrary, we get the result.

Define

$$E_N := \sup_{\alpha \in W} \|\nabla_{\alpha} \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) - \nabla_{\alpha} C(\alpha, \theta^*)\|_2.$$

Similar to D_N , we have

$$\begin{aligned} & \|\nabla_{\alpha} \rho_{\theta \sim \mu_N}(C(\alpha, \theta)) - \nabla_{\alpha} C(\alpha, \theta^*)\|_2 \\ &= \left\| \int_{\Theta} \xi_{\alpha}^*(\theta) \mu_N(\theta) [\nabla_{\alpha} C(\alpha, \theta) - \nabla_{\alpha} C(\alpha, \theta^*)] d\theta \right\|_2 \\ &\leq \int_{U_{\epsilon_0}} \xi_{\alpha}^*(\theta) \mu_N(\theta) \|\nabla_{\alpha} C(\alpha, \theta) - \nabla_{\alpha} C(\alpha, \theta^*)\|_2 d\theta + \int_{V_{\epsilon_0}} \xi_{\alpha}^*(\theta) \mu_N(\theta) \|\nabla_{\alpha} C(\alpha, \theta) - \nabla_{\alpha} C(\alpha, \theta^*)\|_2 d\theta \\ &\leq \delta L_{\theta,2} + 2B_q \kappa(\epsilon_2)^{-1/p} e^{-N(\epsilon_1 - \epsilon_2)/p} \text{Vol}(\Theta)^{1/p} \sup_{\alpha \in W, \theta \in \Theta} \|\nabla_{\alpha} C(\alpha, \theta)\|, \end{aligned}$$

which implies $E_N \rightarrow 0$ as $N \rightarrow \infty$ since δ is arbitrary. \square

G.7 PROOF OF THEOREM 6 AND COROLLARY 1

Assumption G.3. (Assumption 5.11 in (Zhang et al., 2021)) For policy parameterization π_{α} , α overparametrizes the set of policies in the following sense. (i). For any α and $\lambda(\alpha)$ under the true environment P_{θ^*} , there exist (relative) neighbourhoods $\alpha \in \mathcal{U}_{\alpha} \subset W$ and $\lambda(\alpha) \in \mathcal{V}_{\lambda(\alpha)} \subset \lambda(W, \theta^*)$ s.t. $(\lambda|_{\mathcal{U}_{\alpha}})(\cdot)$ forms a bijection between \mathcal{U}_{α} and $\mathcal{V}_{\lambda(\alpha)}$, where $(\lambda|_{\mathcal{U}_{\alpha}})(\cdot)$ is the confinement of λ onto \mathcal{U}_{α} . We assume $(\lambda|_{\mathcal{U}_{\alpha}})^{-1}(\cdot)$ is ℓ_{α} -Lipschitz continuous and $(\lambda|_{\mathcal{U}_{\alpha}})(\cdot)$ is L_{λ} -Lipschitz smooth for any α . (ii). Let π_{α^*} be the optimal policy under the true environment. Assume there exists $\bar{\epsilon}$ small enough, s.t. $(1 - \epsilon)\lambda(\alpha) + \epsilon\lambda(\alpha^*) \in \mathcal{V}_{\lambda(\alpha)}$ for $\forall \epsilon \leq \bar{\epsilon}, \forall \alpha$.

Under the true environment P_{θ^*} , the set of all occupancy measures is a convex set, and there is a bijection between all policies and all occupancy measures. More discussions can be found in Section 5.2 in (Zhang et al., 2021). Based on this observation, we have $\min_{\pi} F(\lambda^{\pi}, \theta^*) = \min_{\lambda} F(\lambda, \theta^*)$, which turns the non-convex policy optimization problem into a convex optimization problem over occupancy measure. Then any stationary point is also globally optimal, which is shown in the following lemma.

Lemma G.3. Assume that Assumption G.3 holds. Then $C(\bar{\alpha}, \theta^*) - C(\alpha^*, \theta^*) \leq \mathcal{O}(\|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2), \forall \bar{\alpha} \in W$.

Proof. Notice that

$$\partial(F \circ \lambda)(\alpha) = \nabla_{\alpha} \lambda(\alpha)^{\top} \partial F(\lambda), \forall \alpha \in W.$$

So there exists $\bar{w} \in \partial F(\bar{\lambda}, \theta^*)$ s.t. $\nabla_{\alpha} C(\bar{\alpha}, \theta^*) = \nabla_{\alpha} \lambda(\bar{\alpha})^{\top} \bar{w}$. Then for any $\lambda(\alpha) \in \mathcal{V}_{\lambda(\bar{\alpha})}$, we have

$$\begin{aligned} & \langle \bar{w}, \lambda - \bar{\lambda} \rangle \\ &= \langle \bar{w}, \nabla_{\alpha} \lambda(\bar{\alpha})(\alpha - \bar{\alpha}) \rangle + \langle \bar{w}, \lambda - \bar{\lambda} - \nabla_{\alpha} \lambda(\bar{\alpha})(\alpha - \bar{\alpha}) \rangle \\ &= I_1 + I_2 \end{aligned} \tag{16}$$

For the first term, we have

$$I_1 = \langle \nabla_{\alpha} C(\bar{\alpha}, \theta^*), \alpha - \bar{\alpha} \rangle \geq -\|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2 \|\alpha - \bar{\alpha}\|_2 \geq -\|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2 \ell_{\alpha} \|\lambda - \bar{\lambda}\|_2.$$

For the second term, we have

$$I_2 \geq -\|\bar{w}\|_2 \cdot \frac{L_{\lambda}}{2} \|\alpha - \bar{\alpha}\|_2^2 \geq -\frac{L_{\lambda} \ell_{\alpha}^2}{2} \|\bar{w}\|_2 \|\lambda - \bar{\lambda}\|_2^2$$

Then we have

$$\langle \bar{w}, \lambda - \bar{\lambda} \rangle \geq -\|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2 \ell_{\alpha} \|\lambda - \bar{\lambda}\|_2 - \frac{L_{\lambda} \ell_{\alpha}^2}{2} \|\bar{w}\|_2 \|\lambda - \bar{\lambda}\|_2^2.$$

1134 Replace λ by $(1 - \epsilon)\lambda(\bar{\alpha}) + \epsilon\lambda(\alpha^*)$ for any $\epsilon \in (0, \bar{\epsilon}]$ and then it holds

$$1135 \quad \epsilon \langle \bar{w}, \lambda(\alpha^*) - \lambda(\bar{\alpha}) \rangle \geq -\epsilon \|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2 \ell_{\alpha} \|\lambda(\alpha^*) - \lambda(\bar{\alpha})\|_2 - \frac{L_{\lambda} \epsilon^2 \ell_{\alpha}^2}{2} \|\bar{w}\|_2 \|\lambda(\alpha^*) - \lambda(\bar{\alpha})\|_2^2.$$

1138 Divide both sides by ϵ and let $\epsilon \rightarrow 0$, we have

$$1139 \quad \langle \bar{w}, \lambda(\alpha^*) - \lambda(\bar{\alpha}) \rangle \geq -\|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2 \ell_{\alpha} \|\lambda(\alpha^*) - \lambda(\bar{\alpha})\|_2.$$

1141 Finally,

$$1142 \quad C(\bar{\lambda}, \theta^*) - C(\lambda^*, \theta^*) \leq -\langle \bar{w}, \lambda(\alpha^*) - \lambda(\bar{\alpha}) \rangle \leq \ell_{\alpha} \|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2 \|\lambda(\alpha^*) - \lambda(\bar{\alpha})\|_2 \leq \ell_{\alpha} D_{\lambda} \|\nabla_{\alpha} C(\bar{\alpha}, \theta^*)\|_2,$$

1144 where $D_{\lambda} = \sup_{\alpha, \alpha' \in W} \|\lambda(\alpha) - \lambda(\alpha')\|_2$. \square

1146 Then we can turn to prove Theorem 6.

1148 *Proof.* If $\mathbb{E} \left[\frac{\sum_{j=0}^{t_i-1} \|\nabla G_i(\alpha_{i,j})\|}{t_i} \right] \leq \epsilon$, choose $\alpha_{i+1,0}$ uniformly from $\alpha_{i,0}, \dots, \alpha_{i,t_i-1}$. Then
1150 $\mathbb{E} \|\nabla G_i(\alpha_{i+1,0})\|_2 \leq \epsilon$ and then

$$1151 \quad \mathbb{E} \|\nabla C(\alpha_{i+1,0}, \theta^*)\|_2 \leq \epsilon + E_i.$$

1153 By Lemma G.3 we have

$$1154 \quad \mathbb{E} C(\alpha_{i+1,0}, \theta^*) - C(\alpha^*, \theta^*) \leq (\epsilon + E_i) \ell_{\alpha} D_{\lambda},$$

1156 and then

$$1157 \quad \mathbb{E} G_{i+1}(\alpha_{i+1,0}) - G_{i+1}(\alpha_{i+1}^*) \leq (\epsilon + E_i) \ell_{\alpha} D_{\lambda} + 2D_{i+1}.$$

1158 By Theorem 4 we have

$$1159 \quad \mathbb{E} \left[\frac{\sum_{j=0}^{t_{i+1}-1} \|\nabla G_{i+1}(\alpha_{i+1,j})\|_2^2}{t_{i+1}} \right] \leq \frac{6L_{G,i+1} [(\epsilon + E_i) \ell_{\alpha} D_{\lambda} + 2D_{i+1}] + 6\frac{\sigma_{\xi}}{r}}{t_{i+1}}$$

1162 Then we can choose

$$1163 \quad t_{i+1} = 12L_{G,i+1} [(\epsilon + E_i) \ell_{\alpha} D_{\lambda} + 2D_{i+1}] \epsilon^{-2}$$

$$1164 \quad r = 12\sigma_{\xi} \epsilon^{-2}$$

1166 to make

$$1167 \quad \mathbb{E} \left[\frac{\sum_{j=0}^{t_{i+1}-1} \|\nabla G_{i+1}(\alpha_{i+1,j})\|_2^2}{t_{i+1}} \right] \leq \epsilon^2.$$

1169 Then by Jensen's inequality we have

$$1170 \quad \mathbb{E} \left[\frac{\sum_{j=0}^{t_{i+1}-1} \|\nabla G_{i+1}(\alpha_{i+1,j})\|_2}{t_{i+1}} \right] \leq \epsilon.$$

1174 When $\mathbb{E} \|\nabla G_N(\alpha_{\text{out}})\|_2 \leq \epsilon$, we have $\|\nabla C(\alpha_{\text{out}}, \theta^*)\|_2 \leq \epsilon + E_N$. Then by Lemma G.3 we complete
1175 the proof of Theorem 6.

1176 As Theorem 5 shows that $E_N \rightarrow \infty$ when $N \rightarrow \infty$, then we complete the proof of Corollary 1. \square

1179 H EXAMPLES OF LOSS FUNCTION

1181 **Example 1** (Risk-Averse Constrained MDP). *In safe RL problems, one usually considers a con-*
1182 *strained MDP (Altman, 2021), where the goal is to minimize the total expected discounted cost under*
1183 *a risk-averse constraint. Given a random vector penalty d , the risk-averse constraint is to control a*
1184 *risk measure of the total expected discounted penalty. This leads to the following constrained MDP*
1185 *formulation:*

$$1186 \quad \min_{\pi} \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t c(s_t, a_t) \mid \pi, s_0 \sim \tau \right] \quad \text{s.t.} \quad \rho \left(\mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t d(s_t, a_t) \mid \pi, s_0 \sim \tau \right] \right) \leq D,$$

where ρ is a coherent risk measure, such as Conditional Value-at-Risk (CVaR)² Using Lagrangian relaxation, we can choose F to be a convex function of λ , i.e., $F(\lambda, P) = \langle \lambda, c \rangle + \ell(\rho(\langle \lambda, d \rangle) - D)$, where ℓ is the Lagrange multiplier.

Example 2 (Imitation Learning). During imitation learning, the agent learn through some demonstrations to behave similarly to an expert. One formulation is minimize the f -divergence between the occupancy measure of the current policy and the target occupancy measure:

$$\min_{\pi} D_f(\lambda^{\pi}, q) = \sum_{s,a} q(s,a) f\left(\frac{\lambda^{\pi}(s,a)}{q(s,a)}\right)$$

I EXAMPLES OF RISK ENVELOP

Example 3. [Conditional Value at Risk] First, Value-at-risk $\text{VaR}_{\beta}(X)$ is defined as the β -quantile of X , i.e., $\text{VaR}_{\beta}(X) := \inf\{t : \mathbb{P}(X \leq t) \geq \beta\}$, where the confidence level $\beta \in (0, 1)$. Assuming there is no probability atom at $\text{VaR}_{\beta}(X)$, CVaR at confidence level β is defined as the mean of the β -tail distribution of X , i.e., $\text{CVaR}_{\beta}(X) = \mathbb{E}[X | X \geq \text{VaR}_{\beta}(X)]$. The envelope set is

$$\mathcal{U}(\mu_N) = \{\xi \in \mathcal{Z}^* : \int_{\Theta} \xi(\theta) \mu_N(\theta) d\theta = 1, \xi(\theta) \in \left[0, \frac{1}{1-\beta}\right] \text{ a.s. } \theta \in \Theta\}$$

Example 4. (Mean-Upper-Semideviation of Order p). For $\mathcal{Z} := \mathcal{L}_p(\Theta, \mathcal{F}, \mu_N)$ and $\mathcal{Z}^* := \mathcal{L}_q(\Theta, \mathcal{F}, \mu_N)$, with $p \in [1, +\infty)$, $c \in [0, 1]$ and \mathcal{F} to be a σ -field on Θ , consider

$$\rho(Z) := \mathbb{E}[Z] + c \left(\mathbb{E}[[Z - \mathbb{E}[Z]]_+^p]\right)^{1/p},$$

where $[a]_+^p = \max\{0, a\}^p$. Then the envelope set is

$$\mathcal{U}(\mu_N) = \{\xi' \in \mathcal{Z}^* : \xi' = 1 + \xi - \mathbb{E}[\zeta], \|\xi\|_q \leq c, \xi \geq 0\}.$$

More examples can be found in Section 6.3.2(Shapiro et al., 2021).

J POLICY GRADIENT FOR MDP WITH CVAR RISK MEASURE : A SPECIAL CASE STUDY

Here we offer an example of gradient estimator with a common coherent risk measure Conditional Value at Risk(CVaR), the definition of which can be found in Example 3. For the considered CVaR risk functional, (Hong & Liu, 2009) shows that the gradient of the CVaR risk functional can be expressed as

$$\nabla \text{CVaR}_{\beta}(X(\alpha)) = \mathbb{E}[\nabla X(\alpha) | X(\alpha) \geq v_{\beta}(\alpha)]$$

where $v_{\beta} = v_{\beta}(\alpha) := \text{VaR}_{\beta}(X(\alpha))$ for a random parameterized variable $X(\alpha)$ satisfying Assumption J.1. Unless otherwise specified, the derivative is assumed to be taken w.r.t. α .

Assumption J.1. (Assumption 1, 2, 3 (Hong & Liu, 2009)) (i) There exists a random variable L with $\mathbb{E}(K) < \infty$ such that $|X(\alpha_2) - X(\alpha_1)| \leq K \|\alpha_2 - \alpha_1\|_2$ for all $\alpha_1, \alpha_2 \in W$, and $\nabla_{\alpha} X(\alpha)$ exists almost surely for all $\alpha \in W$.

(ii) VaR function $v_{\beta}(\alpha)$ is differentiable for any $\alpha \in W$.

(iii) For any $\alpha \in W$, $\mathbb{P}(X(\alpha) = v_{\beta}(\alpha)) = 0$.

Assumption J.1 (i) is commonly used in path-wise derivative estimation; (ii) shows that VaR function is locally Lipschitz; (iii) requires that there is no probability atom at $\text{VaR}(X)$ and implies that $\mathbb{P}(X(\alpha) \geq v_{\beta}(\alpha)) = 1 - \beta$.

Theorem 7. Suppose that Assumption J.1 holds. Then, for any $\alpha \in W$ and $\beta \in (0, 1)$, the policy gradient to the objective function in equation 3 is given by:

$$\begin{aligned} g(\alpha) &= \mathbb{E}_{\theta \sim \mu_N} [\nabla C(\alpha, \theta) | C(\alpha, \theta) \geq v_{\beta}(\alpha)] \\ &= \frac{1}{1-\beta} \mathbb{E}_{\theta \sim \mu_N} [\nabla C(\alpha, \theta) \mathbb{1}_{\{C(\alpha, \theta) \geq v_{\beta}\}}] \end{aligned} \quad (17)$$

where $\mathbb{1}_{\{\cdot\}}$ is the indicator function.

²CVaR(X) = $\mathbb{E}[X | X \geq v_{\beta}(X)]$, where $v_{\beta}(X)$ is a β -quantile of X , i.e. $\mathbb{P}(X \geq v_{\beta}(X)) = 1 - \beta$

If we apply Theorem 2 to CVaR, we will get the same result as Theorem 7. To compute the gradient $g(\alpha)$, we require the cumulative value $C(\alpha, \theta)$ of policy π_α and its gradient $\nabla C(\alpha, \theta)$, value-at-risk v_β , as well as the evaluation of the expectation taken w.r.t. the posterior distribution μ_N . Here we show how to use zeroth-order method instead of variational approach to estimate $\nabla_\alpha C(\alpha, \theta)$. Since there is no closed-form expression for the expectation, we estimate the gradient $g(\alpha)$ with samples $\{\theta^i\}_{i=1}^n$ generated from μ_N . We construct the gradient estimator as follows:

$$\widehat{g}(\alpha) = \frac{1}{n(1-\beta)} \sum_{i=1}^n \widehat{\nabla} C(\alpha, \theta^i) \mathbb{1}_{\{\widehat{C}(\alpha, \theta^i) \geq \widehat{v}_\beta\}}. \quad (18)$$

For a fixed α and θ^i , we first estimate the occupancy measure λ^i by making a truncation of horizon K in equation 1 with error

$$\|\widehat{\lambda}^i - \lambda^i\|_\infty \leq \epsilon_\lambda := \gamma^K / (1 - \gamma) \quad (19)$$

for some $K > 0$. The cumulative value with the truncated occupancy measure $\widehat{\lambda}^i$ is denoted by $\widehat{C}(\alpha, \theta^i) = F(\widehat{\lambda}, P_{\theta^i})$. The value-at-risk estimate is $\widehat{v}_\beta := \widehat{C}(\alpha, \theta)_{[n\beta]:n}$, where $\widehat{C}(\alpha, \theta)_{[n\beta]:n}$ is the $[n\beta]$ -th smallest quantity in $\{\widehat{C}(\alpha, \theta^i)\}_{i=1}^n$.

Here we adopt the Gaussian smoothing approach of estimating gradients from function evaluations (Nesterov & Spokoiny, 2017; Balasubramanian & Ghadimi, 2022). When there is no oracle to the first-order information or it is not efficient to calculate the gradient directly, Gaussian smoothing approach is a useful technique in zeroth-order method. Compared with finite difference method, Gaussian smoothing approach requires weaker smoothness condition of objective function. For a fixed α and θ^i , generate $\{u^{i,j}\}_{j=1}^{m_i}$, where $u^{i,j} \sim \mathcal{N}(0, I_d)$. Then $\widehat{\nabla} C$ can be constructed as:

$$\widehat{\nabla} C(\alpha, \theta^i) = \frac{1}{m_i} \sum_{j=1}^{m_i} \frac{\widehat{C}(\alpha + \nu u^{i,j}, \theta^i) - \widehat{C}(\alpha, \theta^i)}{\nu} u^{i,j} \quad (20)$$

where $\nu > 0$ is the smoothing parameter.

For ease of notation, let $\widehat{G}(\alpha)$ denote the sample estimate of $\rho_{\theta \sim \mu_N}(C(\alpha, \theta))$. We use the following gradient descent step in the t -th iteration:

$$\begin{aligned} \alpha_{t+1} &= \arg \min_{\alpha \in W} \widehat{G}(\alpha_t) + \langle \widehat{g}(\alpha_t), \alpha - \alpha_t \rangle + \frac{\eta_t}{2} \|\alpha - \alpha_t\|^2 \\ &= \text{Proj}_W \left(\alpha_t - \frac{1}{\eta_t} \widehat{g}(\alpha_t) \right) \end{aligned} \quad (21)$$

where η_t is the stepsize and $\text{Proj}_W(x) = \arg \min_{y \in W} \|y - x\|_2^2$ projects x into the parameter space W . We summarize the full algorithm in Algorithm 4.

J.1 CONVERGENCE ANALYSIS FOR CVAR RISK MEASURE

Here we only show the estimation error of the policy gradient. To get a finite-step convergence result similar to Theorem 4, we only need to substitute $\mathcal{O}(r^{-1/4})$ in Theorem 4 with $\mathcal{O}(R^{1/2})$, where $R^2 = \mathcal{O}\left(dn^{-1} + \epsilon_\lambda + \frac{d\epsilon_\lambda^2}{\nu^2} + \frac{d+\nu^2 d^3}{m}\right)$ is the bound for $\mathbb{E}\|g - \widehat{g}\|_2^2$ in Theorem 8.

Here we still adopt the Assumption 3.2 about the smoothness for the considered loss functions, which are commonly used in gradient descent analysis. The error bound for the zeroth-order estimation for ∇C is then shown in the next lemma.

Lemma J.1. *Suppose Assumption J.1 and Assumption 3.2 hold. Then we have for each $i \in [n]$*

$$\begin{aligned} \mathbb{E}\|\widehat{\nabla} C(\alpha, \theta_i) - \nabla C(\alpha, \theta_i)\|_2^2 &\leq \frac{8d}{\nu^2} L_{F,\infty}^2 \epsilon_\lambda^2 \\ &+ \frac{8(d+5)B^2}{m_i} + \frac{2\nu^2 L_{C,2}^2 (d+6)^3}{m_i}, \end{aligned} \quad (22)$$

where $L_{F,\infty}, L_{C,2}, B$ are constants in Assumption 3.2, ϵ_λ is the truncation error defined in equation 19, d is the dimension of the policy parameter α , m_i is the number of samples used to construct the zeroth-order estimator in equation 20.

Algorithm 4 BR-PG: Bayesian Risk Policy Gradient for CVaR

1296
1297
1298 **input:** initial α_0 , data $\zeta^{(N)}$ of size N , prior distribution $\mu_0(\theta)$, iteration number T , truncation
1299 horizon K ;
1300 calculate the posterior $\mu_N(\theta) = \frac{P_\theta(\zeta^{(N)})\mu_0(\theta)}{\int_{\theta'} P_{\theta'}(\zeta^{(N)})\mu_0(\theta')}$;
1301 **for** $t = 0$ to $T - 1$ **do**
1302 sample $\{\theta_t^i\}_{i=1}^n$ from $\mu_N(\theta)$;
1303 **for** $i = 1$ to n **do**
1304 calculate $\hat{\lambda}_t^i$ using the truncation of horizon K specified in equation 1;
1305 calculate $\hat{C}(\alpha_t, \theta_t^i) := F(\hat{\lambda}_t^i, P_{\theta_t^i})$;
1306 generate $\{u^{i,j}\}_{j=1}^{m_i}$, where $u^{i,j} \sim \mathcal{N}(0, I_d)$;
1307 calculate $\widehat{\nabla C}(\alpha_t, \theta_t^i)$ by equation 20;
1308 **end for**
1309 calculate $\hat{v}_\beta(\alpha_t) := \hat{C}(\alpha_t, \theta_t^i)_{[n\beta]:n}$.
1310 calculate $\hat{g}(\alpha_t)$ by equation 18;
1311 update α_{t+1} by equation 10.
1312 **end for**
1313 **output:** α_T .

1314
1315
1316 **Assumption J.2.** (Assumptions 4 and 5 in (Hong & Liu, 2009))

1317 (1) For all $\alpha \in W$, $C(\alpha, \theta)$ is a continuous random variable with a density function $f_{C,\alpha}(y)$.
1318 Furthermore, $f_{C,\alpha}(y)$ and $g_{C,\alpha}(y) := \mathbb{E}_\theta[\nabla C(\alpha, \theta) \mid C(\alpha, \theta) = y]$ are continuous at $y = v_\alpha$, and
1319 $f_{C,\alpha}(v_\alpha) > 0$.

1320
1321 (2) $\mathbb{E}_\theta [C(\alpha, \theta)^2] < \infty$ for all $\alpha \in W$.

1322 Now we are ready to show the error for our gradient estimator given in equation 18.

1323 **Theorem 8.** Suppose that Assumption J.1, Assumption 3.2 and Assumption J.2 hold. Also assume
1324 that the cumulative distribution function of $C(\alpha, \theta)$ w.r.t θ is ℓ_C - Lipschitz continuous for each
1325 $\alpha \in W$. Let $m_i = m \forall i \in [n]$. Then for each $\alpha \in W$,

$$1326 \mathbb{E} \|g - \hat{g}\|_2^2 \leq \mathcal{O} \left(dn^{-1} + \epsilon_\lambda + \frac{d\epsilon_\lambda^2}{\nu^2} + \frac{d + \nu^2 d^3}{m} \right),$$

1327 where n is the number of samples of θ .

1328
1329 *Proof.* First recall that the true gradient and our gradient estimator are $g =$
1330 $\frac{1}{1-\beta} \mathbb{E} [\nabla C(\alpha, \theta) \mathbb{1}_{\{C(\alpha, \theta) \geq v_\beta\}}]$ and $\hat{g} = \frac{1}{n(1-\beta)} \sum_{i=1}^n \widehat{\nabla C}(\alpha, \theta_i) \mathbb{1}_{\{\hat{C}(\alpha, \theta_i) \geq \hat{v}_\beta\}}$. Let

$$1331 \tilde{g} = \frac{1}{n(1-\beta)} \sum_{i=1}^n \nabla C(\alpha, \theta_i) \mathbb{1}_{\{C(\alpha, \theta_i) \geq \tilde{v}_\beta\}},$$

1332 and

$$1333 \hat{g}_1 = \frac{1}{n(1-\beta)} \sum_{i=1}^n \nabla C(\alpha, \theta_i) \mathbb{1}_{\{\hat{C}(\alpha, \theta_i) \geq \hat{v}_\beta\}},$$

1334 where $\tilde{v}_\beta := C(\alpha, \theta_i)_{[n\beta]:n}$. Then we have the decomposition $g - \hat{g} = (g - \tilde{g}) + (\tilde{g} - \hat{g}_1) + (\hat{g}_1 - \hat{g}) :=$
1335 $R_1 + R_2 + R_3$. For R_1 , it is the error in the estimation of expectation taken w.r.t. θ . Suppose that
1336 Assumption J.1 and Assumption J.2 hold, Theorem 4.2 from (Hong & Liu, 2009) shows that

$$1337 \| \mathbb{E} R_1 \|_2 = \| \mathbb{E} [\tilde{g}] - g \|_2 = o(n^{-1/2} d^{-1/2}).$$

1338 Notice that

$$1339 \|g - \tilde{g}\|_2^2 \leq 2\|g - \mathbb{E}\tilde{g}\|_2^2 + 2\|\mathbb{E}\tilde{g} - \tilde{g}\|_2^2.$$

1340 By Theorem 4.3 from (Hong & Liu, 2009), $Var(\tilde{g}) = \mathcal{O}(dn^{-1})$. Thus

$$1341 \mathbb{E} \|R_1\|_2^2 = \mathcal{O}(dn^{-1}). \tag{23}$$

For R_3 , it is the error in the estimation of $C(\alpha, \theta)$. By Lemma J.1, $\mathbb{E}[\|\widehat{\nabla C}(\alpha, \theta_i) - \nabla C(\alpha, \theta_i)\|_2^2] \leq \frac{8d}{\nu^2} L_{F,\infty}^2 \epsilon_\lambda^2 + \frac{8(d+5)B^2}{m_i} + \frac{2\nu^2 L_{C,2}^2 (d+6)^3}{m_i}$. If we choose all m_i to be the same m , then

$$\begin{aligned} \mathbb{E}[\|\widehat{g}_1 - \widehat{g}\|_2^2] &\leq \frac{1}{n(1-\beta)^2} \sum_{i=1}^n \|\widehat{\nabla C}(\alpha, \theta_i) - \nabla C(\alpha, \theta_i)\|_2^2 \\ &\leq \mathcal{O}\left(\frac{d\epsilon_\lambda^2}{\nu^2} + \frac{d+5}{m} + \frac{\nu^2(d+6)^3}{m}\right). \end{aligned}$$

Thus

$$\mathbb{E}[\|R_3\|_2^2] \leq \mathcal{O}\left(\frac{d\epsilon_\lambda^2}{\nu^2} + \frac{d+5}{m} + \frac{\nu^2(d+6)^3}{m}\right). \quad (24)$$

Now we consider R_2 . Define the event $A_i = \{C(\alpha, \theta_i) \geq \tilde{v}_\beta\}$, $\widehat{A}_i = \{\widehat{C}(\alpha, \theta_i) \geq \widehat{v}_\beta\}$ and $A_i \Delta \widehat{A}_i := (A_i \setminus \widehat{A}_i) \cup (\widehat{A}_i \setminus A_i)$. Then

$$\begin{aligned} \|R_2\|_2 &\leq \frac{1}{n(1-\beta)} \sum_{i=1}^n \|\nabla C(\alpha, \theta_i)\|_2 \cdot \mathbb{1}_{A_i \Delta \widehat{A}_i} \\ &\leq \frac{1}{n(1-\beta)} \sum_{i=1}^n B \mathbb{1}_{A_i \Delta \widehat{A}_i}, \end{aligned}$$

and

$$\begin{aligned} \|R_2\|_2^2 &\leq \frac{1}{n^2(1-\beta)^2} \left(\sum_{i=1}^n B \mathbb{1}_{A_i \Delta \widehat{A}_i}\right)^2 \\ &\leq \frac{1}{n(1-\beta)^2} B^2 \sum_{i=1}^n \mathbb{1}_{A_i \Delta \widehat{A}_i}. \end{aligned}$$

Notice that

$$\mathbb{P}(\mathbb{1}_{A_i \Delta \widehat{A}_i}) = \mathbb{P}(A_i \setminus \widehat{A}_i) + \mathbb{P}(\widehat{A}_i \setminus A_i).$$

As the estimation error of λ , i.e. $\|\widehat{\lambda} - \lambda\|_\infty$, is bounded by ϵ_λ and F is $L_{F,\infty}$ -Lipschitz continuous w.r.t $\|\cdot\|_\infty$, we have $|\widehat{C}(\alpha, \theta_i) - C(\alpha, \theta_i)| \leq L_{F,\infty} \epsilon_\lambda$. As a result, $|\tilde{v}_\beta - \widehat{v}_\beta| \leq L_{F,\infty} \epsilon_\lambda$. Notice that $\{C(\alpha, \theta_i) \geq \tilde{v}_\beta + 2L_{F,\infty} \epsilon_\lambda\} \subseteq \{\widehat{C}(\alpha, \theta_i) \geq \widehat{v}_\beta\} \subseteq \{C(\alpha, \theta_i) \geq \tilde{v}_\beta - 2L_{F,\infty} \epsilon_\lambda\}$. Then we have $\mathbb{P}(A_i \setminus \widehat{A}_i) + \mathbb{P}(\widehat{A}_i \setminus A_i) \leq 4\ell_C L_{F,\infty} \epsilon_\lambda$, by the assumption on the cumulative distribution function of C , and thus

$$\mathbb{E}\|R_2\|_2^2 \leq \frac{4}{(1-\beta)^2} B^2 \ell_C L_{F,\infty} \epsilon_\lambda = \mathcal{O}(\epsilon_\lambda). \quad (25)$$

Combining equation 23, equation 24 and equation 25, we have

$$\mathbb{E}\|g - \widehat{g}\|_2^2 \leq \mathcal{O}\left(dn^{-1} + \epsilon_\lambda + \frac{d\epsilon_\lambda^2}{\nu^2} + \frac{d + \nu^2 d^3}{m}\right).$$

□

Theorem 8 implies that the error of the gradient estimator can be reduced to arbitrarily small by increasing the sample size n, m or decreasing the truncation error ϵ_γ .

K IMPLEMENTING DETAILS

Frozen lake problem. Consider moving from the Start (S) to the Goal (G) on an 5×5 frozen lake with 6 holes (H). Then there are 18 ices (F) (involving Start). The agent may not move in the intended direction as the ice is slippery. The position is the row-column coordinate (i, j) with $i, j \in \{0, 1, 2, 3, 4\}$ and the state is the $5 * i + j$. The state space is $\{0, 1, \dots, 24\}$. The action set consists of moving in four directions. The unknown slippery probability is θ_s . Before reaching the goal and standing on the ice, the agent may move in the intended direction with unknown probability

1404 $1 - \theta_s$ and move in either perpendicular direction with probability $\theta_s/2$. When falling into the hole,
 1405 the agent may try to escape from the hole and move to the intended direction. Each time the agent
 1406 will succeed in escaping from the hole with unknown probability θ_e . After reaching the Goal, the
 1407 agent will always stay in the Goal whatever the action is. We set the cost to be 1 for each action on
 1408 ice before reaching goal. Also, stronger efforts may be made when it is harder to escape from the
 1409 hole. So we set the per-action cost in hole to be uniformly distributed between $[1, 1 + 2(1 - \theta_e)]$. We
 1410 aim to find a policy with the minimum general loss function. The data set consists of N historical
 1411 slippery movements and escapement trials.

1412 **Linear Loss.** For each of the considered formulations, we obtain the corresponding optimal policy
 1413 for the same data set and evaluate the actual performance of the obtained policy on the true system,
 1414 i.e. MDP with the true parameter θ^* . Specifically, we use the linear loss function, which corresponds
 1415 to the total discounted cost in a classical MDP problem. This is referred to as one replication, and
 1416 we repeat the experiments for 50 replications using different independent data sets. As the random
 1417 sampling of output policy in Algorithm 1 is for the purpose of proof, we just choose $\alpha_{t,N}$ as the output
 1418 for convenience in implementation. Results for the frozen lake problem are presented in Table 1, with
 1419 varying data size $N = 5$ and $N = 50$, slippery probability $\theta_s = 0.3$ and escape probability $\theta_e = 0.02$.
 1420 Note that we report the positive-sided variance, which corresponds to the second order moment of
 1421 the positive component of the difference between the actual loss and the expected loss. Intuitively, a
 1422 high positive-sided variance indicates more replications with higher costs than the average, which is
 1423 undesirable.

1424 **Episodic Case.** We consider the episodic setting where the data collection and policy update are
 1425 alternatively conducted. Similar with the previous case with fixed data size, we consider the mean
 1426 loss function with slippery probability $\theta_s = 0.3$, escape probability $\theta_e = 0.02$, and $5 \times 20, 10 \times 10,$
 1427 20×5 iterations in total. We repeat the experiments for 50 replications on different independent
 1428 data sets. Figure 1 shows the decrease of the loss function by different methods. As the random
 1429 sampling of output policy in Algorithm 2 is for the purpose of proof, we just choose $\alpha_{i+1,0} = \alpha_{i,t_i}$
 1430 for convenience in implementation.

1431 Results for the frozen lake problem with escape probability $\theta_e = 0.7$ can be found in Table 2 and
 1432 Table 3.

1433 Table 2: Results for frozen lake problem. Expected loss and positive-sided variance at different risk
 1434 levels α are reported for different algorithms. Standard errors are reported in parentheses. Escape
 1435 probability $\theta_e = 0.7$ and number of data points is $N = 5$.

Approach	loss function: mean	
	expected loss	positive-sided variance
BR-PG ($\beta = 0$)	10.322 (0.0182)	0.0153
BR-PG ($\beta = 0.5$)	10.520(0.105)	0.502
BR-PG ($\beta = 0.9$)	11.718 (0.357)	4.982
Empirical	11.667 (0.0687)	0.156
DRQL (radius=0.05)	11.223(0.185)	1.283
DRQL (radius=1)	20.751(1.438)	69.514
DRQL (radius=20)	23.181(1.396)	57.495

1447
 1448 Figure 3 shows the map of the frozen lake problem with 1 Start(S), 1 Goal(G), 6 holes(H) and
 1449 remaining frozen(F) parts. We design such a map so that the agent has to avoid falling in the hole
 1450 when the escape probability is very small and cross the hole when the escape probability is high.
 1451 Detailed parameters are set as follows. The true slippery probability is 0.3. The iteration number for
 1452 gradient descent is 100, the stepsize is 0.5, and the sample number in each iteration is $r = 30$. we set
 1453 the discount factor to be $\gamma = 0.97$, the truncation horizon for occupancy measure to be $K = 130$.
 1454 equation 20.

1455 For the "mean" loss function, we use the maximum likelihood estimator (MLE) of θ as the empirical
 1456 measure to be compared with BR-PG. Also, we use the distributionally robust Q-learning (DRQL)(Liu
 1457 et al., 2022) with different radius for the KL divergence ball as another benchmark. We also use the
 MLE of θ as the parameter for the center of the KL divergence ball in DRQL with different radius.

Table 3: Results for frozen lake problem. Expected loss and positive-sided variance at different risk levels α are reported for different algorithms. Standard errors are reported in parentheses. Escape probability $\theta_e = 0.7$ and number of data points is $N = 50$.

Approach	loss function: mean	
	expected loss	positive-sided variance
BR-PG ($\beta = 0$)	10.271 (0.00227)	0.000197
BR-PG ($\beta = 0.5$)	10.256 (0.00211)	0.000188
BR-PG ($\beta = 0.9$)	10.230(0.00294)	0.000398
Empirical	11.316 (0.0235)	0.017
DRQL (radius=0.05)	10.888(0.171)	1.235
DRQL (radius=1)	20.990(1.324)	56.027
DRQL (radius=20)	23.500(1.282)	51.915

For BR-PG, the sample number from posterior in each iteration is 30, the total iteration number is 100, the step size of SGD is chosen to be 1, and the prior distributions are chosen to be Beta(1, 1) for two parameters. We show the histogram of total cost over 50 replications for all methods in Figure 4 with the risk level 0.8 for CVaR over replications, which visualize the measures of dispersion.

Mimicking a policy. Here we consider a different problem of mimicking an expert policy still using Frozen Lake environment. Given an expert policy, we have access to the state distribution of the expert policy under the true environment, which is denoted by a nonnegative function J satisfying $\sum_{s \in \mathcal{S}} J(s) = 1$. The loss function we want to minimize is defined as the KL divergence between state occupancy measure under the current policy and the expert state distribution $F(\lambda) = \text{KL}((1 - \gamma) \sum_{a \in \mathcal{A}} \lambda_a || J) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} (1 - \gamma) \lambda_{sa} \log \left(\frac{\sum_{a \in \mathcal{A}} (1 - \gamma) \lambda_{sa}}{J(s)} \right)$. We compare the BR-PG algorithm with CVaR risk measure under different risk levels $\beta = 0, 0.5, 0.9$, respectively, with the benchmark empirical approach using the MLE estimator for the parameter as before. Figure 2 shows the decrease of the loss function by different methods. It should be noticed that DRQL can only be applied to the "mean" loss function, thus we don't use it as a benchmark. The performance of the 50 replications is shown in figure 5, where the shown results start from the 30-th iteration.



Figure 3: Map of frozen lake problem

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

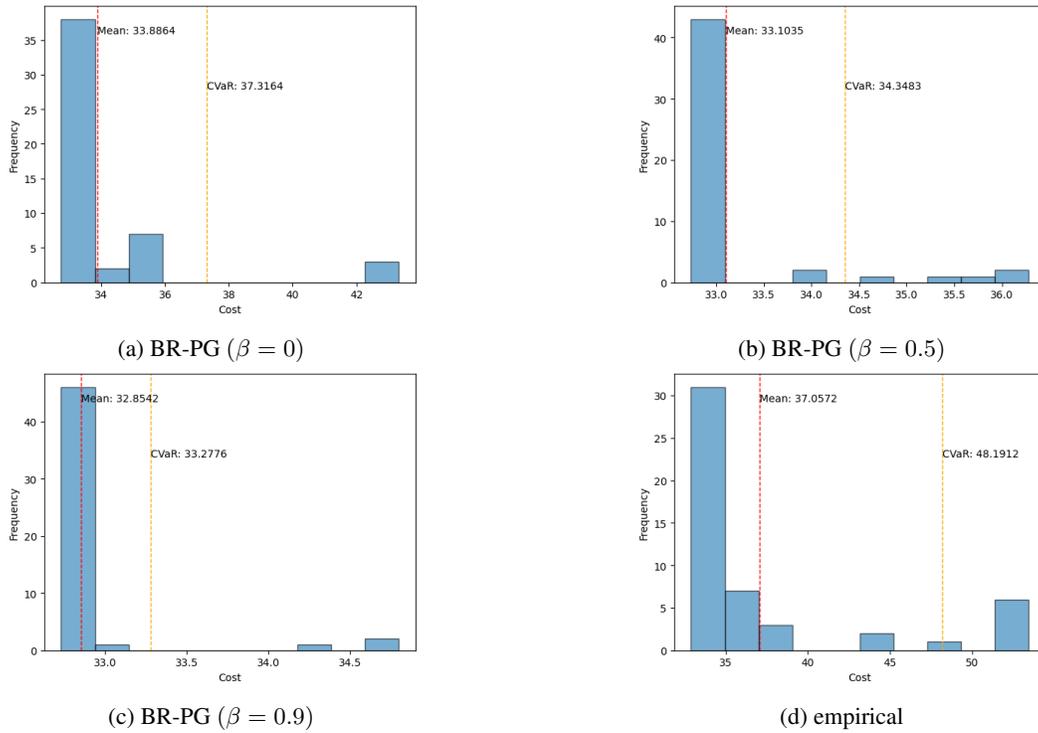


Figure 4: Result for utility function "mean" with data size $N = 5$ and escape probability $\theta_e = 0.02$

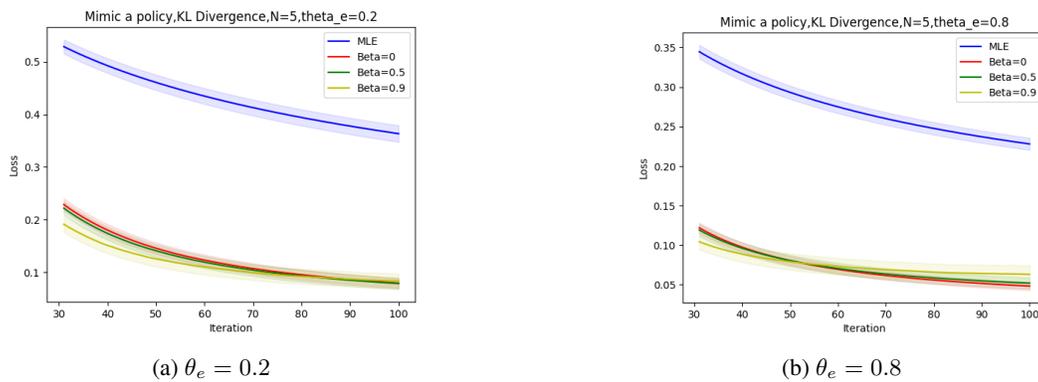


Figure 5: Results for utility function "KL divergence" with data size $N = 5$ and escape probability $\theta_e = 0.2$ and $\theta_e = 0.8$