

# NEAR-OPTIMAL POLICY IDENTIFICATION IN ROBUST CONSTRAINED MARKOV DECISION PROCESSES VIA EPIGRAPH FORM

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Designing a safe policy for uncertain environments is crucial in real-world control systems. However, this challenge remains inadequately addressed within the Markov decision process (MDP) framework. This paper presents the first algorithm guaranteed to identify a near-optimal policy in a robust constrained MDP (RCMDP), where an optimal policy minimizes cumulative cost while satisfying constraints in the worst-case scenario across a set of environments. We first prove that the conventional policy gradient approach to the Lagrangian max-min formulation can become trapped in suboptimal solutions. This occurs when its inner minimization encounters a sum of conflicting gradients from the objective and constraint functions. To address this, we leverage the epigraph form of the RCMDP problem, which resolves the conflict by selecting a single gradient from either the objective or the constraints. Building on the epigraph form, we propose a binary search algorithm with a policy gradient subroutine and prove that it identifies an  $\varepsilon$ -optimal policy in an RCMDP with  $\tilde{O}(\varepsilon^{-4})$  robust policy evaluations.

## 1 INTRODUCTION

In real-world decision-making, it is crucial to design policies that satisfy safety constraints even in uncertain environments. For example, self-driving cars must drive efficiently while maintaining a safe distance from obstacles, regardless of environmental uncertainties such as road conditions, weather, or the state of the vehicle’s components. Traditionally, within the Markov decision process (MDP) framework, constraint satisfaction and environmental uncertainty have been addressed separately—through constrained MDP (CMDP; e.g., Altman (1999)), which aims to minimize cumulative cost while satisfying constraints, and robust MDP (RMDP; e.g., Iyengar (2005)), which aims to minimize the worst-case cumulative cost in an uncertainty set of possible environments. However, in practice, both robustness and constraint satisfaction are important. The recent robust constrained MDP (RCMDP) framework addresses this dual need by aiming to minimize the worst-case cost while robustly satisfying the constraints. Despite the significant theoretical progress made in CMDPs and RMDPs (see Appendix A), theoretical results on RCMDPs are currently scarce. Even in the tabular setting, where the state and action spaces are finite, there exists no algorithm with guarantees to find a near-optimal policy in an RCMDP.

The difficulty of RCMDPs arises from the challenging optimization process, which simultaneously considers robustness and constraints. The dynamic programming (DP) approach, popular in unconstrained RMDPs, is unsuitable for constrained settings where Bellman’s principle of optimality can be violated (Haviv, 1996; Bellman et al., 1957). Similarly, the linear programming (LP) approach, commonly used for CMDPs, is inadequate due to the nonconvexity of the robust formulation (Iyengar, 2005; Grand-Clément & Petrik, 2024). Consequently, the policy gradient method with the Lagrangian formulation has been studied as the primary remaining option (Russel et al., 2020; Wang et al., 2022). The Lagrangian formulation approximates the RCMDP problem  $\min_{\pi} \{f(\pi) \mid h(\pi) \leq 0\}$  by  $\max_{\lambda \geq 0} \min_{\pi} f(\pi) + \lambda h(\pi)$ , where  $f(\pi)$  and  $h(\pi)$  represent the worst-case cumulative cost—called the (cost) return<sup>1</sup>—and the worst-case constraint violation of

<sup>1</sup>We commonly use the term *return* to refer specifically to the objective cost return. When discussing a return value in the context of RCMDP’s constraints, we refer to it as the *constraint return*.

Approach	MDP	CMDP	RMDP	RCMDP
Dynamic Programming	✓ Bellman et al. (1957)	✗	✓ Iyengar (2005)	✗
Linear Programming	✓ Denardo (1970)	✓ Altman (1999)	✗	✗
Lagrangian + PG	✓ Agarwal et al. (2021)	✓ Ding et al. (2020)	✓ Wang et al. (2023)	✗
<b>Epigraph + PG (Ours)</b>	✓	✓	✓	✓

Table 1: Summary of approaches and the problem settings. “PG” denotes Policy Gradient. Each cell displays a “✓” indicating the presence of an algorithm with this approach that guarantees yielding an  $\varepsilon$ -optimal policy. Representative works supporting each “✓” are listed below it. Conversely, “✗” denotes settings where the approach either isn’t suitable or lacks performance guarantees.

a policy  $\pi$ , respectively. There have been a few attempts to provide theoretical guarantees for the Lagrangian approach (Wang et al., 2022; Zhang et al., 2024); however, no existing studies offer rigorous and satisfactory guarantees that the max-min problem yields the same solution as the original RCMDP problem. As a result, existing Lagrangian-based algorithms lack theoretical performance guarantees. This leaves us with a fundamental question:

*How can we identify a near-optimal policy in a tabular RCMDP?*

We address this question by presenting three key contributions, which are summarized as follows:

**Gradient conflict in the Lagrangian formulation (Section 3).** We first show that solving the Lagrangian formulation is inherently difficult, even when its max-min problem can yield an optimal policy. Given the limitations of DP and LP approaches as discussed, the policy gradient method might seem like a viable alternative to solve the max-min. However, our Theorem 1 reveals that policy gradient methods can get trapped in a local minimum during the inner minimization of the Lagrangian formulation. This occurs when the gradients,  $\nabla f(\pi)$  and  $\nabla h(\pi)$ , conflict with each other, causing their sum  $\nabla f(\pi) + \lambda \nabla h(\pi)$  to cancel out, even when the policy  $\pi$  is not optimal. Consequently, the Lagrangian approach for RCMDPs may not reliably lead to a near-optimal policy.

**Epigraph form of RCMDP (Section 4).** We then demonstrate that the *epigraph form*, commonly used in constrained optimization literature (Boyd & Vandenberghe, 2004; Beyer & Sendhoff, 2007; Rahimian & Mehrotra, 2019), entirely circumvents the challenges associated with the Lagrangian formulation. The epigraph form transforms the RCMDP problem into  $\min_y \{y \mid \min_{\pi} \max\{f(\pi) - y, h(\pi)\} \leq 0\}$ , introducing an auxiliary minimization problem of  $\min_{\pi} \max\{f(\pi) - y, h(\pi)\}$  and minimizing its threshold variable  $y$ . Unlike the Lagrangian approach, which necessitates **summing**  $\nabla f(\pi)$  and  $\nabla h(\pi)$ , policy gradient methods for the auxiliary problem update the policy by **selecting either**  $\nabla f(\pi)$  **or**  $\nabla h(\pi)$ , thanks to the maximum operator in the problem. As a result, the epigraph form avoids the problem of conflicting gradient sums, preventing policy gradient methods from getting stuck in suboptimal minima (Theorem 4).

**A new RCMDP algorithm (Section 5).** Finally, we propose a tabular RCMDP algorithm called **Epigraph Robust Constrained Policy Gradient Search (EpiRC-PGS)**, pronounced as “Epic-P-G-S”). The algorithm employs a double-loop structure: the inner loop verifies the feasibility of the threshold variable  $y$  by performing policy gradients on the auxiliary problem, while the outer loop employs binary search to determine the minimal feasible  $y$ . EpiRC-PGS is guaranteed to find an  $\varepsilon$ -optimal policy<sup>2</sup> with  $\tilde{O}(\varepsilon^{-4})$  robust policy evaluations (Corollary 1), where  $\tilde{O}(\cdot)$  represents the conventional big-O notation excluding polylogarithmic terms. Since RCMDP generalizes plain MDP, CMDP, and RMDP, our EpiRC-PGS is applicable to all these types of MDPs, ensuring a near-optimal policy for each. Table 1 compares existing approaches in various MDP settings. Due to the page limitation, more related work is provided in Appendix A. We discuss limitations and potential future directions in Section 7.

<sup>2</sup>The definition of an  $\varepsilon$ -optimal policy is provided in Definition 1

## 2 PRELIMINARY

We use the shorthand  $\mathbb{R}_+ := [0, \infty)$ . The set of probability distributions over  $\mathcal{S}$  is denoted by  $\mathcal{P}(\mathcal{S})$ . For two integers  $a \leq b$ , we define  $\llbracket a, b \rrbracket := \{a, \dots, b\}$ . If  $a > b$ ,  $\llbracket a, b \rrbracket := \emptyset$ . For a vector  $x \in \mathbb{R}^N$ , its  $n$ -th element is denoted by  $x_n$  or  $x(n)$ , and we use the convention that  $\|x\|_2 = \sqrt{\sum_i x_i^2}$  and  $\|x\|_\infty = \max_i |x_i|$ . For two vectors  $x, y \in \mathbb{R}^N$ , we denote  $\langle x, y \rangle = \sum_i x_i y_i$ . We let  $\mathbf{0} := (0, \dots, 0)^\top$  and  $\mathbf{1} := (1, \dots, 1)^\top$ , with their dimensions being clear from the context. All scalar operations and inequalities should be understood point-wise when applied to vectors and functions. Given a finite set  $\mathcal{S}$ , we often treat a function  $f : \mathcal{S} \rightarrow \mathbb{R}$  as a vector  $f \in \mathbb{R}^{\mathcal{S}}$ . Both notations,  $f : \mathcal{S} \rightarrow \mathbb{R}$  and  $f \in \mathbb{R}^{\mathcal{S}}$ , are used depending on notational convenience. Finally,  $\partial f(x)$  and  $\nabla f(x)$  denote the (Fréchet) subgradient and gradient of  $f : \mathcal{X} \rightarrow \mathbb{R}$  at a point  $x$ , respectively. Their formal definitions are deferred to Definition 2.

### 2.1 CONSTRAINED MARKOV DECISION PROCESS

Let  $N \in \mathbb{Z}_{>0}$  be the number of constraints. An infinite-horizon discounted constrained MDP (CMDP) is defined as a tuple  $(\mathcal{S}, \mathcal{A}, \gamma, P, \mathcal{C}, b, \mu)$ , where  $\mathcal{S}$  denotes the finite state space with size  $S$ ,  $\mathcal{A}$  denotes the finite action space with size  $A$ ,  $\gamma \in (0, 1)$  denotes the discount factor, and  $\mu \in \mathcal{P}(\mathcal{S})$  denotes the initial state distribution. For notational brevity, let  $H := (1 - \gamma)^{-1}$  denotes the effective horizon. Further,  $b := (b_1, \dots, b_N) \in [0, H]^N$  denotes the constraint threshold vector, where  $b_n$  is the threshold scalar for the  $n$ -th constraint,  $\mathcal{C} := \{c_n\}_{c_n \in [0, N]}$  denotes the set of cost functions, where  $c_n : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  denotes the  $n$ -th cost function and  $c_n(s, a)$  denotes the  $n$ -th cost when taking an action  $a$  at a state  $s$ .  $c_0$  is for the objective to optimize and  $\{c_1, \dots, c_N\}$  are for the constraints.  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$  denotes the transition probability kernel, which can be interpreted as the environment with which the agent interacts.  $P(s' | s, a)$  denotes the state transition probability to a new state  $s'$  from a state  $s$  when taking an action  $a$ .

### 2.2 POLICY AND VALUE FUNCTIONS

A (Markovian stationary) policy is defined as  $\pi \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  such that  $\pi(s, \cdot) \in \mathcal{P}(\mathcal{A})$  for any  $s \in \mathcal{S}$ .  $\pi(s, a)$  denotes the probability of taking an action  $a$  at state  $s$ . The set of all the policies is denoted as  $\Pi$ , which corresponds to the *direct parameterization* policy class presented in Agarwal et al. (2021). Although non-Markovian policies can yield better performance in general RMDP problems (Wiesemann et al., 2013), for simplicity, we focus on Markovian stationary policies in this paper. With an abuse of notation, for two functions  $\pi, g \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ , we denote  $\langle \pi, g \rangle = \sum_{s, a \in \mathcal{S} \times \mathcal{A}} \pi(s, a) g(s, a)$ .

For a policy  $\pi$  and transition kernel  $P$ , let  $d_P^\pi : \mathcal{S} \rightarrow \mathbb{R}_+$  denote the occupancy measure of  $\pi$  under  $P$ .  $d_P^\pi(s)$  represents the expected discounted number of times  $\pi$  visits state  $s$  under  $P$ , such that  $d_P^\pi(s) = (1 - \gamma) \mathbb{E}[\sum_{h=0}^{\infty} \gamma^h \mathbb{1}\{s_h = s\} \mid s_0 \sim \mu, \pi, P]$ . Here, the notation means that the expectation is taken over all possible trajectories, where  $a_h \sim \pi(s_h, \cdot)$  and  $s_{h+1} \sim P(\cdot \mid s_h, a_h)$ .

For a  $\pi \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  and a cost  $c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ , let  $Q_{c, P}^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  be the action-value function such that<sup>3</sup>

$$Q_{c, P}^\pi(s, a) = c(s, a) + \gamma \sum_{s' \in \mathcal{S}} P(s' \mid s, a) \sum_{a' \in \mathcal{A}} \pi(s', a') Q_{c, P}^\pi(s', a') \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

Let  $V_{c, P}^\pi : \mathcal{S} \rightarrow \mathbb{R}$  be the state-value function such that  $V_{c, P}^\pi(s) = \sum_{a \in \mathcal{A}} \pi(s, a) Q_{c, P}^\pi(s, a)$  for any  $s \in \mathcal{S}$ . If  $\pi \in \Pi$ ,  $V_{c, P}^\pi(s)$  represents the expected cumulative cost of  $\pi$  under  $P$  with an initial state  $s$ . We denote the (cost) return function as  $J_{c, P}(\pi) := \sum_{s, a \in \mathcal{S} \times \mathcal{A}} \mu(s) V_{c, P}^\pi(s)$ .

**Policy gradient method.** For a problem  $\min_{\pi \in \Pi} f(\pi)$  where  $f : \Pi \rightarrow \mathbb{R}$  is differentiable at  $\pi \in \Pi$ , policy gradient methods with direct parameterization update  $\pi$  to a new policy  $\pi'$  as follows:

$$\pi' := \text{Proj}_\Pi(\pi - \alpha \nabla f(\pi)) \stackrel{(a)}{=} \arg \min_{\pi' \in \Pi} \langle \nabla f(\pi), \pi' - \pi \rangle + \frac{1}{2\alpha} \|\pi' - \pi\|_2^2, \quad (1)$$

where  $\alpha > 0$  is the learning rate and  $\text{Proj}_\Pi$  denotes the Euclidean projection operator onto  $\Pi$ . The equality (a) is a standard result (see, e.g., Parikh et al. (2014)). The following lemma provides the gradient of  $J_{c, P}(\pi)$  for the direct parameterization policy class  $\Pi$  (e.g., Agarwal et al. (2021)).

<sup>3</sup>The domain of  $Q_{c, P}^\pi$  is not restricted to  $\Pi$  to ensure well-defined policy gradients over  $\pi \in \Pi$ .

**Lemma 1** (Policy gradient theorem). *For any  $\pi \in \Pi$ , transition kernel  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ , and cost  $c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ , the gradient is given by  $(\nabla J_{c,P}(\pi))(s, a) = Hd_P^\pi(s) Q_{c,P}^\pi(s, a) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$ .*

### 2.3 ROBUST CONSTRAINED MARKOV DECISION PROCESS

An infinite-horizon discounted robust constrained MDP (RCMDP) is defined as a tuple  $(\mathcal{S}, \mathcal{A}, \gamma, \mathcal{U}, \mathcal{C}, b, \mu)$ , where  $\mathcal{U}$  is a compact set of transition kernels, called the uncertainty set, which can be either finite or infinite. The infinite uncertainty set typically requires some structural assumptions. A common structure is the  $(s, a)$ -rectangular set (Iyengar, 2005; Nilim & El Ghaoui, 2005), defined as  $\mathcal{U} = \times_{s,a} \mathcal{U}_{s,a}$ , where  $\mathcal{U}_{s,a} \subseteq \mathcal{P}(\mathcal{S})$  and  $\times_{s,a}$  denotes a Cartesian product over  $\mathcal{S} \times \mathcal{A}$ . We remark that our work is **not** limited to any specific structural assumption, but rather considers a general, tractable uncertainty set (see Assumptions 1, 2 and 3). When  $N = 0$ , an RCMDP reduces to an RMDP. When  $\mathcal{U} = \{P\}$ , an RCMDP becomes a CMDP.

For a cost  $c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ , let  $J_{c,\mathcal{U}}(\pi) := \max_{P \in \mathcal{U}} J_{c,P}(\pi)$  denote the worst-case (cost) return function, which represents the return of  $\pi$  under the most adversarial environment within  $\mathcal{U}$ . The goal of an RCMDP is to find a solution to the following constrained optimization problem:

$$\text{(RCMDP)} \quad J^* := \min_{\pi \in \Pi} J_{c_0,\mathcal{U}}(\pi) \quad \text{such that} \quad J_{c_n,\mathcal{U}}(\pi) \leq b_n \quad \forall n \in \llbracket 1, N \rrbracket. \quad (2)$$

Let  $\Pi_F := \{\pi \in \Pi \mid \max_{c_n \in \llbracket 1, N \rrbracket} J_{c_n,\mathcal{U}}(\pi) - b_n \leq 0\}$  be the set of all the feasible policies. We assume that  $\Pi_F$  is non-empty. An optimal policy  $\pi^* \in \Pi_F$  is a solution to Equation (2).

**Definition 1.**  $\pi \in \Pi$  is  $\varepsilon$ -optimal<sup>4</sup> if  $J_{c_0,\mathcal{U}}(\pi) - J^* \leq \varepsilon$  and  $\max_{c_n \in \llbracket 1, N \rrbracket} J_{c_n,\mathcal{U}}(\pi) - b_n \leq \varepsilon$ .

## 3 CHALLENGES OF LAGRANGIAN FORMULATION

To motivate our formulations and algorithms presented in subsequent sections, this section illustrates the limitations of using the conventional Lagrangian approach for RCMDPs. By introducing Lagrangian multipliers  $\lambda := (\lambda_1, \dots, \lambda_N) \in \mathbb{R}_+^N$ , Equation (2) is equivalent to

$$J^* = \min_{\pi \in \Pi} \max_{\lambda \in \mathbb{R}_+^N} J_{c_0,\mathcal{U}}(\pi) + \sum_{n=1}^N \lambda_n (J_{c_n,\mathcal{U}}(\pi) - b_n).$$

As this  $\min_{\pi}$  is hard to solve due to the inner maximization, Russel et al. (2020); Mankowitz et al. (2020); Wang et al. (2022) swap the min-max and consider the following Lagrangian formulation:

$$\text{(Lagrange)} \quad L^* := \max_{\lambda \in \mathbb{R}_+^N} \min_{\pi \in \Pi} L_\lambda(\pi) \quad \text{where} \quad L_\lambda(\pi) := J_{c_0,\mathcal{U}}(\pi) + \sum_{n=1}^N \lambda_n (J_{c_n,\mathcal{U}}(\pi) - b_n). \quad (3)$$

Let  $\lambda^*$  be a solution to Equation (3). The Lagrangian approach aims to solve Equation (3) by expecting  $\pi^* \in \arg \min_{\pi \in \Pi} L_{\lambda^*}(\pi)$ . However, this expectation may not hold, as swapping the min-max is not necessarily equivalent to the original min-max problem (Boyd & Vandenberghe, 2004). Therefore, to guarantee the performance of the Lagrange approach, the two questions must be addressed:

- (i) *Can we ensure  $\pi^* \in \arg \min_{\pi \in \Pi} L_{\lambda^*}(\pi)$ ?*    (ii) *If so, is it tractable to solve  $\min_{\pi \in \Pi} L_\lambda(\pi)$ ?*

However, answering these questions affirmatively is challenging due to the following issues:

(i)  **$\pi^*$  solution challenge.** To ensure that  $\pi^* \in \arg \min_{\pi \in \Pi} L_{\lambda^*}(\pi)$ , the standard approach is to establish the strong duality, i.e.,  $J^* = L^*$  (Boyd & Vandenberghe, 2004). In the CMDP setting, where  $\mathcal{U} = \{P\}$ , strong duality has been proven to hold (Altman, 1999; Paternain et al., 2019; 2022). However, proving strong duality for RCMDPs is highly non-trivial compared to CMDPs.

When  $\mathcal{U} = \{P\}$ , a typical proof strategy is to combine the sum of returns  $J_{c_0,P}, \dots, J_{c_N,P}$  in  $L_\lambda(\pi)$  into a single return function. For example, when  $N = 1$ , we have  $L_\lambda(\pi) = J_{c',P}(\pi)$ , where

<sup>4</sup>Strict constraint satisfaction is straightforward by using a slightly stricter threshold  $b' := b - \varepsilon$ .

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269

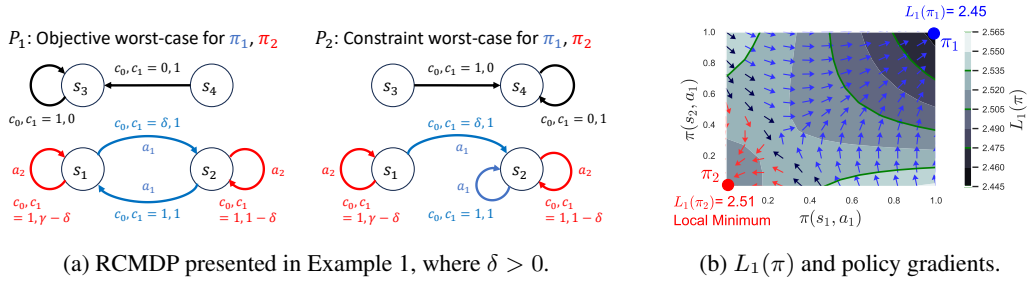


Figure 1: **(a)**: An RCMDP example illustrating the **gradient conflict challenge**. Action labels are omitted when transitions are action-independent. **(b)**: Policy gradients in the example with  $(\gamma, \delta, b_1) = (0.4, 0.09, 0)$ . Arrows represent the gradient to decrease  $L_1(\pi)$ .  $\pi_2$  attracts policy gradients but is a local minimum since  $L_1(\pi_2) > L_1(\pi_1)$ , where  $\pi_1(\cdot, a_1) = 1$  and  $\pi_2(\cdot, a_2) = 1$ .

$c' := c_0 + \lambda(c_1 - b_1/H)$ . Since  $J_{c',P}(\pi)$  is linearly represented as  $J_{c',P}(\pi) = \langle c', d_P^\pi \rangle^5$ , it is easy to show  $\min_{\pi} \max_{\lambda} J_{c',P}(\pi) = \max_{\lambda} \min_{\pi} J_{c',P}(\pi)$  by Sion's minimax theorem (Sion, 1958).

However, applying this return-combining strategy to RCMDPs is difficult. For  $n \in \llbracket 0, N \rrbracket$ , let  $P_n \in \arg \max_{P \in \mathcal{U}} J_{c_n, P}(\pi)$ . When  $|\mathcal{U}| \neq 1$ ,  $L_\lambda(\pi)$  has a sum of returns  $J_{c_0, P_0}(\pi) \dots J_{c_N, P_N}(\pi)$ , where  $P_0 \neq \dots \neq P_N$  may differ. Thus,  $L_\lambda(\pi)$  is no longer a form of  $\langle c', d_P^\pi \rangle$  for some single environment  $P$ , hindering the use of well-established duality proof techniques in CMDP literature.

(ii) **Gradient conflict challenge**. Unfortunately, even if strong duality holds and  $\pi^*$  can be found by  $\pi^* \in \arg \min_{\pi \in \Pi} L_{\lambda^*}(\pi)$ , solving this minimization remains challenging. Since the sum of robust returns in  $L_{\lambda^*}(\pi)$  excludes the use of DP and convex-optimization approaches (Iyengar, 2005; Altman, 1999; Grand-Clément & Petrik, 2024), the policy gradient method such as Equation (1) is the primary remaining option to solve  $\min_{\pi \in \Pi} L_\lambda(\pi)$ .

In CMDP setting when  $|\mathcal{U}| = 1$ , the problem  $\min_{\pi \in \Pi} L_\lambda(\pi)$  reduces to solving a standard MDP and thus the policy gradient method is ensured to solve  $\min_{\pi \in \Pi} L_\lambda(\pi)$  (Agarwal et al., 2021). However, the following Theorem 1 shows that even when  $|\mathcal{U}| = 2$ , RCMDP can trap the policy gradient in a local minimum that does not solve  $\min_{\pi \in \Pi} L_\lambda(\pi)$ :

**Theorem 1.** For any  $\gamma \in (0, 1)$ , there exist a  $\bar{\lambda} > 0$ , a policy  $\bar{\pi} \in \Pi$  and an RCMDP with  $\mu > 0$  satisfying the following condition: There exists a positive constant  $R > 0$  such that, for any  $b_1 \in \mathbb{R}$ ,

$$L_{\bar{\lambda}}(\bar{\pi}) < L_{\bar{\lambda}}(\pi) \quad \forall \pi \in \{\pi \in \Pi \mid \|\pi - \bar{\pi}\|_2 \leq R, \pi \neq \bar{\pi}\} \quad \text{but} \quad L_{\bar{\lambda}}(\bar{\pi}) \geq \min_{\pi \in \Pi} L_{\bar{\lambda}}(\pi) + \frac{3\gamma H}{16}. \quad (4)$$

Moreover, there exists a  $b_1 \in (0, H)$  where  $\bar{\lambda}$  satisfies  $\bar{\lambda} \in \arg \max_{\lambda \in \mathbb{R}_+^N} \min_{\pi \in \Pi} L_\lambda(\pi)$ .

The detailed proof is deferred to Appendix H. Essentially, the proof constructs a simple RCMDP where **the policy gradients for the objective and the constraint are in conflict**.

**Example 1.** Consider the RCMDP with  $\mathcal{U} = \{P_1, P_2\}$  presented in Figure 1a with  $\delta = 0$  and set  $\lambda = 1$  for simplicity. Let  $\pi_1$  and  $\pi_2$  be policies that select  $a_1$  and  $a_2$  for all states, respectively. For both policies, the objective worst-case is  $P_1$  and the constraint worst-case is  $P_2$  (see Appendix H). Hence, switching from policy  $\pi_2$  to taking action  $a_1$  decreases the objective return in  $P_1$  but increases the constraint return in  $P_2$ . This conflict causes the gradients of  $\pi_2$  for the objective ( $\nabla J_{c_0, P_1}(\pi_2)$ ) and for the constraint ( $\nabla J_{c_1, P_2}(\pi_2)$ ) to sum to a constant vector, i.e.,

$$(\nabla L_1(\pi_2))(s, \cdot) = (\nabla J_{c_0, P_1}(\pi_2) + \nabla J_{c_1, P_2}(\pi_2))(s, \cdot) = \text{constant} \cdot \mathbf{1} \quad \forall s \in \mathcal{S}, \quad (5)$$

showing that  $\pi_2$  is a stationary point. However,  $\pi_2$  cannot solve  $\min_{\pi \in \Pi} L_1(\pi)$  because  $\pi_1$  would clearly result in a smaller  $L_1(\pi)$ . This stationary point becomes a strict local minimum when  $\delta > 0$ , where  $\pi_2$  slightly prefers  $a_2$  over  $a_1$  (see Appendix H for details).

Figure 1b computationally illustrates this negative result by plotting the landscape of  $L_1(\pi)$  in the RCMDP example across all possible policies for  $(\gamma, \delta) = (0.4, 0.09)$ . We set  $b_1 = 0$  as it does not influence the landscape of  $L_1(\pi)$ . In this example,  $\pi_2$  becomes a local minimum that attracts the policy gradient but fails to solve  $\min_{\pi \in \Pi} L_1(\pi)$ , as  $\pi_1$  achieves  $L_\lambda(\pi_1) < L_\lambda(\pi_2)$ .

<sup>5</sup>With an abuse of notation, here we denote  $d_P^\pi(s, a) = \mathbb{E}[\sum_{h=0}^{\infty} \gamma^h \mathbb{1}\{s_h = s, a_h = a\} \mid s_0 \sim \mu, \pi, P]$ .

## 4 EPIGRAPH FORM OF RCMDP

This section introduces the epigraph form of RCMDP, which overcomes the challenges discussed in Section 3. For any constrained optimization problem of the form  $\min_x \{f(x) \mid h(x) \leq 0\}$  with  $x \in \mathbb{R}^n$  and  $f, h : \mathbb{R}^n \rightarrow \mathbb{R}$ , its epigraph form is defined as:

$$\min_{x,y} y \text{ such that } f(x) \leq y \text{ and } h(x) \leq 0 \quad (6)$$

with variables  $x \in \mathbb{R}^n$  and  $y \in \mathbb{R}$ . It is well-known that  $(x, y)$  is optimal for Equation (6) if and only if  $x$  is optimal for the original problem and  $y = f(x)$  (see, e.g., Boyd & Vandenberghe (2004)).

Using Equation (6) and by introducing a new optimization variable  $b_0 \in [0, H]$ , an RCMDP becomes

$$J^* = \min_{b_0 \in [0, H], \pi \in \Pi} b_0 \text{ such that } \underbrace{J_{c_0, \mathcal{U}}(\pi) \leq b_0}_{\text{constraint to minimize objective}} \text{ and } \underbrace{J_{c_n, \mathcal{U}}(\pi) \leq b_n \forall n \in \llbracket 1, N \rrbracket}_{\text{constraints for } \pi \in \Pi_{\mathbb{F}}} . \quad (7)$$

Intuitively, Equation (7) seeks the smallest objective threshold value  $b_0$  such that there exists a feasible policy  $\pi \in \Pi_{\mathbb{F}}$  that achieves cumulative objective cost less than  $b_0$ , i.e.,  $J_{c_0, \mathcal{U}}(\pi) \leq b_0$ .

We transform Equation (7) into a more convenient form. Define  $\Delta_{b_0}(\pi) : \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \rightarrow \mathbb{R}$  such that

$$\Delta_{b_0}(\pi) = \max_{n \in \llbracket 0, N \rrbracket} J_{c_n, \mathcal{U}}(\pi) - b_n \quad \forall \pi \in \Pi, \quad (8)$$

which denotes the maximum violation of the constraints  $\max_{n \in \llbracket 1, N \rrbracket} J_{c_n, \mathcal{U}}(\pi) - b_n \leq 0$  with the additional constraint  $J_{c_0, \mathcal{U}}(\pi) - b_0 \leq 0$ . By moving  $\min_{\pi \in \Pi}$  in Equation (7) to its constraint and using  $\Delta_{b_0}(\pi)$ , Equation (7) can be transformed as follows:

**Theorem 2.** Let  $\Delta_{b_0}^* := \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$ , where  $\Delta_{b_0}(\pi)$  is defined in Equation (8). Then,

$$\text{(Epigraph Form)} \quad J^* = \min_{b_0 \in [0, H]} b_0 \text{ such that } \Delta_{b_0}^* \leq 0. \quad (9)$$

Furthermore, if  $b_0 = J^*$ , any  $\pi \in \arg \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  is optimal.

The proof is provided in Appendix I.2. Instead of Equation (7), we call Equation (9) the epigraph form of RCMDP. Since the epigraph form provides  $J^*$  and  $\pi^*$ , it overcomes the  $\pi^*$  **solution challenge** discussed in Section 3. The remaining task is to develop an algorithm to solve Equation (9).

## 5 EPIRC-PGS ALGORITHM

According to Theorem 2, we can solve an RCMDP by first identifying the optimal return value  $b_0 = J^*$  and then solving  $\min_{\pi \in \Pi} \Delta_{b_0}(\pi)$ . Our **Epigraph Robust Constrained Policy Gradient Search (EpiRC-PGS)** algorithm implements these two steps using a double-loop structure: (i) an outer loop that determines  $b_0 = J^*$  through a binary search over  $b_0 \in [0, H]$  (Section 5.1), and (ii) an inner loop that solves  $\min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  using a policy gradient subroutine (Section 5.2).

Note that without any assumptions about the uncertainty set  $\mathcal{U}$ , solving an RCMDP is NP-hard (Wiesemann et al., 2013). However, imposing concrete structures on  $\mathcal{U}$  can restrict the applicability of EpiRC-PGS. To enable EpiRC-PGS to handle a broader class of  $\mathcal{U}$ , we consider  $\mathcal{U}$  where we can approximate the robust return value  $J_{c_n, \mathcal{U}}(\pi)$  and its subgradient  $\partial J_{c_n, \mathcal{U}}(\pi)$  as follows:

**Assumption 1** (Robust policy evaluator). For each  $n \in \llbracket 0, N \rrbracket$ , we have an algorithm  $\hat{J}_n : \Pi \rightarrow \mathbb{R}$  such that  $|\hat{J}_n(\pi) - J_{c_n, \mathcal{U}}(\pi)| \leq \varepsilon_{\text{est}}$  for any  $\pi \in \Pi$ , where  $\varepsilon_{\text{est}} \geq 0$ .

**Assumption 2.**  $\mathcal{U}$  is either (i) a finite set or (ii) a compact set such that, for any  $\pi \in \Pi$ ,  $\nabla J_{c_n, P}(\pi)$  is continuous with respect to  $P \in \mathcal{U}$ .

**Assumption 3** (Subgradient evaluator). For each  $n \in \llbracket 0, N \rrbracket$ , we have an algorithm  $\hat{J}_n^\partial : \Pi \rightarrow \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$  such that  $\min_{g \in \partial J_{c_n, \mathcal{U}}(\pi)} \|\hat{J}_n^\partial(\pi) - g\|_2 \leq \varepsilon_{\text{grd}}$  for any  $\pi \in \Pi$ , where  $\varepsilon_{\text{grd}} \geq 0$ .

Assumptions 1, 2 and 3 are satisfied for most tractable uncertainty sets, such as finite, ball (Kumar et al., 2024),  $R$ -contamination (Wang & Zou, 2022),  $L_1$ ,  $\chi^2$ , and Kullback–Leibler (KL) sets (Yang et al., 2022). For these tractable uncertainty sets  $\mathcal{U}$ , the robust policy evaluator ( $\hat{J}_n$ ) and subgradient evaluator ( $\hat{J}_n^\partial$ ) can be efficiently implemented using robust DP methods (Iyengar, 2005; Kumar et al., 2022; 2024; Wang & Zou, 2022). As concrete examples, we provide detailed implementations of  $\hat{J}_n$  and  $\hat{J}_n^\partial$  for finite and KL sets in Appendix C. We assumed Assumption 2 because Danskin’s theorem (Lemmas 9 and 10), together with this assumption, guarantees that  $\partial J_{c_n, \mathcal{U}}(\pi)$  is well-defined.

---

**Algorithm 1** Double-Loop Optimization with  $\min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  Subroutine  
(also referred to as `EpIRC-PGS` when using Algorithm 2 as the subroutine)

---

- 1: **Input:** Iteration length  $K \in \mathbb{N}$ , evaluator  $\widehat{J}_n$  and subroutine  $\mathcal{A}$  (see Assumptions 1 and 4)
  - 2: Initialize the search space:  $i^{(0)} := 0$  and  $j^{(0)} := H$
  - 3: **for**  $k = 0, \dots, K - 1$  **do**
  - 4:    $\pi^{(k)} := \mathcal{A}(b_0^{(k)})$  where  $b_0^{(k)} := (i^{(k)} + j^{(k)})/2$                    // Compute policy by subroutine
  - 5:    $\widehat{\Delta}^{(k)} := \max_{n \in \llbracket 0, N \rrbracket} \widehat{J}_n(\pi^{(k)}) - b_n$  where  $b_0 = b_0^{(k)}$                    // Robust policy evaluation
  - 6:   Compute  $i^{(k+1)}$  and  $j^{(k+1)}$  by Equation (11) using  $\widehat{\Delta}^{(k)}$                    // Update search space
  - 7: **end for**
  - 8: **return**  $\pi_{\text{ret}}$  computed by  $\mathcal{A}(j^{(K)})$
- 

### 5.1 BINARY SEARCH WITH $\min_{\pi \in \Pi} \Delta_{b_0}(\pi)$ SUBROUTINE

This section describes the outer loop of `EpIRC-PGS` to identify  $b_0 = J^*$ . The outer-loop utilizes the following properties of  $\Delta_{b_0}^*$  for the identification. The proof is deferred to Appendix I.1:

**Lemma 2.**  $\Delta_{b_0}^*$  is monotonically decreasing in  $b_0$  and  $\Delta_{J^*}^* = 0$ .

Thanks to this monotonicity of  $\Delta_{b_0}^*$ , **if  $\Delta_{b_0}^*$  can be efficiently computed, a line search over  $b_0 \in [0, H]$  will readily find  $b_0 = J^*$ .** Increase  $b_0$  if  $\Delta_{b_0}^* > 0$ , and decrease it if  $\Delta_{b_0}^* \leq 0$ . Figure 2 summarizes this idea to solve the epigraph form.

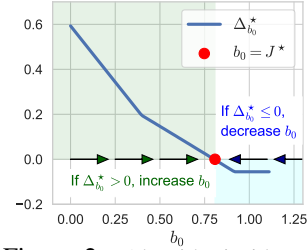


Figure 2: Algorithmic idea to find  $b_0 = J^*$  in Example 1 with  $(\gamma, \delta, b_1) = (0.1, 0, 2/3)$ .

To implement this idea, let us assume for now that we have a subroutine algorithm  $\mathcal{A}$  that computes  $\Delta_{b_0}^* = \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  with sufficient accuracy. We will implement  $\mathcal{A}$  in Section 5.1.

**Assumption 4** (Subroutine algorithm). We have an algorithm  $\mathcal{A} : \mathbb{R}_+ \rightarrow \Pi$  that takes a value  $b_0 \geq 0$  and returns  $\pi \in \Pi$  such that  $\Delta_{b_0}(\pi) \leq \min_{\pi' \in \Pi} \Delta_{b_0}(\pi') + \varepsilon_{\text{opt}}$ , where  $\varepsilon_{\text{opt}} \geq 0$ .

Using this subroutine  $\mathcal{A}$ , the outer loop conducts a binary search over  $b_0 \in [0, H]$ . Let  $K \in \mathbb{N}$  be the number of iterations. For each iteration  $k$ , let  $[i^{(k)}, j^{(k)}] \subseteq [0, H]$  be the search space where  $i^{(k)} \leq j^{(k)}$ . Set  $i^{(0)} = 0$  and  $j^{(0)} = H$ , and define  $b_0^{(k)} := (i^{(k)} + j^{(k)})/2$ . Additionally, given  $b_0^{(k)}$ , we denote the returned policy from  $\mathcal{A}$  as  $\pi^{(k)} := \mathcal{A}(b_0^{(k)})$  and its estimated  $\Delta$  value as

$$\widehat{\Delta}^{(k)} := \max_{n \in \llbracket 0, N \rrbracket} \widehat{J}_n(\pi^{(k)}) - b_n \quad \text{where } b_0 = b_0^{(k)}. \quad (10)$$

Based on Figure 2, our binary search increases  $b_0^{(k)}$  if  $\widehat{\Delta}^{(k)} > 0$ ; otherwise, it decreases  $b_0^{(k)}$ :

$$i^{(k+1)} := \begin{cases} b_0^{(k)} & \text{if } \widehat{\Delta}^{(k)} > 0 \\ i^{(k)} & \text{otherwise} \end{cases} \quad \text{and} \quad j^{(k+1)} := \begin{cases} j^{(k)} & \text{if } \widehat{\Delta}^{(k)} > 0 \\ b_0^{(k)} & \text{otherwise} \end{cases} \quad (11)$$

We summarize the pseudocode of this binary search in Algorithm 1. The following Theorem 3 ensures that Algorithm 1 returns a near-optimal policy. We provide the proof in Appendix J.1.

**Theorem 3.** Suppose that Algorithm 1 is run with algorithms  $\widehat{J}_n$  and  $\mathcal{A}$  that satisfy Assumptions 1 and 4. Then, Algorithm 1 returns an  $\tilde{\varepsilon}$ -optimal policy, where  $\tilde{\varepsilon} := 2(\varepsilon_{\text{opt}} + \varepsilon_{\text{est}}) + 2^{-K}H$ .

### 5.2 SUBROUTINE ALGORITHM TO SOLVE $\min_{\pi \in \Pi} \Delta_{b_0}(\pi)$

The remaining task is to implement the subroutine  $\mathcal{A}$  which satisfies Assumption 4. In other words, for a given  $b_0$ , we need to solve the following auxiliary problem:

**(Epigraph’s Auxiliary Problem)**  $\min_{\pi \in \Pi} \Delta_{b_0}(\pi) = \min_{\pi \in \Pi} \max_{n \in \llbracket 0, N \rrbracket} \max_{P \in \mathcal{U}} J_{c_n, P}(\pi) - b_n. \quad (12)$

The right-hand side of Equation (12) can be seen as an RMDP with additional robustness over the set of modified cost functions  $\mathcal{C}_{b_0} := \{c_n - b_n/H\}_{n \in \llbracket 0, N \rrbracket}$ . Note that since  $\mathcal{C}_{b_0}$  is not a rectangular set,

**Algorithm 2** Projected Policy Gradient Subroutine

---

```

1: Input: Threshold parameter  $b_0 \geq 0$ , learning rate  $\alpha > 0$ , iteration length  $T \in \mathbb{N}$ , evaluator  $\widehat{J}_n$ 
   and subgradient evaluator  $\widehat{J}_n^\partial$  (see Assumptions 1 and 3)
2: Set an arbitrary initial policy  $\pi^{(0)} \in \Pi$ 
3: for  $t = 0, \dots, T - 1$  do
4:    $n^{(t)} \in \arg \max_{n \in \llbracket 0, N \rrbracket} \widehat{J}_n(\pi^{(t)}) - b_n$ . // Select the most violated constraint
5:    $\pi^{(t+1)} := \text{Proj}_\Pi(\pi^{(t)} - \alpha g^{(t)})$  where  $g^{(t)} := \widehat{J}_{n^{(t)}}^\partial(\pi^{(t)})$  // Update policy
6: end for
7: return  $\pi^{(t^*)}$  where  $t^* \in \arg \min_{t \in \llbracket 0, T-1 \rrbracket} \widehat{\Delta}^{(t)}$  and  $\widehat{\Delta}^{(t)} := \max_{n \in \llbracket 0, N \rrbracket} \widehat{J}_n(\pi^{(t)}) - b_n$ 

```

---

even if  $\mathcal{U}$  is rectangular, the combination  $\mathcal{C}_{b_0} \times \mathcal{U}$  does not retain rectangularity. As a result, existing RMDP algorithms designed for rectangular sets, such as DP (Iyengar, 2005), natural policy gradient (Li et al., 2022), and convex optimization (Grand-Clément & Petrik, 2024), are inapplicable.

Due to this non-rectangularity issue, we employ the projected policy gradient method, together with the following subgradient of  $\partial \Delta_{b_0}(\pi)$ . The proof is deferred to Appendix I.3:

**Lemma 3.** Define  $\mathcal{G}_{b_0}(\pi) := \left\{ \nabla J_{c_n, P}(\pi) \mid n, P \in \arg \max_{(n, P) \in \llbracket 0, N \rrbracket \times \mathcal{U}} J_{c_n, P}(\pi) - b_n \right\}$  for any  $b_0 \in \mathbb{R}$  and  $\pi \in \Pi$ . Let  $\text{conv } B$  denote the convex hull of a set  $B \subset \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ . Under Assumption 2, for any  $\pi \in \Pi$  and  $b_0 \in \mathbb{R}$ , the subgradient of  $\Delta_{b_0}(\cdot)$  at  $\pi$  is given by  $\partial \Delta_{b_0}(\pi) = \text{conv } \mathcal{G}_{b_0}(\pi)$ .

We implement the subroutine  $\mathcal{A}(b_0)$  based on this subgradient lemma. Starting from an arbitrary policy  $\pi^{(0)}$ , let  $\pi^{(1)}, \dots, \pi^{(T)}$  be the updated policies where  $T \in \mathbb{N}$  is the iteration length. Using the evaluators  $\widehat{J}_n, \widehat{J}_n^\partial$ , and a learning rate  $\alpha > 0$ , for a given  $b_0$ , our subroutine updates policy as follows:

$$\pi^{(t+1)} := \text{Proj}_\Pi(\pi^{(t)} - \alpha g^{(t)}) \text{ where } g^{(t)} := \widehat{J}_{n^{(t)}}^\partial(\pi^{(t)}) \text{ and } n^{(t)} \in \arg \max_{n \in \llbracket 0, N \rrbracket} \widehat{J}_n(\pi^{(t)}) - b_n. \quad (13)$$

We summarize the pseudocode of this policy update subroutine in Algorithm 2.

**Remark 1** (Comparison to Lagrange). Recall Equation (5) that the subgradient of the Lagrangian’s auxiliary problem,  $\partial L_\lambda(\pi)$ , involves a summation of policy gradients over different environments. On the other hand, Equation (13) focuses on the policy gradient of a single worst-case environment by taking  $\max_{n \in \llbracket 0, N \rrbracket}$ . Intuitively, our policy update avoids the sum of conflicting policy gradients, thereby circumventing the **gradient conflict challenge** discussed in Section 3. Indeed, when the initial distribution satisfies the following coverage assumption<sup>6</sup>, **there is no local minimum in  $\Delta_{b_0}(\pi)$** :

**Assumption 5** (Initial distribution coverage). The initial distribution  $\mu \in \mathcal{P}(\mathcal{S})$  satisfies  $\mu > \mathbf{0}$ .

**Theorem 4** (Optimality of stationary points). Under Assumptions 2 and 5, for any  $(\pi, b_0) \in \Pi \times \mathbb{R}$ ,

$$\Delta_{b_0}(\pi) - \min_{\pi' \in \Pi} \Delta_{b_0}(\pi') \leq DH \max_{\pi' \in \Pi} \langle \pi - \pi', g \rangle \quad \forall g \in \partial \Delta_{b_0}(\pi),$$

where  $D := \max_{n, P \in \llbracket 0, N \rrbracket \times \mathcal{U}} \left\| d_P^{\pi_n^*, P} / \mu \right\|_\infty$  with  $\pi_n^*, P \in \arg \min_{\pi' \in \Pi} J_{c_n, P}(\pi')$ .

The detailed proof can be found in Appendix I.4. Our proof is similar to **Theorem 3.2** in Wang et al. (2023), but it is more rigorous and corrects a crucial error that can invalidate their result<sup>7</sup>. Moreover, while their proof is limited to cases where  $\arg \max_{P \in \mathcal{U}} J_{c_0, P}(\pi)$  is finite, ours is not. We leverage Sion’s minimax theorem (Sion, 1958) for this refinement.

Thanks to the optimality of stationary points of  $\Delta_{b_0}(\pi)$  (Theorem 4), Algorithm 2 is guaranteed to solve  $\min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  and satisfies the requirement of  $\mathcal{A}(b_0)$  in Assumption 4 as follows:

**Theorem 5.** Suppose Assumptions 1, 2, 3 and 5 hold. Then, there exist problem-dependent constants  $C_\partial, C_J, C_\alpha, C_T > 0$  that do not depend on  $\varepsilon$  such that, when Algorithm 2 is run with  $\alpha = C_\alpha \varepsilon^2$  and  $T = C_T \varepsilon^{-4}$ , if the evaluators are sufficiently accurate such that  $\varepsilon_{\text{grd}} = C_\partial \varepsilon^2$  and  $\varepsilon_{\text{est}} = C_J \varepsilon^2$ , Algorithm 2 returns a policy  $\pi^{(t^*)}$  satisfying  $\Delta_{b_0}(\pi^{(t^*)}) - \min_{\pi \in \Pi} \Delta_{b_0}(\pi) \leq \varepsilon$ .

<sup>6</sup>Such coverage assumption is necessary to ensure the global convergence of policy gradient methods (Mei et al., 2020). Additionally, note that the Lagrange performs poorly even under Assumption 5 (see Theorem 1).

<sup>7</sup>For example, their proof around **Equation (32)** incorrectly bounds a positive value by a negative value.



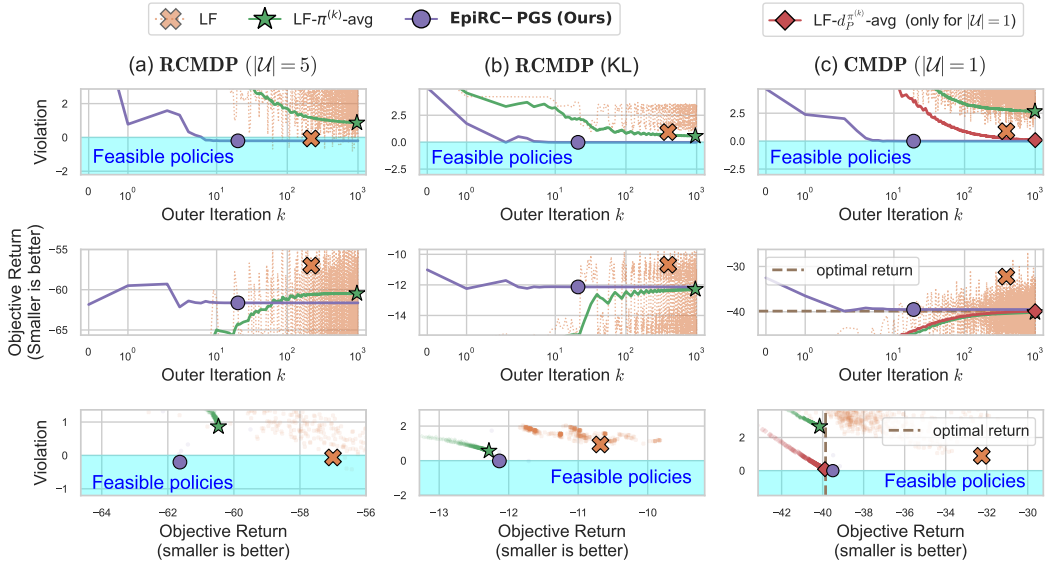


Figure 3: Comparison of the algorithms in different settings (a), (b), and (c), defined in Section 6. The feasible  $\pi^{(k)}$  with the smallest return is marked; if none is feasible, the one with the smallest violation is marked. In all the settings, EpiRC-PGS quickly identifies a feasible and low-return policy (●). **Top row:** Constraint violation ( $y$ -axis:  $J_{c_1, \mathcal{U}}(\pi^{(k)}) - b_1$ ). Policies in the blue area satisfy the constraints. **Middle row:** Objective return relative to the uniform policy ( $y$ -axis:  $J_{c_0, \mathcal{U}}(\pi^{(k)}) - J_{c_0, \mathcal{U}}(\pi_{\text{unif}})$ ). Negative values indicate that the policies achieve non-trivial low cumulative objective costs. **Bottom row:** Constraint violation vs. relative objective return.

We provide the proof and the concrete values of  $C_{\partial}$ ,  $C_J$ ,  $C_{\alpha}$  and  $C_T$  in Appendix J.2. The proof is primarily based on the weakly convex function analysis (Beck, 2017; Wang et al., 2022).

### 5.3 COMBINING BINARY SEARCH WITH THE POLICY GRADIENT SUBROUTINE

Finally, we combine Algorithm 1 with Algorithm 2 subroutine and refer to the combination as EpiRC-PGS. According to Theorems 3 and 5, EpiRC-PGS is ensured to find an  $\varepsilon$ -optimal policy:

**Corollary 1.** Consider the same settings and notations as in Theorem 5. Set Algorithm 2 as the subroutine  $\mathcal{A}$  with parameters  $(\hat{J}_n, \hat{J}_n^{\partial}, \alpha, T)$ , where we set  $\alpha = C_{\alpha}\varepsilon^2/4$ ,  $T = 16C_T\varepsilon^{-4}$ . Then, given inputs  $\hat{J}_n$  and  $\mathcal{A}_n$ , Algorithm 1 returns an  $\varepsilon$ -optimal policy after  $K = \lceil \log(2H\varepsilon^{-1}) \rceil$  iteration.

**Remark 2** (Computational complexity). EpiRC-PGS outputs an  $\varepsilon$ -optimal policy by querying  $\hat{J}_n$  and  $\hat{J}_n^{\partial}$  a total of  $\tilde{O}((N+1)KT)$  times. Thus, the computational complexity of EpiRC-PGS can be expressed as  $\tilde{O}((N+1)KT \times [\text{querying cost}])$ . As a simple example, consider the case where  $\mathcal{U}$  is finite, where querying  $\hat{J}_n(\pi)$  and  $\hat{J}_n^{\partial}(\pi)$  require  $\tilde{O}(S^2A|\mathcal{U}|)$  operations Kumar et al. (2024). Using the concrete value of  $KT$ , the computational complexity of EpiRC-PGS for finite  $\mathcal{U}$  becomes  $\tilde{O}(D^4S^5A^4H^{14}|\mathcal{U}|(N+1)\varepsilon^{-4})$ . Similar analyses can be applied to other types of uncertainty sets.

**Remark 3** (Last-iterate convergence). Lagrangian-based algorithms for CMDPs typically require the average of past policies (e.g., Li et al. (2024); Liu et al. (2021)). However, they encounter difficulties in scenarios where policy averaging is impractical, such as in deep RL applications. In contrast, Corollary 1 does not require policy averaging and ensures that the final policy output is near-optimal.

## 6 EXPERIMENTS

To support the theoretical guarantees of EpiRC-PGS and demonstrate the limitations of the Lagrangian methods in identifying near-optimal policies, this section empirically evaluates EpiRC-PGS in three settings with five constraints ( $N = 5$ ):

- 486 (a) RCMDP with  $\mathcal{U} = \{P_1, \dots, P_5\}$ , where each environment  $P$  is randomly generated.  
 487  
 488 (b) RCMDP with a  $(s, a)$ -rectangular KL uncertainty set (Iyengar, 2005), which considers  
 489  $\mathcal{U} = \times_{s,a} \{p \in \mathcal{P}(\mathcal{S}) \mid \text{KL}[p \parallel P(\cdot \mid s, a)] \leq C_{\text{KL}}\}$ <sup>8</sup> where  $P$  is a nominal environment.  
 490  
 491 (c) CMDP, which is equivalent to RCMDP with  $\mathcal{U} = \{P\}$ .

492 **Environment construction.** In (a), (b), and (c), we randomly generate environments following the  
 493 experimental setup of Dann et al. (2017). For each state-action pair  $(s, a)$ , the transition probabilities  
 494  $P(\cdot \mid s, a)$  are independently sampled from a Dirichlet(0.1, ..., 0.1) distribution. This produces  
 495 a concentrated yet non-deterministic transition model, resembling the widely used GARNET  
 496 benchmark with a *branching factor* of 1 (Archibald et al., 1995). Cost values  $c_n(s, a)$  are assigned  
 497 as 0 with a probability of 0.1 and 1 otherwise. Initial state probabilities  $\mu(\cdot)$  are sampled from a  
 498 Dirichlet(0.5, ..., 0.5) distribution. Constraint thresholds  $b_1, \dots, b_5$  are configured to guarantee  
 499 the existence of a feasible policy. Other environmental parameters are described in Appendix D.

500 **Baseline algorithms.** We compare `EpiRC-PGS` to a Lagrangian counterpart, denoted `LF`, which  
 501 abstracts most of the existing Lagrangian-based algorithms for RCMDPs (e.g., (Russel et al., 2020;  
 502 Wang et al., 2022)). `LF` aims to solve the problem  $\max_{\lambda \in \mathbb{R}_+^n} \min_{\pi \in \Pi} L_\lambda(\pi)$  in Equation (3) by  
 503 performing gradient ascent on  $\lambda$  while using a policy gradient subroutine to solve  $\min_{\pi \in \Pi} L_\lambda(\pi)$ .

504 We also evaluate the averaged policies generated by `LF`, defined as  $\frac{1}{k} \sum_{j=0}^k \pi^{(k)}$ . Such policy  
 505 averaging is employed in Lagrangian methods (Miryoosefi et al., 2019; Zhang et al., 2024), though it  
 506 often lacks theoretical guarantees (see Appendix A.4). In the CMDP setting (c), we further evaluate  
 507 the averaged occupancy measures, where the  $k$ -th policy is derived from  $\frac{1}{k} \sum_{j=0}^k d_P^{\pi^{(k)}}$ . Averaging  
 508  $d_P^{\pi^{(k)}}$  is ensured to identify a near-optimal policy (Zahavy et al., 2021), but is well-defined only  
 509 when  $|\mathcal{U}| = 1$ . We refer to these two averagings as `LF- $\pi^{(k)}$ -avg` and `LF- $d_P^{\pi^{(k)}}$ -avg`, respectively.  
 510 Moreover, (c) reports the optimal return value computed by an LP method (Altman, 1999). The  
 511 detailed implementation of algorithms are provided in Appendix D.

512 **Results.** Figure 3 illustrates the performance of the algorithms averaged over 10 random seeds.  
 513 In all the settings, `EpiRC-PGS` rapidly converges to a feasible policy with a low objective  
 514 return, while both `LF` and `LF- $\pi^{(k)}$ -avg` fail to identify feasible policies in certain settings  
 515 (e.g., `LF- $\pi^{(k)}$ -avg` in (a) and `LF` in (b)). In the CMDP (c), `EpiRC-PGS` and `LF- $d_P^{\pi^{(k)}}$ -avg`  
 516 converge to a near-optimal policy, but `LF- $\pi^{(k)}$ -avg` does not. These results empirically validate  
 517 that `EpiRC-PGS` yields a near-optimal policy in RCMDPs, contrasting with the conventional  
 518 Lagrangian-based algorithm’s inability in robust settings.  
 519  
 520  
 521

## 522 7 CONCLUSION AND LIMITATIONS

523  
 524 In this work, we propose `EpiRC-PGS`, the first algorithm guaranteed to find a near-optimal policy  
 525 in an RCMDP (Corollary 1). At the core of `EpiRC-PGS` is the use of the epigraph form for  
 526 RCMDP. Remarkably, the epigraph form produces the optimal policy  $\pi^*$  (Section 4) and supports a  
 527 policy gradient algorithm to find it (Theorem 4). These features effectively address the optimization  
 528 challenges encountered in the conventional Lagrangian formulation (Section 3).  
 529

530 **Limitations and future work.** A double-loop algorithm like `EpiRC-PGS` is often inefficient  
 531 when the inner problem requires high computational cost (Lin et al., 2024). Developing a single-loop  
 532 algorithm is a promising direction for future research, and we discuss the challenges in Appendix B.

533 Another research avenue is improving the iteration complexity of our  $\tilde{\mathcal{O}}(\varepsilon^{-4})$ . This may not be  
 534 tight, since for RMDPs with  $(s, a)$ -rectangularity, the natural policy gradient method is ensured to  
 535 find an  $\varepsilon$ -optimal policy with  $\tilde{\mathcal{O}}(\varepsilon^{-2})$  iterations Li et al. (2022).  
 536

537 Finally, the coverage assumption on the initial distribution (Assumption 5) is not necessary in  
 538 CMDPs (Ding et al., 2024). We leave the removal of Assumption 5 in RCMDPs for future work.

539 <sup>8</sup> $\text{KL}[p \parallel q] = \sum_{s \in \mathcal{S}} p(s) \ln p(s)/q(s)$  represents the KL divergence between two probability distributions  
 $p > \mathbf{0}$  and  $q > \mathbf{0}$  defined over  $\mathcal{S}$ .  $C_{\text{KL}} > 0$  is a positive constant.

## REFERENCES

- 540  
541  
542 Jacob D Abernethy and Jun-Kun Wang. On Frank-Wolfe and Equilibrium Computation. In *Advances*  
543 *in Neural Information Processing Systems*, 2017.
- 544  
545 Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained Policy Optimization. In  
546 *International Conference on Machine Learning*, 2017.
- 547  
548 Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the Theory of Policy Gra-  
549 *Research*, 22(98):1–76, 2021.
- 550  
551 Eitan Altman. *Constrained Markov Decision Processes*, volume 7. CRC Press, 1999.
- 552  
553 James Anderson, John C Doyle, Steven H Low, and Nikolai Matni. System Level Synthesis. *Annual*  
554 *Reviews in Control*, 47:364–393, 2019.
- 555  
556 TW Archibald, KIM McKinnon, and LC Thomas. On the Generation of Markov Decision Processes.  
557 *Journal of the Operational Research Society*, 46(3):354–361, 1995.
- 558  
559 Amir Beck. *First-Order Methods in Optimization*. SIAM, 2017.
- 560  
561 R. Bellman, R.E. Bellman, and Rand Corporation. *Dynamic Programming*. Princeton University  
562 *Press*, 1957.
- 563  
564 Alberto Bemporad and Manfred Morari. Robust Model Predictive Control: A Survey. In *Robustness*  
565 *in Identification and Control*, pp. 207–226. Springer, 2007.
- 566  
567 Hans-Georg Beyer and Bernhard Sendhoff. Robust Optimization—a Comprehensive Survey. *Com-*  
568 *puter Methods in Applied Mechanics and Engineering*, 196(33-34):3190–3218, 2007.
- 569  
570 Stephen P Boyd and Lieven Vandenbergh. *Convex Optimization*. Cambridge University Press,  
571 2004.
- 572  
573 Archana Bura, Aria HasanzadeZonuzy, Dileep Kalathil, Srinivas Shakkottai, and Jean-Francois  
574 Chamberland. DOPE: Doubly Optimistic and Pessimistic Exploration for Safe Reinforcement  
575 *Learning*. In *Advances in Neural Information Processing Systems*, 2022.
- 576  
577 Hyeong Soo Chang. Approximate Constrained Discounted Dynamic Programming with Uniform  
578 *Feasibility and Optimality*. *arXiv preprint arXiv:2308.03297*, 2023.
- 579  
580 Richard C Chen and Gilmer L Blankenship. Dynamic Programming Equations for Discounted  
581 *Constrained Stochastic Control*. *IEEE Transactions on Automatic Control*, 49(5):699–709, 2004.
- 582  
583 Richard C Chen and Eugene A Feinberg. Non-Randomized Control of Constrained Markov Decision  
584 *Processes*. In *American Control Conference*, 2006.
- 585  
586 Yi Chen, Jing Dong, and Zhaoran Wang. A Primal-Dual Approach to Constrained Markov Decision  
587 *Processes*. *arXiv preprint arXiv:2101.10895*, 2021.
- 588  
589 Zhi Chen, Pengqian Yu, and William B Haskell. Distributionally Robust Optimization for Sequential  
590 *Decision-Making*. *Optimization*, 68(12):2397–2426, 2019.
- 591  
592 Francis H Clarke, Yuri S Ledyaev, Ronald J Stern, and Peter R Wolenski. *Nonsmooth analysis and*  
593 *Control Theory*, volume 178. Springer Science & Business Media, 2008.
- Frank H Clarke. *Optimization and Nonsmooth Analysis*. SIAM, 1990.
- Pierre Clavier, Erwan Le Pennec, and Matthieu Geist. Towards minimax Optimality of Model-based  
Robust Reinforcement Learning. *arXiv preprint arXiv:2302.05372*, 2023.
- Christoph Dann, Tor Lattimore, and Emma Brunskill. Unifying PAC and Regret: Uniform PAC  
Bounds for Episodic Reinforcement Learning. In *Advances in Neural Information Processing*  
*Systems*, 2017.

- 594 Eric V Denardo. On Linear Programming in a Markov Decision Problem. *Management Science*, 16  
595 (5):281–288, 1970.
- 596
- 597 Esther Derman, Matthieu Geist, and Shie Mannor. Twice Regularized MDPs and the Equivalence  
598 Between Robustness and Regularization. In *Advances in Neural Information Processing Systems*,  
599 2021.
- 600 Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo Jovanovic. Natural Policy Gradient  
601 Primal-Dual Method for Constrained Markov Decision Processes. In *Advances in Neural Infor-*  
602 *mation Processing Systems*, 2020.
- 603
- 604 Dongsheng Ding, Chen-Yu Wei, Kaiqing Zhang, and Alejandro Ribeiro. Last-Iterate Convergent  
605 Policy Gradient Primal-Dual Methods for Constrained MDPs. In *Advances in Neural Information*  
606 *Processing Systems*, 2024.
- 607 John Doyle. Analysis of Feedback Systems with Structured Uncertainties. In *IEE Proceedings*  
608 *D-Control Theory and Applications*, 1982.
- 609
- 610 Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong  
611 Wang. Bilinear Classes: A Structural Framework for Provable Generalization in RL. In *Interna-*  
612 *tional Conference on Machine Learning*, 2021.
- 613 Simon S Du, Yining Wang, Sivaraman Balakrishnan, Pradeep Ravikumar, and Aarti Singh.  
614 Robust Nonparametric Regression under Huber’s  $\epsilon$ -contamination Model. *arXiv preprint*  
615 *arXiv:1805.10406*, 2018.
- 616 Yonathan Efroni, Shie Mannor, and Matteo Pirotta. Exploration-Exploitation in Constrained MDPs.  
617 *arXiv preprint arXiv:2003.02189*, 2020.
- 618
- 619 Arnob Ghosh. Sample Complexity for Obtaining Sub-optimality and Violation Bound for Distribu-
- 620 tionally Robust Constrained MDP. In *Reinforcement Learning Safety Workshop*, 2024.
- 621
- 622 Vineet Goyal and Julien Grand-Clement. Robust Markov Decision Processes: Beyond Rectangular-
- 623 ity. *Mathematics of Operations Research*, 48(1):203–226, 2023.
- 624 Julien Grand-Clément and Christian Kroer. Scalable First-Order Methods for Robust MDPs. In *AAAI*  
625 *Conference on Artificial Intelligence*, 2021.
- 626
- 627 Julien Grand-Clément and Marek Petrik. On the Convex Formulations of Robust Markov Decision  
628 Processes. *Mathematics of Operations Research*, 0(0), 2024.
- 629 Aria HasanzadeZonuzi, Archana Bura, Dileep Kalathil, and Srinivas Shakkottai. Learning with  
630 Safety Constraints: Sample Complexity of Reinforcement Learning for Constrained MDPs. In  
631 *AAAI Conference on Artificial Intelligence*, 2021.
- 632
- 633 Moshe Haviv. On Constrained Markov Decision Processes. *Operations Research Letters*, 19(1):  
634 25–28, 1996.
- 635 Chin Pang Ho, Marek Petrik, and Wolfram Wiesemann. Partial Policy Iteration for  $L_1$ -Robust  
636 Markov Decision Processes. *Journal of Machine Learning Research*, 22(275):1–46, 2021.
- 637
- 638 Garud N Iyengar. Robust Dynamic Programming. *Mathematics of Operations Research*, 30(2):  
639 257–280, 2005.
- 640 Nan Jiang. PAC Reinforcement Learning with an Imperfect Model. In *AAAI Conference on Artificial*  
641 *Intelligence*, 2018.
- 642
- 643 Toshinori Kitamura, Tadashi Kozuno, Masahiro Kato, Yuki Ichihara, Soichiro Nishimori, Akiyoshi  
644 Sannai, Sho Sonoda, Wataru Kumagai, and Yutaka Matsuo. A Policy Gradient Primal-Dual Algo-  
645 rithm for Constrained MDPs with Uniform PAC Guarantees. *arXiv preprint arXiv:2401.17780*,  
646 2024.
- 647 A Ya Kruger. On Fréchet Subdifferentials. *Journal of Mathematical Sciences*, 116(3):3325–3358,  
2003.

- 648 Navdeep Kumar, Kfir Levy, Kaixin Wang, and Shie Mannor. Efficient Policy Iteration for Robust  
649 Markov Decision Processes via Regularization. *arXiv preprint arXiv:2205.14327*, 2022.  
650
- 651 Navdeep Kumar, Esther Derman, Matthieu Geist, Kfir Y Levy, and Shie Mannor. Policy Gradient for  
652 Rectangular Robust Markov Decision Processes. In *Advances in Neural Information Processing*  
653 *Systems*, 2024.
- 654 Hoang Le, Cameron Voloshin, and Yisong Yue. Batch Policy Learning Under Constraints. In  
655 *International Conference on Machine Learning*, 2019.  
656
- 657 Tianjiao Li, Ziwei Guan, Shaofeng Zou, Tengyu Xu, Yingbin Liang, and Guanghui Lan. Faster  
658 Algorithm and Sharper Analysis for Constrained Markov Decision Process. *Operations Research*  
659 *Letters*, 54:107107, 2024.
- 660 Yan Li, Guanghui Lan, and Tuo Zhao. First-Order Policy Optimization for Robust Markov Decision  
661 Process. *arXiv preprint arXiv:2209.10579*, 2022.  
662
- 663 Zhenwei Lin, Chenyu Xue, Qi Deng, and Yinyu Ye. A Single-Loop Robust Policy Gradient Method  
664 for Robust Markov Decision Processes. *arXiv preprint arXiv:2406.00274*, 2024.
- 665 Tao Liu, Ruida Zhou, Dileep Kalathil, Panganamala Kumar, and Chao Tian. Learning Policies with  
666 Zero or Bounded Constraint Violation for Constrained MDPs. In *Advances in Neural Information*  
667 *Processing Systems*, 2021.  
668
- 669 Tien Mai and Patrick Jaillet. Robust Entropy-Regularized Markov Decision Processes. *arXiv*  
670 *preprint arXiv:2112.15364*, 2021.
- 671 Daniel J Mankowitz, Dan A Calian, Rae Jeong, Cosmin Paduraru, Nicolas Heess, Sumanth  
672 Dathathri, Martin Riedmiller, and Timothy Mann. Robust Constrained Reinforcement Learning  
673 for Continuous Control with Model Misspecification. *arXiv preprint arXiv:2010.10644*, 2020.  
674
- 675 Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the Global Convergence  
676 Rates of Softmax Policy Gradient Methods. In *International Conference on Machine Learning*,  
677 2020.
- 678 VS Mikhalevich, AM Gupal, and VI Norkin. Methods of Nonconvex Optimization. *arXiv preprint*  
679 *arXiv:2406.10406*, 2024.
- 680 Sobhan Miryoosefi, Kianté Brantley, Hal Daume III, Miro Dudik, and Robert E Schapire. Reinforce-  
681 ment Learning with Convex Constraints. In *Advances in Neural Information Processing Systems*,  
682 2019.  
683
- 684 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Belle-  
685 mare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level  
686 Control through Deep Reinforcement Learning. *Nature*, 518(7540):529–533, 2015.
- 687 Adrian Müller, Pragnya Alatur, Volkan Cevher, Giorgia Ramponi, and Niao He. Truly No-Regret  
688 Learning in Constrained MDPs. *arXiv preprint arXiv:2402.15776*, 2024.  
689
- 690 Ofir Nachum and Bo Dai. Reinforcement Learning via Fenchel-Rockafellar Duality. *arXiv preprint*  
691 *arXiv:2001.01866*, 2020.
- 692 Arnab Nilim and Laurent El Ghaoui. Robust Control of Markov Decision Processes with Uncertain  
693 Transition Matrices. *Operations Research*, 53(5):780–798, 2005.  
694
- 695 Kishan Panaganti and Dileep Kalathil. Sample Complexity of Robust Reinforcement Learning with  
696 a Generative Model. In *International Conference on Artificial Intelligence and Statistics*, 2022.
- 697 Neal Parikh, Stephen Boyd, et al. Proximal algorithms. *Foundations and Trends® in Optimization*,  
698 1(3):127–239, 2014.  
699
- 700 Santiago Paternain, Luiz Chamon, Miguel Calvo-Fullana, and Alejandro Ribeiro. Constrained Re-  
701 inforcement Learning Has Zero Duality Gap. In *Advances in Neural Information Processing*  
*Systems*, 2019.

- 702 Santiago Paternain, Miguel Calvo-Fullana, Luiz FO Chamon, and Alejandro Ribeiro. Safe Policies  
703 for Reinforcement Learning via Primal-Dual Methods. *IEEE Transactions on Automatic Control*,  
704 68(3):1321–1336, 2022.
- 705  
706 Matteo Pirodda, Marcello Restelli, Alessio Pecorino, and Daniele Calandriello. Safe Policy Iteration.  
707 In *International Conference on Machine Learning*, 2013.
- 708 Hamed Rahimian and Sanjay Mehrotra. Distributionally Robust Optimization: A Review. *arXiv*  
709 *preprint arXiv:1908.05659*, 2019.
- 710  
711 R Tyrrell Rockafellar and Roger J-B Wets. *Variational Analysis*, volume 317. Springer Science &  
712 Business Media, 2009.
- 713 Reazul Hasan Russel, Mouhacine Benosman, and Jeroen Van Baar. Robust Constrained-  
714 MDPs: Soft-Constrained Robust Policy Optimization under Model Uncertainty. *arXiv preprint*  
715 *arXiv:2010.04870*, 2020.
- 716  
717 Maurice Sion. On General Minimax Theorems. *Pacific Journal of Mathematics*, 8(1):171 – 176,  
718 1958.
- 719  
720 Oswin So and Chuchu Fan. Solving Stabilize-Avoid Optimal Control via Epigraph Form and Deep  
721 Reinforcement Learning. In *Proceedings of Robotics: Science and Systems*, 2023.
- 722  
723 Oswin So, Cheng Ge, and Chuchu Fan. Solving Minimum-Cost Reach Avoid using Reinforcement  
724 Learning. In *Advances in Neural Information Processing Systems*, 2024.
- 725  
726 Zhongchang Sun, Sihong He, Fei Miao, and Shaofeng Zou. Constrained Reinforcement Learning  
727 Under Model Mismatch. *arXiv preprint arXiv:2405.01327*, 2024.
- 728  
729 Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward Constrained Policy Optimization. In  
730 *International Conference on Learning Representations*, 2018.
- 731  
732 Qiu hao Wang, Chin Pang Ho, and Marek Petrik. Policy Gradient in Robust MDPs with Global  
733 Convergence Guarantee. In *International Conference on Machine Learning*, 2023.
- 734  
735 Yue Wang and Shaofeng Zou. Online Robust Reinforcement Learning with Model Uncertainty. In  
736 *Advances in Neural Information Processing Systems*, 2021.
- 737  
738 Yue Wang and Shaofeng Zou. Policy Gradient Method for Robust Reinforcement Learning. In  
739 *International Conference on Machine Learning*, 2022.
- 740  
741 Yue Wang, Fei Miao, and Shaofeng Zou. Robust Constrained Reinforcement Learning. *arXiv*  
742 *preprint arXiv:2209.06866*, 2022.
- 743  
744 Honghao Wei, Xin Liu, and Lei Ying. A Provably-Efficient Model-Free Algorithm for Constrained  
745 Markov Decision Processes. *arXiv preprint arXiv:2106.01577*, 2021.
- 746  
747 Wolfram Wiesemann, Daniel Kuhn, and Berç Rustem. Robust Markov Decision Processes. *Mathe-*  
748 *matics of Operations Research*, 38(1):153–183, 2013.
- 749  
750 Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Toward Theoretical Understandings of Robust  
751 Markov Decision Processes: Sample Complexity and Asymptotics. *The Annals of Statistics*, 50  
752 (6):3223–3248, 2022.
- 753  
754 Wenhao Yang, Han Wang, Tadashi Kozuno, Scott M Jordan, and Zhihua Zhang. Robust Markov  
755 Decision Processes without Model Estimation. *arXiv preprint arXiv:2302.01248*, 2023.
- 756  
757 Donghao Ying, Yuhao Ding, and Javad Lavaei. A Dual Approach to Constrained Markov Decision  
758 Processes with Entropy Regularization. In *International Conference on Artificial Intelligence and*  
759 *Statistics*, 2022.
- 760  
761 Tom Zahavy, Brendan O’Donoghue, Guillaume Desjardins, and Satinder Singh. Reward is Enough  
762 for Convex MDPs. In *Advances in Neural Information Processing Systems*, 2021.

George Zames. Feedback and Optimal Sensitivity: Model Reference Transformations, Multiplicative Seminorms, and Approximate Inverses. *IEEE Transactions on Automatic Control*, 26(2): 301–320, 1981.

Zhengfei Zhang, Kishan Panaganti, Laixi Shi, Yanan Sui, Adam Wierman, and Yisong Yue. Distributionally Robust Constrained Reinforcement Learning under Strong Duality. *arXiv preprint arXiv:2406.15788*, 2024.

Liyuan Zheng and Lillian Ratliff. Constrained Upper Confidence Reinforcement Learning. In *Learning for Dynamics and Control*, 2020.

## A ADDITIONAL RELATED WORK

This section reviews existing approaches for CMDPs (Appendix A.1), RMDPs (Appendix A.2), and RCMDPs (Appendix A.3). It also highlights their inherent limitations and the challenges they face when applied to RCMDPs.

### A.1 CONSTRAINED MARKOV DECISION PROCESSES

CMDP is a specific subclass of RCMDP where the uncertainty set consists of a single element, i.e.,  $\mathcal{U} = \{P\}$ . This section describes the primary approaches to the CMDP problem: the linear programming (LP) approach, the Lagrangian approach, and the epigraph approach.

**Linear programming approach.** The LP approach has been extensively studied in the theoretical literature (Efroni et al., 2020; Liu et al., 2021; Bura et al., 2022; HasanzadeZonuzuy et al., 2021; Zheng & Ratliff, 2020). Although it is a fundamental method in CMDP, it is less popular in practice due to its difficulty in scaling to high-dimensional problem settings, such as those encountered in deep RL. Additionally, incorporating environmental uncertainty into the LP approach for CMDPs is challenging. The LP approach utilizes the fact that the return minimization problem of an MDP can be formulated as a convex optimization problem with respect to the occupancy measure (Altman, 1999; Nachum & Dai, 2020). However, RMDPs do not permit a convex formulation in terms of occupancy measures (Iyengar, 2005; Grand-Clément & Petrik, 2024). While Grand-Clément & Petrik (2024) recently introduced a convex optimization approach for RMDPs, their formulation is convex for the transformed objective value function, not for the occupancy measure, making it challenging to incorporate constraints as seen in RCMDPs.

**Lagrangian approach.** The Lagrangian approach is perhaps the most popular approach to CMDPs both in theory (Ding et al., 2020; Wei et al., 2021; HasanzadeZonuzuy et al., 2021; Kitamura et al., 2024) and practice (Achiam et al., 2017; Tessler et al., 2018; Wang et al., 2022; Le et al., 2019; Russel et al., 2020). This popularity stems from its compatibility with policy gradient methods, making it readily extendable to deep RL. The Lagrangian approach benefits from the strong duality in CMDPs. When  $\mathcal{U}$  consists of a single element, it is well established that strong duality holds, meaning that  $L^* = J^*$  holds, where  $L^*$  is from Equation (3) and  $J^*$  is from Equation (2) (Altman, 1999; Paternain et al., 2019; 2022).

The challenge with the Lagrangian method is the identification of an optimal policy. Even if Equation (3) is solved, there’s no guarantee that the solution to the inner minimization problem will represent an optimal policy. In some CMDPs, where feasible policies in  $\Pi_F$  must be stochastic (Altman, 1999), the inner minimization may yield a deterministic solution that is infeasible. Zahavy et al. (2021); Miryoosefi et al. (2019); Chen et al. (2021); Li et al. (2024); Liu et al. (2021) addressed this challenge by averaging policies (or occupancy measures) obtained during the optimization process. However, policy averaging can be impractical for large-scale algorithms (e.g., deep RL) because it necessitates storing all past policies, which is often infeasible. On the other hand, Ying et al. (2022); Ding et al. (2024); Müller et al. (2024); Kitamura et al. (2024) tackled the issue by introducing entropy regularization into the objective return. However, the regularization can lead to biased solutions and result in a policy design that may deviate from what is intended by the cost function.

In contrast,  $\text{EpiRC-PGS}$  requires neither policy averaging nor regularization, thereby offering advantageous properties even in CMDP settings.

**Epigraph approach.** Few studies have investigated the epigraph form in the CMDP setting. So & Fan (2023) proposed a deep RL algorithm aimed at system stabilization under constraints, and So et al. (2024) developed a deep RL algorithm for goal-reaching tasks with risk-avoidance constraints. Although these studies empirically demonstrated the effectiveness of the epigraph form in constrained RL problems, they did not establish theoretical guarantees, such as global convergence. Moreover, they did not consider robust settings, whereas our `EpIRC-PGS` accounts for robustness.

On the other hand, this is the first work to provide theoretical guarantees for the epigraph form in CMDPs. Furthermore, unlike existing constrained RL studies, we consider not only constraints but also the robustness against the transition kernel.

## A.2 ROBUST MARKOV DECISION PROCESSES

RMDP is a specific subclass of RCMDP where there are no constraints, i.e.,  $N = 0$ . RMDP is a crucial research area for the practical success of RL applications, where the environmental mismatch between the training phase and the testing phase is almost unavoidable. Without robust policy design, even a small mismatch can lead to poor performance of the trained policy in the testing phase (Li et al., 2022; Jiang, 2018).

**Dynamic programming approach.** Since the seminal work by Iyengar (2005), numerous studies have explored dynamic programming (DP) approaches for RMDPs (Nilim & El Ghaoui, 2005; Clavier et al., 2023; Panaganti & Kalathil, 2022; Mai & Jaillet, 2021; Grand-Clément & Kroer, 2021; Derman et al., 2021; Wang & Zou, 2021; Kumar et al., 2022; Yang et al., 2023). The DP approach decomposes the original problem into smaller sub-problems using Bellman’s principle of optimality (Bellman et al., 1957). To apply this principle, DP approaches enforce rectangularity on the uncertainty set, which assumes independent worst-case transitions at each state or state-action pair. However, as pointed out by Goyal & Grand-Clement (2023), the rectangularity assumption can result in a very conservative optimal policy. Moreover, applying DP to constrained settings is challenging since CMDPs typically do not satisfy the principle of optimality (Haviv, 1996). Although several studies have attempted to apply DP to CMDPs, they face issues such as excessive memory consumption, due to the use of non-stationary policy classes, or are restricted to deterministic policy classes (Chang, 2023; Chen & Blankenship, 2004; Chen & Feinberg, 2006).

**Epigraph application to DP approach.** Chen et al. (2019); Wiesemann et al. (2013); Ho et al. (2021) employed the epigraph form to implement a robust DP algorithm for the  $s$ -rectangular uncertainty set setting. Specifically, they showed that the  $s$ -rectangular robust Bellman operator, which is the  $s$ -rectangular counterpart of Equation (19), can be efficiently implemented using the epigraph form. However, since their algorithms rely on Bellman’s principle of optimality, similar to standard DP, they are likely to encounter the same challenges in CMDP settings as those discussed above.

**Policy gradient approach.** Another promising approach for RMDPs is the policy gradient method. Similar to the DP approach, most existing policy gradient algorithms also work only under the rectangularity assumption (Kumar et al., 2024; Wang & Zou, 2022; Li et al., 2022), and thus suffer from the same conservativeness issue. It is important to note that robust policy evaluation can be NP-hard without any structural assumptions on the uncertainty set (Wiesemann et al., 2013), but such assumptions are potentially not required for the robust policy optimization step. Our policy gradient algorithm abstracts the evaluation step by Assumption 1 and avoids the need for the rectangularity assumption during the policy optimization phase, similar to the recent work by Wang et al. (2023).

## A.3 ROBUST CONSTRAINED MARKOV DECISION PROCESSES

Russel et al. (2020); Mankowitz et al. (2020) proposed heuristic algorithms for RCMDPs, but their approaches lack theoretical guarantees for convergence to a near-optimal policy. Wang et al. (2022) introduced a Lagrangian approach with convergence guarantee to a stationary point. However, they do not ensure the optimality of this stationary point. Moreover, their method is heavily dependent on the restrictive  $R$ -contamination set assumption (Du et al., 2018; Wang & Zou, 2021; 2022).

Sun et al. (2024) applied a trust-region method to RCMDPs. The policy is updated to remain sufficiently similar to the previous one, ensuring that performance and constraint adherence do not



864 degrade, even in the face of environmental uncertainty. However, while they ensure that each policy  
865 update step maintains performance, convergence to a near-optimal policy is not guaranteed.  
866

867 Ghosh (2024) employed a penalty approach which considers the optimization problem of the form  
868  $\min_{\pi} f(\pi) + \lambda \max\{h(\pi), 0\}$ , where  $f$  and  $h$  are defined in Section 1. While this approach can  
869 yield a near-optimal policy for a sufficiently large value of  $\lambda > 0$ , the author does not provide a  
870 concrete optimization method for the minimization and instead assumes the availability of an oracle  
871 to solve it. As we will demonstrate in Section 3, this minimization is intrinsically difficult, making  
872 the practical implementation of such an oracle challenging.

873 Finally, Zhang et al. (2024) tackled RCMDPs using the policy-mixing technique (Miryoosefi et al.,  
874 2019; Le et al., 2019). In this technique, a policy is sampled from a finite set of deterministic policies  
875 according to a sampling distribution at the start of each episode, and it remains fixed throughout  
876 the episode. However, even if a good sampling distribution is determined, there is no guarantee that  
877 the resulting expected policy will be optimal due to the non-convexity of the return function with  
878 respect to policies (Agarwal et al., 2021). We discuss the limitations of the policy-mixing technique  
879 in Appendix A.4. Additionally, Zhang et al. (2024) assume an  $R$ -contamination uncertainty set,  
880 limiting its applicability similarly to the work of Wang et al. (2022).

881 Although the RCMDP problem remains unsolved, the control theory community has long studied  
882 the computation of safe controllers under environmental uncertainties. Notable methods include  
883 robust model predictive control (Bemporad & Morari, 2007) and  $H_{\infty}$  optimal control (Anderson  
884 et al., 2019; Zames, 1981; Doyle, 1982). These approaches are specifically tailored for a specialized  
885 class of MDPs, known as the linear quadratic regulator (LQR, Du et al. (2021)). However, because  
886 LQR and tabular MDPs operate within distinct frameworks, these control methods are unsuitable  
887 for tabular RCMDPs. Given that most modern reinforcement learning (RL) algorithms, such as  
888 DQN (Mnih et al., 2015), are based on the tabular MDP framework, our results bridge the gap  
889 between the RL and control theory communities, laying the foundation for the development of  
890 reliable RL applications in the future.

#### 891 A.4 NOTES ON THE POLICY-MIXING TECHNIQUE

892 This section explains the theoretical limitations of the policy-mixing technique (Zhang et al., 2024;  
893 Miryoosefi et al., 2019; Le et al., 2019) for identifying a near-optimal policy.  
894

895 **Policy-mixing technique.** Let  $\tilde{\Pi} := \{\pi_1, \dots, \pi_m\}$  be a finite set of policies with  $m \in \mathbb{N}$ . Con-  
896 sider a non-robust, single-constraint CMDP  $(\mathcal{S}, \mathcal{A}, \gamma, P, \mathcal{C} = \{c_0, c_1\}, b, \mu)$ . Given a distribution  
897  $\rho \in \mathcal{P}(\tilde{\Pi})$ , define  
898

$$899 \tilde{J}_{c_0, P}(\rho) := \sum_{\pi \in \tilde{\Pi}} \rho(\pi) J_{c_0, P}(\pi) \quad \text{and} \quad \tilde{J}_{c_1, P}(\rho) := \sum_{\pi \in \tilde{\Pi}} \rho(\pi) J_{c_1, P}(\pi).$$

900 The policy-mixing technique considers the following optimization problem:  
901

$$902 \tilde{J}^* := \min_{\rho \in \mathcal{P}(\tilde{\Pi})} \tilde{J}_{c_0, P}(\rho) \quad \text{such that} \quad \tilde{J}_{c_1, P}(\rho) \leq b_1 \quad (14)$$

$$903 = \min_{\rho \in \mathcal{P}(\tilde{\Pi})} \max_{\lambda \in \mathbb{R}_+} \sum_{\pi \in \tilde{\Pi}} \rho(\pi) (J_{c_0, P}(\pi) + \lambda (J_{c_1, P}(\pi) - b_1)) =: \min_{\rho \in \mathcal{P}(\tilde{\Pi})} \max_{\lambda \in \mathbb{R}_+} \tilde{L}(\rho, \lambda).$$

904 Let  $\rho^*$  be the solution of Equation (14) such that  $\rho^* \in \arg \min_{\rho \in \mathcal{P}(\tilde{\Pi})} \max_{\lambda \in \mathbb{R}_+} \tilde{L}(\rho, \lambda)$ .  
905

906 In this setting, a policy is sampled from  $\rho$  at the start of each episode and remains fixed throughout  
907 the episode. The term  $\tilde{J}_{c_0, P}(\rho)$  represents the expected return under the distribution  $\rho$ . Since  
908  $\tilde{L}(\rho, \lambda)$  is convex in  $\rho$  and concave in  $\lambda$ , under some mild assumptions, Equation (14) can be  
909 solved efficiently by the following standard optimization procedure for min-max problems: At each  
910 iteration  $t = 1, \dots, T$ , with initial values  $\lambda^{(0)} \in \mathbb{R}_+$  and  $\rho^{(0)} \in \mathcal{P}(\tilde{\Pi})$ ,  
911  
912  
913  
914  
915  
916  
917

1. Update  $\lambda^{(t)}$  using a no-regret algorithm. For example, with gradient ascent and a learning rate  $\alpha > 0$ :

$$\lambda^{(t)} := \max\left\{\lambda^{(t-1)} + \alpha\left(\tilde{J}_{c_1,P}(\rho^{(t-1)}) - b_1\right), 0\right\}.$$

2. Update  $\rho^{(t)}$  as  $\rho^{(t)}(\pi) = \mathbb{1}\{\pi = \pi^{(t)}\}$  where

$$\pi^{(t)} \in \arg \min_{\pi \in \tilde{\Pi}} J_{c_0,P}(\pi) + \lambda^{(t)}(J_{c_1,P}(\pi) - b_1).$$

Then, the averaged distribution  $\bar{\rho}^{(T)} := \frac{1}{T} \sum_{t=0}^T \rho^{(t)}$  converges to  $\rho^*$  as  $T \rightarrow \infty$  (Abernethy & Wang, 2017; Zahavy et al., 2021). When  $\tilde{\Pi}$  is sufficiently large, we can expect that the optimal value of Equation (14) is equivalent to that of the CMDP problem, i.e.,  $\tilde{J}^* = J^*$ , where  $J^*$  is defined in Equation (2) with  $\mathcal{U} = \{P\}$ .

**Limitation of policy-mixing.** Even when  $\tilde{J}^* = J^*$ , it is crucial to note that, while  $\bar{\rho}^{(T)}$  converges to  $\rho^*$ , **there is no guarantee that  $\bar{\pi}^{(T)} := \frac{1}{T} \sum_{t=0}^T \pi^{(t)}$  will converge to  $\pi^*$ .**

Let  $\bar{\lambda}^{(T)} := \frac{1}{T} \sum_{t=0}^T \lambda^{(t)}$ . Zhang et al. (2024); Miryoosefi et al. (2019); Le et al. (2019) argued for the convergence of  $\bar{\pi}^{(T)}$  by asserting that the equality (a) in the following equation holds:

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \tilde{L}(\rho^{(t)}, \lambda^{(t)}) &= \frac{1}{T} \sum_{t=1}^T \left( J_{c_0,P}(\pi^{(t)}) + \lambda^{(t)} \left( J_{c_1,P}(\pi^{(t)}) - b_1 \right) \right) \\ &\stackrel{(a)}{=} J_{c_0,P}(\bar{\pi}^{(T)}) + \bar{\lambda}^{(T)} \left( J_{c_1,P}(\bar{\pi}^{(T)}) - b_1 \right) \end{aligned} \quad (15)$$

(see, for example, **Equation (14)** in Zhang et al. (2024), **Equation (1)** in Le et al. (2019), and around **Equation (13)** in Miryoosefi et al. (2019)).

However, (a) in Equation (15) does not hold in general because the return function is neither convex nor concave in policy. Even when  $T = 2$ , there is an example where Equation (15) fails (see **Proof of Lemma 3.1** in Agarwal et al. (2021)). This invalidates the results of Miryoosefi et al. (2019); Le et al. (2019); Zhang et al. (2024), thus illustrating the theoretical limitations of the policy-mixing approach for near-optimal policy identification.

## B DISCUSSION ON SINGLE-LOOP ALGORITHM

Although Algorithm 1 can identify a near-optimal policy, it uses a double-loop structure that repetitively solves  $\min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  by Algorithm 2. In practice, single-loop algorithms, such as primal-dual algorithms for CMDPs (e.g., Efroni et al. (2020); Ding et al. (2024)), are typically more efficient and preferable compared to double-loop algorithms. This section discusses the challenge of designing a single-loop algorithm for the epigraph form.

Since the epigraph form is a constrained optimization problem, we can further transform it using a Lagrangian multiplier  $\lambda \in \mathbb{R}_+$ , yielding:

$$J^* = \min_{b_0 \in [0, H]} \max_{\lambda \in \mathbb{R}_+} L_{\text{epi}}(b_0, \lambda) \quad \text{where} \quad L_{\text{epi}}(b_0, \lambda) := b_0 + \lambda \Delta_{b_0}^*. \quad (16)$$

Similar to the typical Lagrangian approach, let's swap the min-max order. We call the resulting formulation the ‘‘epigraph-Lagrange’’ formulation:

$$\text{(Epigraph-Lagrange)} \quad L_{\text{epi}}^* = \max_{\lambda \in \mathbb{R}_+} \min_{b_0 \in [0, H]} \min_{\pi \in \Pi} b_0 + \lambda \Delta_{b_0}(\pi). \quad (17)$$

Does the strong duality,  $J^* = L_{\text{epi}}^*$ , hold? If it does, we could design a single-loop algorithm similar to primal-dual CMDP algorithms, performing gradient ascent and descent on Equation (17). Unfortunately, proving the strong duality is challenging.

Essentially, the min-max can be swapped when  $L_{\text{epi}}(b_0, \lambda)$  in Equation (16) is quasiconvex-quasiconcave (Sion, 1958). While  $L_{\text{epi}}(b_0, \lambda)$  is clearly concave in  $\lambda$ , the quasiconvexity in  $b_0$  is

**Algorithm 3** Evaluators for Finite Uncertainty Set

---

```

1: Input: Policy  $\pi$ , uncertainty set  $\mathcal{U} = \{P_1, \dots, P_M\}$ , and index  $n \in \llbracket 0, N \rrbracket$ .
2: for  $m \in \llbracket 1, M \rrbracket$  do
3:    $Q_{c_n, P_m}^\pi := (I - \gamma P_m \Pi^\pi)^{-1} c_n \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ .
4:    $J_{c_n, P_m}(\pi) := \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mu(s) \pi(s, a) Q_{c_n, P_m}^\pi(s, a)$ 
5: end for
6: Let  $m^* \in \arg \max_{m \in \llbracket 1, M \rrbracket} J_{c_n, P_m}(\pi)$ 
7:  $d_{P_{m^*}}^\pi = (1 - \gamma) \mu (I - \gamma \Pi^\pi P_{m^*})^{-1} \in \mathbb{R}^{\mathcal{S}}$ 
8: return (for  $\widehat{J}_n$ ):  $J_{c_n, P_{m^*}}(\pi)$ 
9: return (for  $\widehat{J}_n^\partial$ ):  $H d_{P_{m^*}}^\pi(s) Q_{c_n, P_{m^*}}^\pi(s, a) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$ 

```

---

not obvious. Although  $\Delta_{b_0}^*$  is decreasing due to Lemma 2 and thus a quasi-convex function, there is no guarantee on the quasi-convexity of  $b_0 + \lambda \Delta_{b_0}^*$ . The situation would be resolved if  $\Delta_{b_0}^*$  were convex in  $b_0$ . However, since  $\Delta_{b_0}^* = \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  is a pointwise minimum and  $\Delta_{b_0}(\pi)$  may not be convex in  $\pi$  (Agarwal et al., 2021),  $\Delta_{b_0}^*$  may not be convex in  $b_0$  (Boyd & Vandenberghe, 2004).

Therefore, algorithms for the epigraph-Lagrange formulation face a problem similar to the  $\pi^*$  **solution challenge** of the Lagrangian formulation (Section 3). Proving strong duality or finding alternative ways to circumvent this challenge is a promising direction for future RCMDP research.

## C UNCERTAINTY SETS AND ALGORITHMS FOR $\Delta_{b_0}$ AND SUBGRADIENT EVALUATORS

This section provides examples of uncertainty set structures and algorithms that realize  $\widehat{J}_n$  in Assumption 1 and  $\widehat{J}_n^\partial$  in Assumption 3.  $\widehat{J}_n$  evaluates  $J_{c_n, \mathcal{U}}(\pi)$ , and  $\widehat{J}_n^\partial$  evaluates one of the elements in  $\partial J_{c_n, \mathcal{U}}(\pi)$ . In this section, we frequently use the following useful matrix formulations (Pierotta et al., 2013): for a cost  $c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ ,

$$\begin{aligned}
Q_{c, P}^\pi &= (I - \gamma P \Pi^\pi)^{-1} c \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}} \\
J_{c, P}(\pi) &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mu(s) \pi(s, a) Q_{c, P}^\pi(s, a) \in \mathbb{R} \\
d_P^\pi &= \mu^\top (I - \gamma \Pi^\pi P)^{-1} \in \mathbb{R}^{\mathcal{S}},
\end{aligned} \tag{18}$$

where  $I$  is an identity matrix and  $\Pi^\pi$  is a  $\mathbb{R}^{\mathcal{S} \times \mathcal{S} \times \mathcal{A}}$  matrix such that  $\Pi^\pi(s, (s, a)) = \pi(s, a)$ . Due to Lemma 1, an element in  $\partial J_{c_n, \mathcal{U}}(\pi)$  takes the form of:

$$(\nabla J_{c_n, P}(\pi))(s, a) = H d_P^\pi(s) Q_{c_n, P}^\pi(s, a) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

### C.1 FINITE UNCERTAINTY SET

Let  $\mathcal{U}$  be a finite uncertainty set such that  $\mathcal{U} = \{P_1, \dots, P_M\}$  where  $M \in \mathbb{N}$ .  $\mathcal{U}$  clearly satisfies Assumption 2. The implementation of  $\widehat{J}_n$  is trivial by Equation (18). The implementation of  $\widehat{J}_n^\partial$  is also straightforward due to Lemma 10. Specifically, it can be implemented by the following subgradient representation:

$$\partial J_{c_n, \mathcal{U}} = \text{conv} \left\{ \nabla J_{c_n, P_m}(\pi) \left| m \in \arg \max_{m \in \llbracket 1, M \rrbracket} J_{c_n, P_m}(\pi) \right. \right\}.$$

Algorithm 3 summarizes the implementations of  $\widehat{J}_n$  and  $\widehat{J}_n^\partial$ . Our experiment (a) in Section 6 uses Algorithm 3.

### C.2 $(s, a)$ -RECTANGULAR KL UNCERTAINTY SET

**$(s, a)$ -rectangularity.** An uncertainty set  $\mathcal{U}$  is called  $(s, a)$ -rectangular if it satisfies:

$$\mathcal{U} = \times_{s,a} \mathcal{U}_{s,a} \quad \text{where } \mathcal{U}_{s,a} \subseteq \mathcal{P}(\mathcal{S}).$$

**Algorithm 4** Evaluators for KL Uncertainty Set

- 
- 1: **Input:** Policy  $\pi$ , nominal transition kernel  $P$ , regularization parameter  $C'_{\text{KL}} > 0$ , and index  $n \in \llbracket 0, N \rrbracket$ .
  - 2: Repeat Equation (25) and compute its fixed point  $Q_{c_n}^{(\infty)}$ .
  - 3:  $P_n^*(\cdot | s, a) \propto P(\cdot | s, a) \exp\left(V_{c_n}^{(\infty)}(\cdot)/C'_{\text{KL}}\right)$  // Compute the worst-case environment
  - 4:  $Q_{c_n, P_n^*}^\pi := (I - \gamma P_n^* \Pi^\pi)^{-1} c_n \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ .
  - 5:  $J_{c_n, P_n^*}(\pi) := \sum_{(s, a) \in \mathcal{S} \times \mathcal{A}} \mu(s) \pi(s, a) Q_{c_n, P_n^*}^\pi(s, a)$
  - 6:  $d_{P_n^*}^\pi := (1 - \gamma) \mu (I - \gamma \Pi^\pi P_n^*)^{-1} \in \mathbb{R}^{\mathcal{S}}$
  - 7: **return (for  $\hat{J}_n$ ):**  $J_{c_n, P_n^*}(\pi)$
  - 8: **return (for  $\hat{J}_n^\partial$ ):**  $H d_{P_n^*}^\pi(s) Q_{c_n, P_n^*}^\pi(s, a) \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$
- 

For such  $(s, a)$ -rectangular uncertainty sets, the following **robust DP** update is a widely-used approach to compute the worst-case  $Q$ -function:

$$\begin{aligned}
 \text{(Robust DP)} \quad Q_{c_n}^{(t+1)}(s, a) &= c_n(s, a) + \gamma \max_{p \in \mathcal{U}_{s, a}} \sum_{s' \in \mathcal{S}} p(s') V_{c_n}^{(t)}(s') \\
 \text{where } V_{c_n}^{(t)}(s') &:= \sum_{a' \in \mathcal{A}} \pi(s', a') Q_{c_n}^{(t)}(s', a').
 \end{aligned} \tag{19}$$

By repeatedly applying Equation (19),  $Q_{c_n}^{(t)}$  converges linearly to  $Q_{c_n, P_n^*}^\pi$ , where  $P_n^* \in \arg \max_{P \in \mathcal{U}} J_{c_n, P}(\pi)$  (see, e.g., **Corollary 2** in Iyengar (2005)).

Once we have  $Q_{c_n, P_n^*}^\pi$ , thanks to the rectangularity, the worst-case environment can be computed by:

$$P_n^*(\cdot | s, a) = \arg \max_{p \in \mathcal{U}_{s, a}} \sum_{s' \in \mathcal{S}} p(s') V_{c_n, P_n^*}^\pi(s') \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \tag{20}$$

**KL uncertainty set.** Both Equation (19) and Equation (20) require an efficient computation of  $\max_{p \in \mathcal{U}_{s, a}} \langle p, v \rangle$  for some vector  $v \in \mathbb{R}^{\mathcal{S}}$ . The KL uncertainty set is one of the most popular choices, as it allows for efficient computation of this maximization (Iyengar, 2005; Yang et al., 2022).

For some positive constant  $C_{\text{KL}} > 0$ , if  $\mathcal{U}_{s, a}$  satisfies

$$\mathcal{U}_{s, a} = \{p \in \mathcal{P}(\mathcal{S}) \mid \text{KL}[p \parallel P(\cdot | s, a)] \leq C_{\text{KL}}\},$$

we call  $\mathcal{U}$  a  $(s, a)$ -rectangular KL uncertainty set. Here,  $\text{KL}[p \parallel q] = \sum_{s \in \mathcal{S}} p(s) \ln \frac{p(s)}{q(s)}$  denotes the KL divergence between  $p \in \mathcal{P}(\mathcal{S})$  and  $q \in \mathcal{P}(\mathcal{S})$ . Due to Lemma 1, this  $\mathcal{U}$  clearly satisfies Assumption 2.

The following lemma is useful for implementing the robust DP update with a KL uncertainty set.

**Lemma 4 (Lemma 4 in Iyengar (2005)).** *Let  $v \in \mathbb{R}^{\mathcal{S}}$  and  $\mathbf{0} < q \in \mathcal{P}(\mathcal{S})$ . The value of the optimization problem:*

$$\min_{p \in \mathcal{P}(\mathcal{S})} \langle p, v \rangle \text{ such that } \text{KL}[p \parallel q] \leq C_{\text{KL}} \tag{21}$$

is equal to

$$- \min_{\theta \geq 0} \theta \cdot C_{\text{KL}} + \theta \ln \left\langle q, \exp\left(-\frac{v}{\theta}\right) \right\rangle. \tag{22}$$

Let  $\theta^*$  be the solution of Equation (22). Then, the solution of Equation (21) is

$$p \propto q \exp\left(-\frac{v}{\theta^*}\right).$$

Using this lemma, Equation (19) can be implemented by

$$Q_{c_n}^{(t+1)}(s, a) = c_n(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{s,a}^*(s') V_{c_n}^{(t)}(s')$$

$$\text{where } P_{s,a}^* \propto P(\cdot | s, a) \exp\left(\frac{V_{c_n}^{(t)}(\cdot)}{\theta_{s,a}^*}\right) \quad (23)$$

$$\text{and } \theta_{s,a}^* := \arg \min_{\theta \geq 0} \theta \cdot C_{\text{KL}} + \theta \ln \left\langle P(\cdot | s, a), \exp\left(\frac{V_{c_n}^{(t)}(\cdot)}{\theta}\right) \right\rangle.$$

**Regularized alternative.** While Equation (22) is a convex optimization problem, solving it for all  $P(\cdot | s, a) \forall (s, a) \in \mathcal{S} \times \mathcal{A}$  in Equation (23) is computationally extensive in practice.

Rather than the exact constrained problem of Equation (21), Yang et al. (2023) proposed the following regularized robust DP update:

$$Q_{c_n}^{(t+1)}(s, a) = c_n(s, a) + \gamma \max_{p \in \mathcal{P}(\mathcal{S})} \left( \sum_{s' \in \mathcal{S}} p(s') V_{c_n}^{(t)}(s') - C'_{\text{KL}} \text{KL}[p \| P(\cdot | s, a)] \right), \quad (24)$$

where  $C'_{\text{KL}} > 0$  is a constant. This regularized form has the following efficient analytical solution:

$$Q_{c_n}^{(t+1)}(s, a) = c_n(s, a) + \gamma \sum_{s' \in \mathcal{S}} P_{s,a}^*(s') V_{c_n}^{(t)}(s')$$

$$\text{where } P_{s,a}^* \propto P(\cdot | s, a) \exp\left(\frac{V_{c_n}^{(t)}(\cdot)}{C'_{\text{KL}}}\right). \quad (25)$$

While this is a regularized approximation, the following lemma shows that Equation (25) solves some exact robust DP with a KL uncertainty set:

**Lemma 5** (Adaptation of **Proposition 3.1** and **Theorem 3.1** in Yang et al. (2023)). *For any  $C'_{\text{KL}} > 0$ , there exists  $C_{\text{KL}} > 0$  such that Equation (24) converges linearly to the fixed point of Equation (23).*

Algorithm 4 summarizes the regularized DP update to implement the algorithms. Our experiment (b) in Section 6 uses Algorithm 4.

## D EXPERIMENT DETAILS

**Environment construction.** For the settings with finite uncertainty sets (a), the parameters are set as  $S = 7$ ,  $A = 4$ ,  $\gamma = 0.995$ . We set  $\gamma = 0.99$  for the CMDP setting (c). For the setting with KL uncertainty set (b), we set  $S = 5$ ,  $A = 3$ , and  $\gamma = 0.99$ .

**EpiRC-PGS implementation.** For the policy gradient subroutine (Algorithm 2), we set the iteration length  $T$  and the learning rate  $\alpha$  to ensure that Assumption 4 is satisfied with a sufficiently small  $\varepsilon_{\text{opt}}$ . Specifically, for (a) and (c), we set the iteration length to  $T = 10^4$  and the learning rate to  $\alpha = 5 \times 10^{-5}$ . For (b), we set  $T = 10^3$  and  $\alpha = 5 \times 10^{-4}$ .

Notably, parameter tuning for EpiRC-PGS is straightforward, as any sufficiently large  $T$  and small  $\alpha$  should meet the conditions of Assumption 4. Since the initial policy in Algorithm 2 can be chosen arbitrarily, the  $(k - 1)$ -th policy from the outer loop is used as the initial policy for the  $k$ -th policy computation.

For the finite uncertainty set settings (a) and (c), we implement the evaluators  $\hat{J}_n$  and  $\hat{J}_n^\partial$  using Algorithm 3. For the KL uncertainty set setting (b), we implement  $\hat{J}_n$  and  $\hat{J}_n^\partial$  using Algorithm 4 with  $C'_{\text{KL}} = 2.0$ .

**Algorithm 5** Lagrangian Formulation Policy Gradient Search (LF)

---

```

1134 1: Input: Outer iteration length  $K \in \mathbb{N}$ , inner iteration length  $T \in \mathbb{N}$ , learning rate for Lagrangian
1135 multipliers  $\alpha_\lambda > 0$ , learning rate for policy  $\alpha_\pi > 0$ 
1136 2: Initialize the Lagrangian multipliers  $\lambda^{(0)} = \mathbf{0} \in \mathbb{R}^N$ 
1137 3: Set an arbitrary initial policy  $\pi^{(0)} \in \Pi$ 
1138 4: for  $k = 0, \dots, K - 1$  do
1139 5:   Set the initial policy  $\pi^{(k,0)} := \pi^{(k)}$  for the inner loop
1140 6:   for  $t = 0, \dots, T - 1$  do
1141 7:      $g^{(k,t)} \in \partial L_{\lambda^{(k)}}(\pi^{(k,t)})$  // Compute policy gradient
1142 8:      $\pi^{(k,t+1)} := \text{Proj}_\Pi(\pi^{(k,t)} - \alpha_\pi g^{(k,t)})$  // Policy update
1143 9:   end for
1144 10: Set the new policy:  $\pi^{(k+1)} := \pi^{(k,t^*)}$  where  $t^* \in \arg \min_{t \in \llbracket 0, T-1 \rrbracket} L_{\lambda^{(k)}}(\pi^{(k,t)})$ 
1145 11: Update Lagrangian multipliers:  $\lambda_n^{(k+1)} := \max\{\lambda_n^{(k)} + \alpha_\lambda (J_{n,\mathcal{U}}(\pi^{(k+1)}) - b_n), 0\}$  for all
1146  $n \in \llbracket 1, N \rrbracket$ 
1147 12: end for

```

---

**LF implementation.** The pseudocode for LF is shown in Algorithm 5. We set the iteration length and learning rate for the inner policy optimization to  $T = 10^4$  and  $\alpha_\pi = 5 \times 10^{-5}$  in **(a, c)**, and  $T = 10^3$  and  $\alpha_\pi = 5 \times 10^{-4}$  in **(b)**. Similar to EPiRC-PGS, these values are chosen to expect sufficient optimization in the inner loop. We choose  $\alpha_\lambda = 0.01$  from  $\{0.1, 0.01, 0.001\}$  for the outer updates, balancing between the convergence speed and performance.

## E ADDITIONAL DEFINITIONS

Throughout this section, let  $\mathcal{X}$  denote a set such that  $\mathcal{X} \subset \mathbb{R}^d$  with  $d \in \mathbb{N}$ .

**Definition 2** (Subgradient (Kruger, 2003)). Let  $\mathcal{X} \subset \mathbb{R}^d$  be an open set where  $d \in \mathbb{N}$ . The (Fréchet) subgradient of a function  $f : \mathcal{X} \rightarrow \mathbb{R}$  at a point  $x \in \mathcal{X}$  is defined as the set

$$\partial f(x) := \left\{ u \in \mathcal{X} \mid \liminf_{x' \rightarrow x, x' \neq x} \frac{f(x') - f(x) - \langle u, x' - x \rangle}{\|x' - x\|_2} \geq 0 \right\}.$$

Furthermore, if  $\partial f(x)$  is a singleton, its element is denoted as  $\nabla f(x)$  and called the (Fréchet) gradient of  $f$  at  $x$ .

**Definition 3** (Lipschitz continuity). Let  $\ell \geq 0$ . A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $\ell$ -Lipschitz if for any  $x_1, x_2 \in \mathcal{X}$ , we have that

$$\|f(x_1) - f(x_2)\|_2 \leq \ell \|x_1 - x_2\|_2.$$

**Definition 4** (Smoothness). Let  $\ell \geq 0$ . A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $\ell$ -smooth if for any  $x_1, x_2 \in \mathcal{X}$ , we have that

$$\|\nabla f(x_1) - \nabla f(x_2)\|_2 \leq \ell \|x_1 - x_2\|_2.$$

**Definition 5** (Weak convexity). Let  $\ell > 0$ . A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $\ell$ -weakly convex if for any  $g \in \partial f(x)$  and  $x, x' \in \mathcal{X}$ ,

$$f(x') - f(x) \geq \langle g, x' - x \rangle - \frac{\ell}{2} \|x' - x\|_2^2.$$

Note that  $f(x) + \frac{\ell}{2} \|x\|_2^2$  is convex in  $\mathcal{X}$  if and only if  $f$  is  $\ell$ -weakly convex.

**Definition 6** (Moreau envelope of a weakly convex function). Given a  $\ell$ -weakly convex function  $f : \mathcal{X} \rightarrow \mathbb{R}$  and a parameter  $0 < \tau < \ell^{-1}$ , the Moreau envelope function of  $f$  is given by  $M_\tau \circ f : \mathbb{R}^d \rightarrow \mathbb{R}$  such that

$$(M_\tau \circ f)(x) = \min_{x' \in \mathcal{X}} \left\{ f(x') + \frac{1}{2\tau} \|x - x'\|_2^2 \right\}.$$

## F USEFUL LEMMAS

Throughout this section,  $\mathcal{X}$  denotes a compact set such that  $\mathcal{X} \subset \mathbb{R}^d$ , where  $d \in \mathbb{N}$ .

**Lemma 6 (Lemma D.2, in Wang et al. (2023)).** *Let  $\ell \geq 0$  and  $h : \mathcal{X} \rightarrow \mathbb{R}$  be an  $\ell$ -smooth function. Then,  $h$  is a  $\ell$ -weakly convex function.*

**Lemma 7** (e.g., **Proposition 13.37** in Rockafellar & Wets (2009)). *Let  $f : \mathcal{X} \rightarrow \mathbb{R}$  be an  $\ell$ -weakly convex function, and let  $0 < \tau < \ell$  be a parameter. The Moreau envelope function  $M_\tau \circ f : \mathbb{R}^d \rightarrow \mathbb{R}$  is differentiable, and its gradient is given by*

$$\nabla(M_\tau \circ f)(x) = \frac{1}{\tau} \left( x - \arg \min_{x' \in \mathcal{X}} \left( f(x') + \frac{1}{2\tau} \|x - x'\|_2^2 \right) \right).$$

**Lemma 8** (Sion's minimax theorem (Sion, 1958)). *Let  $n, m \in \mathbb{N}$ . Let  $\mathcal{X} \subset \mathbb{R}^n$  be a compact convex set and  $\mathcal{Y} \subset \mathbb{R}^m$  a convex set. Suppose that  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  satisfies the following two properties:*

- $f(x, \cdot)$  is upper semicontinuous and quasi-concave on  $\mathcal{Y}$  for any  $x \in \mathcal{X}$ .
- $f(\cdot, y)$  is lower semicontinuous and quasi-convex on  $\mathcal{X}$  for any  $y \in \mathcal{Y}$ .

Then,  $\min_{x \in \mathcal{X}} \sup_{y \in \mathcal{Y}} f(x, y) = \sup_{y \in \mathcal{Y}} \min_{x \in \mathcal{X}} f(x, y)$ .

**Lemma 9** (e.g., **Problem 9.13, Page 99** in Clarke et al. (2008)). *Let  $\mathcal{Y} \subset \mathbb{R}^m$  be a compact set and  $f : \mathbb{R}^d \times \mathcal{Y} \rightarrow \mathbb{R}$  be a continuous function of two arguments. Consider a point  $\bar{x} \in \mathbb{R}^d$  and let  $\Omega(\bar{x}) \subset \mathbb{R}^d$  be its neighborhood. For any  $(x, y) \in \Omega(\bar{x}) \times \mathcal{Y}$ , suppose that the gradient  $\nabla_x f(x, y)$  exists and is jointly continuous.*

Let  $h(\bar{x}) := \max_{y \in \mathcal{Y}} f(\bar{x}, y)$ . Then, the subgradient of  $h$  at  $\bar{x}$  is given by

$$\partial h(\bar{x}) = \text{conv} \left\{ \nabla_x f(\bar{x}, y) \mid y \in \arg \max_{y \in \mathcal{Y}} f(\bar{x}, y) \right\}.$$

**Lemma 10.** *Let  $N \in \mathbb{N}$ . Let  $f_i : \mathcal{X} \rightarrow \mathbb{R}$  for  $i \in \llbracket 1, N \rrbracket$  be  $\ell$ -weakly convex functions for some  $\ell \geq 0$ . Define the pointwise maximum function  $f : \mathcal{X} \rightarrow \mathbb{R}$  as*

$$f(x) = \max\{f_1(x), \dots, f_N(x)\} \quad \forall x \in \mathcal{X}.$$

Then, for any  $x \in \mathcal{X}$ ,

$$\partial f(x) = \text{conv} \{g \in \mathbb{R}^d \mid g \in \partial f_i(x), f_i(x) = f(x)\}.$$

*Proof.* The claim directly follows from **Theorem 1.3** and **Theorem 1.5** in Mikhalevich et al. (2024).  $\square$

**Lemma 11** (Maximum difference inequality). *Let  $N \in \mathbb{N}$ . For two sets of real numbers  $\{x_i\}_{i \in \llbracket 1, N \rrbracket}$  and  $\{y_i\}_{i \in \llbracket 1, N \rrbracket}$ , where  $x_i, y_i \in \mathbb{R}$ ,*

$$\left| \max_{i \in \llbracket 1, N \rrbracket} x_i - \max_{i' \in \llbracket 1, N \rrbracket} y_{i'} \right| \leq \max_{i \in \llbracket 1, N \rrbracket} |x_i - y_i|.$$

*Proof.* For any  $i \in \llbracket 1, N \rrbracket$ ,

$$\begin{aligned} \max_{i \in \llbracket 1, N \rrbracket} x_i &= \max_{i \in \llbracket 1, N \rrbracket} x_i - y_i + y_i \leq \max_{i \in \llbracket 1, N \rrbracket} (x_i - y_i) + \max_{i' \in \llbracket 1, N \rrbracket} y_{i'} \\ \implies \max_{i \in \llbracket 1, N \rrbracket} x_i - \max_{i' \in \llbracket 1, N \rrbracket} y_{i'} &\leq \max_{i \in \llbracket 1, N \rrbracket} x_i - y_i. \end{aligned}$$

By the symmetry of  $x_i$  and  $y_i$ , we have  $\max_{i \in \llbracket 1, N \rrbracket} y_i - \max_{i' \in \llbracket 1, N \rrbracket} x_{i'} \leq \max_{i \in \llbracket 1, N \rrbracket} y_i - x_i$ . Therefore,

$$\left| \max_{i \in \llbracket 1, N \rrbracket} x_i - \max_{i' \in \llbracket 1, N \rrbracket} y_{i'} \right| \leq \max_{i \in \llbracket 1, N \rrbracket} |x_i - y_i|.$$

$\square$

**Lemma 12** (Point-wise maximum preserves weak convexity). *Let  $h : \mathcal{X} \rightarrow \mathbb{R}$  and  $f : \mathcal{X} \rightarrow \mathbb{R}$  be  $\ell_h$ - and  $\ell_f$ -weakly convex functions, respectively. Then,  $g : \mathcal{X} \rightarrow \mathbb{R}$  defined by  $g(x) = \max\{h(x), f(x)\}$  for any  $x \in \mathcal{X}$  is  $\ell$ -weakly convex, where  $\ell := \max\{\ell_h, \ell_f\}$ .*

*Proof.* By the definition of weak convexity, for any  $\theta \in [0, 1]$  and  $x, y \in \mathcal{X}$ ,

$$h(\theta x + (1 - \theta)y) + \frac{\ell_h}{2} \|\theta x + (1 - \theta)y\|_2^2 \leq \theta \left( h(x) + \frac{\ell_h}{2} \|x\|_2^2 \right) + (1 - \theta) \left( h(y) + \frac{\ell_h}{2} \|y\|_2^2 \right).$$

A similar inequality holds for  $f$ . Then,

$$\begin{aligned} & g(\theta x + (1 - \theta)y) + \frac{\ell}{2} \|\theta x + (1 - \theta)y\|_2^2 \\ &= \max\{h(\theta x + (1 - \theta)y), f(\theta x + (1 - \theta)y)\} + \frac{\ell}{2} \|\theta x + (1 - \theta)y\|_2^2 \\ &= \max \left\{ h(\theta x + (1 - \theta)y) + \frac{\ell}{2} \|\theta x + (1 - \theta)y\|_2^2, f(\theta x + (1 - \theta)y) + \frac{\ell}{2} \|\theta x + (1 - \theta)y\|_2^2 \right\} \\ &\leq \max \left\{ \theta h(x) + (1 - \theta)h(y) + \frac{\ell}{2} (\theta \|x\|_2^2 + (1 - \theta)\|y\|_2^2), \theta f(x) + (1 - \theta)f(y) + \frac{\ell}{2} (\theta \|x\|_2^2 + (1 - \theta)\|y\|_2^2) \right\} \\ &= \max\{\theta h(x) + (1 - \theta)h(y), \theta f(x) + (1 - \theta)f(y)\} + \frac{\ell}{2} (\theta \|x\|_2^2 + (1 - \theta)\|y\|_2^2) \\ &\leq \theta \left( \max\{h(x), f(x)\} + \frac{\ell}{2} \|x\|_2^2 \right) + (1 - \theta) \left( \max\{h(y), f(y)\} + \frac{\ell}{2} \|y\|_2^2 \right) \\ &= \theta \left( g(x) + \frac{\ell}{2} \|x\|_2^2 \right) + (1 - \theta) \left( g(y) + \frac{\ell}{2} \|y\|_2^2 \right). \end{aligned}$$

Therefore,  $g$  is  $\ell$ -weakly convex.  $\square$

**Lemma 13.** *Let  $N_{\mathcal{X}}(x)$  be the normal cone of  $\mathcal{X}$  at  $x \in \mathcal{X}$ , defined as*

$$N_{\mathcal{X}}(x) := \{g \in \mathbb{R}^d \mid \langle g, y \rangle \leq \langle g, x \rangle \quad \forall y \in \mathcal{X}\}.$$

*Define the indicator function  $\mathbb{I}_{\mathcal{X}} : \mathbb{R}^d \rightarrow \mathbb{R}$  such that*

$$\mathbb{I}_{\mathcal{X}}(x) = \begin{cases} 0 & \text{if } x \in \mathcal{X} \\ \infty & \text{otherwise} \end{cases}.$$

*Then,  $\partial \mathbb{I}_{\mathcal{X}}(x) = N_{\mathcal{X}}(x)$  for any  $x \in \mathcal{X}$ .*

*Proof.* Note that any  $g \in \partial \mathbb{I}_{\mathcal{X}}(x)$  satisfies

$$\mathbb{I}_{\mathcal{X}}(y) \geq \mathbb{I}_{\mathcal{X}}(x) + \langle g, y - x \rangle \quad \forall y \in \mathbb{R}^d. \quad (26)$$

Suppose that  $g \notin N_{\mathcal{X}}(x)$ . Then, there exists  $y' \in \mathcal{X}$  such that  $\langle g, x \rangle < \langle g, y' \rangle$ , which contradicts Equation (26). Therefore,  $g \in N_{\mathcal{X}}(x)$  for any  $g \in \partial \mathbb{I}_{\mathcal{X}}(x)$  and thus  $\partial \mathbb{I}_{\mathcal{X}}(x) \subseteq N_{\mathcal{X}}(x)$ .

Consider  $g \in N_{\mathcal{X}}(x)$ . It satisfies  $0 \geq \langle g, y - x \rangle$  for any  $y \in \mathcal{X}$ . Since  $x \in \mathcal{X}$  and by the definition of  $\mathbb{I}_{\mathcal{X}}$ , Equation (26) holds for any  $y \in \mathbb{R}^d$ . Therefore,  $N_{\mathcal{X}}(x) \subseteq \partial \mathbb{I}_{\mathcal{X}}(x)$ . This concludes the proof.  $\square$

**Lemma 14.** *Let  $h : \mathcal{X} \rightarrow \mathbb{R}$  be an  $\ell$ -weakly convex function. For  $0 < \tau < 1/\ell$ , define*

$$\bar{x}_{\tau} \in \arg \min_{x' \in \mathcal{X}} h(x') + \frac{1}{2\tau} \|x - x'\|_2^2.$$

*Then, there exists a subgradient  $g \in \partial h(\bar{x}_{\tau})$  such that, for any  $y \in \mathcal{X}$ ,*

$$\langle g, \bar{x}_{\tau} - y \rangle \leq \langle \nabla(\mathbb{M}_{\tau} \circ h)(x), \bar{x}_{\tau} - y \rangle$$



1296 *Proof.* Let  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  be a function such that  $f(x) = h(x) + \mathbb{I}_{\mathcal{X}}(x)$ .

1297 The Moreau envelope function of  $f$  satisfies that, for any  $x \in \mathbb{R}^d$ ,

$$1298 \quad (\mathbb{M}_\tau \circ f)(x) = \min_{x' \in \mathbb{R}^d} \left\{ h(x') + \mathbb{I}_{\mathcal{X}}(x') + \frac{1}{2\tau} \|x - x'\|_2^2 \right\} = \min_{x' \in \mathcal{X}} \left\{ h(x') + \frac{1}{2\tau} \|x - x'\|_2^2 \right\}.$$

1300 It holds that  $\nabla(\mathbb{M}_\tau \circ f)(x) = \frac{1}{\tau}(x - \bar{x}_\tau)$  due to Lemma 7.

1301 Note that

$$1302 \quad \bar{x}_\tau \in \arg \min_{x' \in \mathcal{X}} h(x') + \frac{1}{2\tau} \|x - x'\|_2^2 = \arg \min_{x' \in \mathbb{R}^d} h(x') + \mathbb{I}_{\mathcal{X}}(x') + \frac{1}{2\tau} \|x - x'\|_2^2.$$

1303 It is clear that  $\bar{x}_\tau$  is a minimizer of the function  $\phi_x(x') := h(x') + \mathbb{I}_{\mathcal{X}}(x') + \frac{1}{2\tau} \|x - x'\|_2^2$ .  
1304 Therefore, it holds that  $\mathbf{0} \in \partial\phi_x(\bar{x}_\tau)$ . Accordingly,

$$1305 \quad \mathbf{0} \in \partial \left( h(y) + \mathbb{I}_{\mathcal{X}}(y) + \frac{1}{2\tau} \|x - y\|_2^2 \right) \Big|_{y=\bar{x}_\tau} \implies -\frac{1}{\tau}(\bar{x}_\tau - x) \in \partial(h(y) + \mathbb{I}_{\mathcal{X}}(y)) \Big|_{y=\bar{x}_\tau}.$$

1306 Due to Lemma 13,  $\partial\mathbb{I}_{\mathcal{X}}(x) = N_{\mathcal{X}}(x)$ . Therefore, there exists a subgradient  $g \in \partial h(\bar{x}_\tau)$  such that

$$1307 \quad -g - \frac{1}{\tau}(\bar{x}_\tau - x) \in N_{\mathcal{X}}(\bar{x}_\tau).$$

1308 Since any  $z \in N_{\mathcal{X}}(\bar{x}_\tau)$  satisfies  $\langle z, y - \bar{x}_\tau \rangle \leq 0$  for any  $y \in \mathcal{X}$ , it holds that

$$1309 \quad \langle -g, y - \bar{x}_\tau \rangle \leq \left\langle \frac{1}{\tau}(\bar{x}_\tau - x), y - \bar{x}_\tau \right\rangle, \quad \forall y \in \mathcal{X}.$$

1310 Then the claim follows from the fact that  $\frac{1}{\tau}(x - \bar{x}_\tau) = \nabla(\mathbb{M}_\tau \circ h)(x)$  due to Lemma 7.  $\square$

1311 **Lemma 15** (Linear optimization on convex hull). *Given  $c \in \mathbb{R}^d$  and a compact set  $\mathcal{X} \subset \mathbb{R}^d$ , it holds that*

$$1312 \quad \min_{x \in \mathcal{X}} \langle c, x \rangle = \min_{x \in \text{conv}\{\mathcal{X}\}} \langle c, x \rangle.$$

1313 *Proof.* Let  $x^* \in \arg \min_{x \in \text{conv}\{\mathcal{X}\}} \langle c, x \rangle$ . The claim holds for  $x^* \in \mathcal{X}$ . Suppose that  $x^* \notin \mathcal{X}$ .  
1314 Then, by the definition of the convex hull, there exist  $y, z \in \mathcal{X}$  and  $\theta \in (0, 1)$  such that  $y \neq z$  and

$$1315 \quad x^* = \theta y + (1 - \theta)z.$$

1316 Since  $x^*$  is a minimizer, we have

$$1317 \quad \langle c, x^* \rangle \leq \langle c, y \rangle \quad \text{and} \quad \langle c, x^* \rangle \leq \langle c, z \rangle.$$

1318 Accordingly,

$$1319 \quad \langle c, x^* \rangle = \theta \langle c, x^* \rangle + (1 - \theta) \langle c, x^* \rangle \leq \theta \langle c, y \rangle + (1 - \theta) \langle c, z \rangle = \langle c, x^* \rangle.$$

1320 The inequality must be an equality, and thus

$$1321 \quad \theta \underbrace{(\langle c, y \rangle - \langle c, x^* \rangle)}_{\geq 0} + (1 - \theta) \underbrace{(\langle c, z \rangle - \langle c, x^* \rangle)}_{\geq 0} = 0.$$

1322 Since  $\theta \in (0, 1)$ , it holds that

$$1323 \quad \langle c, y \rangle = \langle c, z \rangle = \langle c, x^* \rangle.$$

1324 The above equality means that both  $y$  and  $z \in \mathcal{X}$  satisfy  $\langle c, y \rangle = \langle c, z \rangle = \min_{x \in \text{conv}\{\mathcal{X}\}} \langle c, x \rangle$ .  
1325 Therefore,  $\min_{x \in \mathcal{X}} \langle c, x \rangle = \min_{x \in \text{conv}\{\mathcal{X}\}} \langle c, x \rangle$ .  $\square$

## G USEFUL LEMMAS FOR MDPs

**Lemma 16** (**Lemma 3.1** in Wang et al. (2023)). *Let*

$$\ell_{\text{LP}} := H^2\sqrt{A} \quad \text{and} \quad \ell_{\text{sm}} := 2\gamma AH^3 .$$

For any  $\pi, \pi' \in \Pi$ ,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ ,  $\mu \in \mathcal{P}(\mathcal{S})$ , and  $c \in [0, 1]^{\mathcal{S} \times \mathcal{A}}$ ,

$$|J_{c,P}(\pi) - J_{c,P}(\pi')| \leq \ell_{\text{LP}} \|\pi - \pi'\|_2, \quad \|\nabla J_{c,P}(\pi) - \nabla J_{c,P}(\pi')\|_2 \leq \ell_{\text{sm}} \|\pi - \pi'\|_2,$$

$$\text{and } |J_{c,\mathcal{U}}(\pi) - J_{c,\mathcal{U}}(\pi')| \leq \ell_{\text{LP}} \|\pi - \pi'\|_2 .$$

Furthermore,  $J_{c,P}(\pi)$  is  $\ell_{\text{sm}}$ -weakly convex in  $\Pi$ , as follows directly from Lemma 6.

**Lemma 17** (e.g., **Lemma 4.1** in Agarwal et al. (2021) and **Lemma E.2** in Wang et al. (2023)). *Let*  $\mu \in \mathcal{P}(\mathcal{S})$  such that  $\min_{s \in \mathcal{S}} \mu(s) > 0$ . For any  $\pi \in \Pi$ ,  $P : \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{P}(\mathcal{S})$ , and  $c \in [0, 1]^{\mathcal{S} \times \mathcal{A}}$ ,

$$J_{c,P}(\pi) - J_{c,P}(\pi_{c,P}^*) \leq H \left\| \frac{d_P^{\pi_{c,P}^*}}{\mu} \right\|_{\infty} \max_{\pi' \in \Pi} \langle \pi - \pi', \nabla J_{c,P}(\pi) \rangle ,$$

where  $\pi_{c,P}^* \in \arg \min_{\pi' \in \Pi} J_{c,P}(\pi')$ .

## H PROOF OF THEOREM 1

*Proof of Theorem 1.* Consider the deterministic RCMDP shown in Figure 1a with  $N = 1$ ,  $\mathcal{U} = \{P_1, P_2\}$ ,  $\mathcal{S} = \{s_1, s_2, s_3, s_4\}$ , and  $\mathcal{A} = \{a_1, a_2\}$ . Set the initial distribution such that  $\mu(s_1) = \mu(s_2) = \mu(s_3) = \mu(s_4) = 1/4$ .

**First part of Theorem 1.** We set  $\lambda = 1$ . The threshold  $b_1$  can be arbitrary.

Let  $\pi_1$  and  $\pi_2$  be two policies such that  $\pi_1$  always chooses  $a_1$  and  $\pi_2$  always chooses  $a_2$  in any state. For any  $\delta > 0$ , we will show two results:

- Equation (28):  $L_\lambda(\pi_2) - \min_{\pi \in \Pi} L_\lambda(\pi) \geq \frac{H\gamma}{4} - \frac{3H\delta}{4}$ .
- Equation (30):  $(\nabla L_\lambda(\pi_2))(\cdot, a_1) > (\nabla L_\lambda(\pi_2))(\cdot, a_2)$ .

The former shows the suboptimality of  $\pi_2$ , and the latter indicates that  $\pi_2$  is a local minimum.

According to the RCMDP construction, for any  $\pi \in \Pi$ , we have

$$\mu(s_3)V_{c_0,P_1}^\pi(s_3) + \mu(s_4)V_{c_0,P_1}^\pi(s_4) = \frac{H}{4}(1 + \gamma),$$

$$\mu(s_3)V_{c_0,P_2}^\pi(s_3) + \mu(s_4)V_{c_0,P_2}^\pi(s_4) = \frac{H}{4}(1 - \gamma),$$

$$\mu(s_3)V_{c_1,P_1}^\pi(s_3) + \mu(s_4)V_{c_1,P_1}^\pi(s_4) = \frac{H}{4}(1 - \gamma),$$

$$\mu(s_3)V_{c_1,P_2}^\pi(s_3) + \mu(s_4)V_{c_1,P_2}^\pi(s_4) = \frac{H}{4}(1 + \gamma).$$

For  $\pi_1$  and  $\pi_2$ , it is easy to verify that

$$\mu(s_1)V_{c_0,P_1}^{\pi_1}(s_1) + \mu(s_2)V_{c_0,P_1}^{\pi_1}(s_2) = \frac{1}{4}(\delta + \gamma + \gamma^2\delta + \dots) + \frac{1}{4}(1 + \gamma\delta + \gamma^2 + \dots) = \frac{H}{4}(1 + \delta),$$

$$\mu(s_1)V_{c_0,P_2}^{\pi_1}(s_1) + \mu(s_2)V_{c_0,P_2}^{\pi_1}(s_2) = \frac{1}{4}(\delta + \gamma + \gamma^2 + \dots) + \frac{1}{4}(1 + \gamma + \gamma^2 + \dots) = \frac{H}{4}(1 + \gamma) + \frac{\delta}{4},$$

$$\mu(s_1)V_{c_1,P_1}^{\pi_1}(s_1) + \mu(s_2)V_{c_1,P_1}^{\pi_1}(s_2) = \frac{H}{2}, \quad \mu(s_1)V_{c_1,P_2}^{\pi_1}(s_1) + \mu(s_2)V_{c_1,P_2}^{\pi_1}(s_2) = \frac{H}{2},$$

$$\mu(s_1)V_{c_0,P_1}^{\pi_2}(s_1) + \mu(s_2)V_{c_0,P_1}^{\pi_2}(s_2) = \frac{H}{2}, \quad \mu(s_1)V_{c_0,P_2}^{\pi_2}(s_1) + \mu(s_2)V_{c_0,P_2}^{\pi_2}(s_2) = \frac{H}{2},$$

$$\mu(s_1)V_{c_1,P_1}^{\pi_2}(s_1) + \mu(s_2)V_{c_1,P_1}^{\pi_2}(s_2) = \frac{H}{4}(1 + \gamma - 2\delta),$$

$$\mu(s_1)V_{c_1,P_2}^{\pi_2}(s_1) + \mu(s_2)V_{c_1,P_2}^{\pi_2}(s_2) = \frac{H}{4}(1 + \gamma - 2\delta).$$

Therefore,

$$\begin{aligned}
J_{c_0, P_1}(\pi_1) &= \frac{H}{2} + \frac{H}{4}(\gamma + \delta), & J_{c_0, P_2}(\pi_1) &= \frac{H}{2} + \frac{\delta}{4}, \\
J_{c_1, P_1}(\pi_1) &= \frac{H}{4}(3 - \gamma), & J_{c_1, P_2}(\pi_1) &= \frac{H}{4}(3 + \gamma), \\
J_{c_0, P_1}(\pi_2) &= \frac{H}{4}(3 + \gamma), & J_{c_0, P_2}(\pi_2) &= \frac{H}{4}(3 - \gamma), \\
J_{c_1, P_1}(\pi_2) &= \frac{H}{2} - \frac{H\delta}{2}, & J_{c_1, P_2}(\pi_2) &= \frac{H}{2} + \frac{H\gamma}{2} - \frac{H\delta}{2},
\end{aligned} \tag{27}$$

Hence,

$$\begin{aligned}
J_{c_0, \mathcal{U}}(\pi_1) &= J_{c_0, P_1}(\pi_1) = \frac{H}{2} + \frac{H\gamma}{4} + \frac{H\delta}{4}, & J_{c_1, \mathcal{U}}(\pi_1) &= J_{c_1, P_2}(\pi_1) = \frac{H}{4}(3 + \gamma), \\
J_{c_0, \mathcal{U}}(\pi_2) &= J_{c_0, P_1}(\pi_2) = \frac{H}{4}(3 + \gamma), & J_{c_1, \mathcal{U}}(\pi_2) &= J_{c_1, P_2}(\pi_2) = \frac{H}{2} + \frac{H\gamma}{2} - \frac{H\delta}{2}.
\end{aligned}$$

Accordingly, since  $\lambda = 1$ , we have

$$L_\lambda(\pi_2) - \min_{\pi \in \Pi} L_\lambda(\pi) \geq L_\lambda(\pi_2) - L_\lambda(\pi_1) \geq \frac{H\gamma}{4} - \frac{3H\delta}{4}. \tag{28}$$

The next task is to show that  $(\nabla L_\lambda(\pi_2))(\cdot, a_1) > (\nabla L_\lambda(\pi_2))(\cdot, a_2)$ .

By using Lemma 10, it is easy to show that

$$\nabla L_\lambda(\pi_2) = \nabla J_{c_0, P_1}(\pi_2) + \nabla J_{c_1, P_2}(\pi_2).$$

Since  $d_{P_1}^{\pi_2}(s) = d_{P_2}^{\pi_2}(s) = 0.25\mathbf{1}$  and due to Lemma 1, we have

$$\frac{4}{H} \nabla L_\lambda(\pi_2) = Q_{c_0, P_1}^{\pi_2} + Q_{c_1, P_2}^{\pi_2}.$$

Note that

$$\begin{aligned}
\frac{4}{H} (\nabla L_\lambda(\pi_2))(s_1, a_1) &= Q_{c_0, P_1}^{\pi_2}(s_1, a_1) + Q_{c_1, P_2}^{\pi_2}(s_1, a_1) = (\delta + H\gamma) + (H - H\gamma\delta), \\
\frac{4}{H} (\nabla L_\lambda(\pi_2))(s_1, a_2) &= Q_{c_0, P_1}^{\pi_2}(s_1, a_2) + Q_{c_1, P_2}^{\pi_2}(s_1, a_2) = H + H(\gamma - \delta), \\
\frac{4}{H} (\nabla L_\lambda(\pi_2))(s_2, a_1) &= Q_{c_0, P_1}^{\pi_2}(s_2, a_1) + Q_{c_1, P_2}^{\pi_2}(s_2, a_1) = 2H, \\
\frac{4}{H} (\nabla L_\lambda(\pi_2))(s_2, a_2) &= Q_{c_0, P_1}^{\pi_2}(s_2, a_2) + Q_{c_1, P_2}^{\pi_2}(s_2, a_2) = 2H(1 - \delta).
\end{aligned} \tag{29}$$

Therefore, since  $\delta > 0$ ,

$$\begin{aligned}
\frac{4}{H} ((\nabla L_\lambda(\pi_2))(s_1, a_1) - (\nabla L_\lambda(\pi_2))(s_1, a_2)) &= \delta - H\gamma\delta + H\delta = 2\delta > 0, \\
\frac{4}{H} ((\nabla L_\lambda(\pi_2))(s_2, a_1) - (\nabla L_\lambda(\pi_2))(s_2, a_2)) &= 2H\delta > 0.
\end{aligned} \tag{30}$$

Now, with a sufficiently small  $R > 0$ , let  $\tilde{\Pi}_2 := \{\pi \in \Pi \mid \|\pi - \pi_2\|_2 \leq R, \pi \neq \pi_2\}$  be policies near  $\pi_2$ . When  $R$  is sufficiently small, due to the Lipschitz continuity of  $J_{c_n, P}(\pi)$  by Lemma 16, Equation (27) indicates that

$$J_{c_0, \mathcal{U}}(\pi) = J_{c_0, P_1}(\pi) \text{ and } J_{c_1, \mathcal{U}}(\pi) = J_{c_1, P_2}(\pi) \quad \forall \pi \in \tilde{\Pi}_2. \tag{31}$$

Similarly, due to Equation (31) with the Lipschitz continuity of  $\nabla J_{n, P}(\pi)$  by Lemma 16, Equation (29) and Equation (30) indicate that,

$$(\nabla L_\lambda(\pi))(\cdot, a_1) > (\nabla L_\lambda(\pi))(\cdot, a_2) \quad \forall \pi \in \tilde{\Pi}_2.$$

Therefore, since  $\pi_2$  always chooses  $a_2$ , we have  $L_\lambda(\pi_2) < L_\lambda(\pi) \quad \forall \pi \in \tilde{\Pi}_2$  for a sufficiently small  $R > 0$ . The first part of the claim holds by setting  $\delta = \gamma/4$  with Equation (28).

**Second part of Theorem 1.** Consider again the deterministic RCMDP given in the previous part of the proof with  $\delta = \gamma/4$ . For a value  $b_1 \in \mathbb{R}$ , define a function  $\Psi_{b_1} : \mathbb{R} \rightarrow \mathbb{R}$  such that

$$\Psi_{b_1}(\lambda) = \min_{\pi \in \Pi} J_{c_0, \mathcal{U}}(\pi) + \lambda J_{c_1, \mathcal{U}}(\pi) - \lambda b_1 = \min_{\pi \in \Pi} L_\lambda(\pi).$$

We first show that, when  $b_1$  ranges from 0 to  $H$ , the  $\arg \sup_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$  ranges from  $\infty$  to 0.

Since  $\mathcal{U} = \{P_1, P_2\}$ , Lemma 10 and Lemma 1 indicates that, for any  $\lambda \in \mathbb{R}$  and  $b_1 \in \mathbb{R}$ ,

$$\partial_\lambda \Psi_{b_1}(\lambda) \subseteq \text{conv}\{J_{c_1, P}(\pi) - b_1 \mid \pi \in \Pi, P \in \{P_1, P_2\}\}. \quad (32)$$

Since  $\mu = \frac{1}{4} \cdot \mathbf{1}$  and due to the construction of the RCMDP in Figure 1a, it is easy to verify that,

$$\begin{aligned} \min_{\pi \in \Pi} \min_{P \in \mathcal{U}} J_{c_1, P}(\pi) &\geq \min_{\pi \in \Pi} \min_{P \in \mathcal{U}} \min_{s \in \{s_3, s_4\}} \frac{1}{4} V_{c_1, P}^\pi(s) = \frac{1}{4}, \\ \max_{\pi \in \Pi} \max_{P \in \mathcal{U}} J_{c_1, P}(\pi) &\leq \frac{H}{2} + \max_{\pi \in \Pi} \max_{P \in \mathcal{U}} \frac{1}{4} (V_{c_1, P}^\pi(s_3) + V_{c_1, P}^\pi(s_4)) = H - \frac{1}{4}. \end{aligned}$$

By inserting this to Equation (32), for any  $\lambda \in \mathbb{R}$  and  $b_1 \in \mathbb{R}$ , we have

$$\frac{1}{4} - b_1 \leq g \leq H - \frac{1}{4} - b_1 \quad \forall g \in \partial_\lambda \Psi_{b_1}(\lambda). \quad (33)$$

Therefore,

- When  $b_1 \in [0, 1/4)$ ,  $g \in \partial_\lambda \Psi_{b_1}(\lambda)$  must satisfy  $g > 0$  for any  $\lambda$ . Thus,  $\{\infty\} = \arg \sup_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$  for any  $b_1 \in [0, 1/4)$ . Moreover,  $\sup_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda) = \infty$ .
- When  $b_1 \in (H - 1/4, H]$ ,  $g \in \partial_\lambda \Psi_{b_1}(\lambda)$  must satisfy  $g < 0$  for any  $\lambda$ . Thus,  $\{0\} = \arg \max_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$  for any  $b_1 \in (H - 1/4, H]$ . Moreover,  $\max_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda) \leq H$ .

Next, we will show that there exists  $b_1$  such that  $1 \in \arg \max_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$ . From now, we only consider sufficiently large  $b_1$  such that the value of  $\arg \max_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$  becomes finite.

Let  $\Psi_{b_1}^* := \max_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$  and  $f(b_1) := \frac{\Psi_{b_1}^* - \Psi_H^*}{b_1 - H}$ . Since  $f(H - 1/5) = 0$ ,  $f(0) = -\infty$ , and  $f(b_1)$  is continuous in  $b_1$ , the intermediate value theorem ensures that there exists  $b_1' \in [0, H]$  such that  $f(b_1') = -1$ . Moreover, the generalized mean value theorem (**Theorem 2.3.7** in Clarke (1990)) states that there exists  $b_1^* \in [b_1', H]$  such that  $-1 = f(b_1^*) \in \partial_{b_1} \Psi_{b_1^*}^*$ .

Let  $\Lambda_{b_1}$  be the set that provides maximums of  $\Psi_{b_1}$ , i.e.,  $\Lambda_{b_1} := \arg \max_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$ . Using Lemma 9,

$$-1 \in \partial_{b_1} \Psi_{b_1^*}^* = \text{conv}\{\nabla_{b_1} \Psi_{b_1^*}(\lambda) \mid \lambda \in \Lambda_{b_1^*}\} = \text{conv}\{-\lambda \mid \lambda \in \Lambda_{b_1^*}\} = [-\max \Lambda_{b_1^*}, -\min \Lambda_{b_1^*}].$$

Since  $\Psi_{b_1}(\lambda)$  is concave in  $\lambda$ , any  $\lambda \in [\min \Lambda_{b_1}, \max \Lambda_{b_1}]$  provides  $\max_{\lambda \in \mathbb{R}_+} \Psi_{b_1}(\lambda)$ . Thus,  $1 \in \arg \max_{\lambda \in \mathbb{R}_+} \Psi_{b_1^*}(\lambda)$ . This proves the existence of  $b_1$  such that  $1 \in \arg \max_{\lambda \in \mathbb{R}_+} \min_{\pi \in \Pi} L_\lambda(\pi)$ .  $\square$

## I MISSING PROOFS IN SECTION 5

### I.1 PROOF OF LEMMA 2

*Proof of Lemma 2.* We prove the first claim. Recall the definition of  $\Delta_{b_0}^*$ :

$$\Delta_{b_0}^* = \min_{\pi \in \Pi} \Delta_{b_0}(\pi) = \min_{\pi \in \Pi} \max_{n \in [0, N]} J_{c_n, \mathcal{U}}(\pi) - b_n. \quad (34)$$

It is easy to see that  $\Delta_{b_0}(\pi)$  is monotonically decreasing in  $b_0$ . Consider two real numbers  $x \leq y$  and let  $\pi^x \in \arg \min_{\pi \in \Pi} \Delta_x(\pi)$ . Then,

$$\Delta_y^* = \min_{\pi \in \Pi} \Delta_y(\pi) \leq \Delta_y(\pi^x) \leq \Delta_x(\pi^x) = \min_{\pi \in \Pi} \Delta_x(\pi) = \Delta_x^*.$$

1512 Therefore,  $\Delta_{b_0}^*$  is monotonically decreasing in  $b_0$ .

1513 Next, we prove the second claim. Suppose that  $\Delta_{J^*}^* < 0$ . Then, there exists a feasible policy  
 1514  $\pi \in \Pi_{\mathbb{F}}$  such that  $J_{c_0, \mathcal{U}}(\pi) < J^* = J_{c_0, \mathcal{U}}(\pi^*)$ . This contradicts the definition of the optimal policy.  
 1515 Therefore,  $\Delta_{J^*}^* \geq 0$ .

1517 Suppose that  $\Delta_{J^*}^* > 0$ . Since  $\min_{\pi \in \Pi} \Delta_{J^*}(\pi) > 0$ , no feasible policy achieves the objective return  
 1518  $J^*$ . This also contradicts the existence of the optimal policy. Therefore,  $\Delta_{J^*}^* = 0$ .  $\square$

## 1520 I.2 PROOF OF THEOREM 2

1522 *Proof of Theorem 2.* We first prove Equation (9) by contradiction. Let  $x :=$   
 1523  $\min\{b_0 \in [0, H] \mid \Delta_{b_0}^* \leq 0\}$  and suppose that  $x < J^*$ . Since  $\Delta_{J^*}^* = 0$  by Lemma 2, there  
 1524 exists a feasible policy  $\pi \in \Pi_{\mathbb{F}}$  such that  $J_{c_0, \mathcal{U}}(\pi) \leq x < J^* = J_{c_0, \mathcal{U}}(\pi^*)$ . This contradicts the  
 1525 definition of the optimal policy.

1526 We then show that Equation (9) provides  $\pi^*$ . Since  $\Delta_{J^*}^* = 0$  by Lemma 2, any policy  
 1527  $\pi \in \arg \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$  is feasible and satisfies  $J_{c_0, \mathcal{U}}(\pi) = J^*$ . The claim directly follows from  
 1528 the definition of an optimal policy.  $\square$

## 1530 I.3 PROOF OF LEMMA 3

1532 Instead of Lemma 3, we prove the following lemma that includes Lemma 3.

1533 **Lemma 18** (Properties of  $\Delta_{b_0}$ ). *The following properties hold for any  $b_0 \in \mathbb{R}$ .*

- 1535 1. (Lipschitz continuity): For any  $\pi, \pi' \in \Pi$ ,  $|\Delta_{b_0}(\pi) - \Delta_{b_0}(\pi')| \leq \ell_{\text{LP}} \|\pi - \pi'\|_2$  with  
 1536  $\ell_{\text{LP}} := H^2 \sqrt{A}$ .
- 1538 2. (Weak convexity):  $\Delta_{b_0}(\pi) + \frac{\ell_{\text{sm}}}{2} \|\pi\|_2^2$  is convex in  $\pi$  with  $\ell_{\text{sm}} := 2\gamma AH^3$ .
- 1540 3. (Subdifferentiability): For any  $\pi \in \Pi$ , the subgradient of  $\Delta_{b_0}$  at  $\pi$  is given by

$$\partial \Delta_{b_0}(\pi) = \text{conv}\{\nabla_{\pi} J_{c_n, P}(\pi) \mid n, P \in \mathcal{W}\},$$

1543 where  $\text{conv } B$  represents the convex hull of a set  $B \subset \mathbb{R}^{\mathcal{S} \times \mathcal{A}}$ .

1546 *Proof of Lipschitz continuity.*

$$\begin{aligned} 1547 |\Delta_{b_0}(\pi) - \Delta_{b_0}(\pi')| &\leq \left| \max_{n \in \llbracket 0, N \rrbracket} \{J_{c_n, \mathcal{U}}(\pi) - b_n\} - \max_{m \in \llbracket 0, N \rrbracket} \{J_{c_m, \mathcal{U}}(\pi') - b_m\} \right| \\ 1548 &\stackrel{(a)}{\leq} \max_{n \in \llbracket 0, N \rrbracket} |J_{c_n, \mathcal{U}}(\pi) - b_n - (J_{c_n, \mathcal{U}}(\pi') - b_n)| \stackrel{(b)}{\leq} \ell_{\text{LP}} \|\pi - \pi'\|_2 \end{aligned}$$

1552 where (a) uses Lemma 11 and (b) is due to Lemma 16. This concludes the proof of the Lipschitz  
 1553 continuity.  $\square$

1555 *Proof of weak convexity.* The weak convexity of  $\Delta_{b_0}(\pi) = \max_{n \in \llbracket 0, N \rrbracket} J_{c_n, \mathcal{U}}(\pi) - b_n$  immediately  
 1556 follows from the weak convexity of  $J_{c_n, \mathcal{U}}(\pi)$  due to Lemma 16 with Lemma 12.  $\square$

1559 *Proof of subdifferentiability.* Suppose that  $\mathcal{U}$  is a finite set. The claim directly follows from  
 1560 Lemma 10 with the weak convexity of  $J_{c_n, P}(\pi)$  due to Lemma 16.

1562 Suppose that  $\mathcal{U}$  is a compact set such that, for any  $\pi \in \Pi$ ,  $\nabla J_{c_n, P}(\pi)$  is continuous with respect to  
 1563  $P \in \mathcal{U}$ . Danskin's theorem (Lemma 9) indicates that, for any  $n \in \llbracket 0, N \rrbracket$ ,

$$1564 \partial J_{c_n, \mathcal{U}}(\pi) = \text{conv}\left\{ \nabla J_{c_n, P}(\pi) \mid P \in \arg \max_{P \in \mathcal{U}} J_{c_n, P}(\pi) - b_n \right\}.$$

Then, using Lemma 10 with the weak convexity of  $J_{c_n, \mathcal{U}}(\pi)$  due to Lemma 18, we have

$$\begin{aligned} \partial \Delta_{b_0}(\pi) &= \text{conv} \left\{ g \mid g \in \partial J_{c_n, \mathcal{U}}(\pi) \text{ where } n \in \arg \max_{n \in [0, N]} J_{c_n, \mathcal{U}}(\pi) - b_n \right\} \\ &= \text{conv} \left\{ \nabla J_{c_n, P}(\pi) \mid n, P \in \arg \max_{(n, P) \in [0, N] \times \mathcal{U}} J_{c_n, P}(\pi) - b_n \right\}. \end{aligned}$$

□

#### I.4 PROOF OF THEOREM 4

*Proof of Theorem 4.* We introduce shorthands  $\mathcal{G}$  and  $\mathcal{W}$  such that

$$\mathcal{G} = \{ \nabla J_{c_n, P}(\pi) \mid n, P \in \mathcal{W}_{b_0}(\pi) \} \text{ and } \mathcal{W} = \arg \max_{(n, P) \in [0, N] \times \mathcal{U}} J_{c_n, P}(\pi) - b_n. \quad (35)$$

Let  $\pi_{b_0}^* \in \arg \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$ . For any  $\pi \in \Pi$  and  $b_0 \in \mathbb{R}$ , we have

$$\begin{aligned} &\Delta_{b_0}(\pi) - \Delta_{b_0}(\pi_{b_0}^*) \\ &= \left( \max_{n \in [0, N]} \max_{P \in \mathcal{U}} J_{c_n, P}(\pi) - b_n \right) - \left( \max_{n \in [0, N]} \max_{P \in \mathcal{U}} J_{c_n, P}(\pi_{b_0}^*) - b_n \right) \\ &= \left( \min_{n, P \in \mathcal{W}} J_{c_n, P}(\pi) - b_n \right) - \left( \max_{n \in [0, N]} \max_{P \in \mathcal{U}} J_{c_n, P}(\pi_{b_0}^*) - b_n \right) \\ &\leq \min_{n, P \in \mathcal{W}} (J_{c_n, P}(\pi) - b_n) - (J_{c_n, P}(\pi_{b_0}^*) - b_n) \\ &= \min_{n, P \in \mathcal{W}} J_{c_n, P}(\pi) - J_{c_n, P}(\pi_{b_0}^*) \\ &\leq \min_{n, P \in \mathcal{W}} J_{c_n, P}(\pi) - \min_{\pi' \in \Pi} J_{c_n, P}(\pi') \\ &\stackrel{(a)}{\leq} H \min_{n, P \in \mathcal{W}} \left\| \frac{d_P^{\pi_{b_0}^*}}{\mu} \right\|_{\infty} \underbrace{\max_{\pi' \in \Pi} \langle \pi - \pi', \nabla_{\pi} J_{c_n, P}(\pi) \rangle}_{\geq 0 \text{ when } \pi' \text{ is greedy to } \nabla_{\pi} J_{c_n, P}(\pi)} \\ &\leq DH \min_{n, P \in \mathcal{W}} \max_{\pi' \in \Pi} \langle \pi - \pi', \nabla_{\pi} J_{c_n, P}(\pi) \rangle \\ &= DH \min_{g \in \mathcal{G}} \max_{\pi' \in \Pi} \langle \pi - \pi', g \rangle, \end{aligned} \quad (36)$$

where (a) uses Lemma 17.

The claim holds by showing that

$$\min_{g \in \mathcal{G}} \max_{\pi' \in \Pi} \langle \pi - \pi', g \rangle = \min_{g \in \partial \Delta_{b_0}(\pi)} \max_{\pi' \in \Pi} \langle \pi - \pi', g \rangle. \quad (37)$$

Since  $\text{conv}\{\mathcal{G}\} = \partial \Delta_{b_0}(\pi)$  due to Lemma 18, Equation (37) holds when there exists a  $g^* \in \arg \min_{g \in \text{conv}\{\mathcal{G}\}} \max_{\pi' \in \Pi} \langle \pi - \pi', g \rangle$  such that  $g^* \in \mathcal{G}$ .

Let  $z^* \in \arg \max_{\pi' \in \Pi} \min_{g \in \text{conv}\{\mathcal{G}\}} \langle \pi - \pi', g \rangle$ . For any  $g^* \in \arg \min_{g \in \text{conv}\{\mathcal{G}\}} \max_{\pi' \in \Pi} \langle \pi - \pi', g \rangle$ , it holds that

$$\begin{aligned} \max_{\pi' \in \Pi} \langle \pi - \pi', g^* \rangle &= \min_{g \in \text{conv}\{\mathcal{G}\}} \max_{\pi' \in \Pi} \langle \pi - \pi', g \rangle \\ &\stackrel{(a)}{=} \max_{\pi' \in \Pi} \min_{g \in \text{conv}\{\mathcal{G}\}} \langle \pi - \pi', g \rangle \\ &= \min_{g \in \text{conv}\{\mathcal{G}\}} \langle \pi - z^*, g \rangle \\ &\stackrel{(b)}{=} \min_{g \in \mathcal{G}} \langle \pi - z^*, g \rangle \end{aligned} \quad (38)$$

where (a) uses Sion's minimax theorem (Lemma 8) with the convexity of  $\Pi$  and  $\text{conv}\{\mathcal{G}\}$ , and (b) uses Lemma 15.

Note that

$$\langle \pi - z^*, g^* \rangle \leq \max_{\pi' \in \Pi} \langle \pi - \pi', g^* \rangle \stackrel{(a)}{=} \min_{g \in \text{conv}\{\mathcal{G}\}} \langle \pi - z^*, g \rangle \leq \langle \pi - z^*, g^* \rangle, \quad (39)$$

where (a) is due to the third line of Equation (38). The inequality must be equality. Accordingly,

$$\langle \pi - z^*, g^* \rangle \stackrel{(a)}{=} \max_{\pi' \in \Pi} \langle \pi - \pi', g^* \rangle \stackrel{(b)}{=} \min_{g \in \mathcal{G}} \langle \pi - z^*, g \rangle,$$

where (a) uses Equation (39) and (b) uses Equation (38). Therefore,  $g^* \in \mathcal{G}$  and thus Equation (37) holds. This concludes the proof.  $\square$

## J MISSING PROOFS IN SECTION 5

### J.1 PROOF OF THEOREM 3

To facilitate the analysis with estimation error, we present a slightly modified version of the epigraph form. Let  $\varepsilon \in \mathbb{R}$  be an admissible violation parameter. We introduce the following formulation:

$$(\mathbf{Epigraph}_\varepsilon) \quad J_\varepsilon^* := \min_{b_0 \in [0, H]} b_0 \quad \text{such that} \quad \Delta_{b_0}^* \leq \varepsilon. \quad (40)$$

Note that  $J_\varepsilon^*$  is monotonically decreasing in  $\varepsilon$ .

Additionally, we introduce a slightly generalized version of Theorem 2:

**Lemma 19.** *For any  $\varepsilon_1, \varepsilon_2 \geq 0$ , if  $b_0$  and a policy  $\pi \in \Pi$  satisfy  $b_0 \leq J^* + \varepsilon_2$  and  $\Delta_{b_0}(\pi) \leq \varepsilon_1$ , then  $\pi$  is an  $(\varepsilon_1 + \varepsilon_2)$ -optimal policy.*

*Proof.* Note that  $J_{c_0, \mathcal{U}}(\pi) \leq J^* + \varepsilon_1 + \varepsilon_2$  and  $J_{n, \mathcal{U}}(\pi) \leq b_n + \varepsilon_1$  for any  $n \in \llbracket 1, N \rrbracket$ . The claim directly follows from Definition 1 and the fact that  $J^* = J_{c_0, \mathcal{U}}(\pi^*)$ .  $\square$

For any  $\bar{b}_0 \in [J_{\varepsilon_1}^*, J^* + \varepsilon_2]$  with some  $\varepsilon_1, \varepsilon_2 \geq 0$ , the subroutine returns a policy  $\bar{\pi} = \mathcal{A}(\bar{b}_0)$  such that

$$\Delta_{\bar{b}_0}(\bar{\pi}) \stackrel{(a)}{\leq} \min_{\pi' \in \Pi} \Delta_{\bar{b}_0}(\pi') + \varepsilon_{\text{opt}} \stackrel{(b)}{\leq} \min_{\pi' \in \Pi} \Delta_{J_{\varepsilon_1}^*}(\pi') + \varepsilon_{\text{opt}} \stackrel{(c)}{\leq} \varepsilon_1 + \varepsilon_{\text{opt}},$$

where (a) is due to Assumption 4, (b) holds since  $\Delta_{b_0}(\pi)$  is monotonically decreasing in  $b_0$ , and (c) follows from Equation (40). Consequently, by applying Lemma 19,  $\bar{\pi}$  is  $(\varepsilon_1 + \varepsilon_2 + \varepsilon_{\text{opt}})$ -optimal.

The following intermediate lemma guarantees that the search space of Algorithm 1 always contains such  $\bar{b}_0$  with  $\varepsilon_1 = \varepsilon_{\text{est}}$  and  $\varepsilon_2 = \varepsilon_{\text{est}} + \varepsilon_{\text{opt}}$ .

**Lemma 20.** *Suppose that Algorithm 1 is run with algorithms  $\widehat{J}_n$  and  $\mathcal{A}$  that satisfy Assumptions 1 and 4. For any  $k \in \llbracket 0, K \rrbracket$ ,  $[i^{(k)}, j^{(k)}] \cap [J_{\varepsilon_{\text{est}}}^*, J^* + \varepsilon_{\text{est}} + \varepsilon_{\text{opt}}] \neq \emptyset$ .*

*Proof.* The claim holds for  $k = 0$ . Suppose that the claim holds for a fixed  $k$ . Recall  $\widehat{\Delta}^{(k)}$  defined in Equation (10). If  $\widehat{\Delta}^{(k)} > 0$ , it holds that

$$-\varepsilon_{\text{est}} - \varepsilon_{\text{opt}} \stackrel{(a)}{<} \widehat{\Delta}^{(k)} - \left| \widehat{\Delta}^{(k)} - \Delta_{b_0^{(k)}}(\pi^{(k)}) \right| - \varepsilon_{\text{opt}} \leq \Delta_{b_0^{(k)}}(\pi^{(k)}) - \varepsilon_{\text{opt}} \stackrel{(b)}{\leq} \Delta_{b_0^{(k)}}(\pi^*) \stackrel{(c)}{=} J^* - b_0^{(k)} \quad (41)$$

where (a) is due to Assumption 1 with  $\widehat{\Delta}^{(k)} > 0$ , (b) is due to Assumption 4, and (c) holds since  $\pi^*$  is a feasible policy. Combining this with the induction assumption and the update rule of Equation (11), we have  $i^{(k+1)} = b_0^{(k)} \leq J^* + \varepsilon_{\text{est}} + \varepsilon_{\text{opt}}$  and  $J_{\varepsilon_{\text{est}}}^* \leq j^{(k+1)}$ . Hence,  $[i^{(k+1)}, j^{(k+1)}] \cap [J_{\varepsilon_{\text{est}}}^*, J^* + \varepsilon_{\text{est}} + \varepsilon_{\text{opt}}] \neq \emptyset$  when  $\widehat{\Delta}^{(k)} > 0$ .

On the other hand, if  $\widehat{\Delta}^{(k)} \leq 0$ , we have

$$\min_{\pi} \Delta_{b_0^{(k)}}(\pi) \leq \Delta_{b_0^{(k)}}(\pi^{(k)}) \leq \widehat{\Delta}^{(k)} + \varepsilon_{\text{est}} \leq \varepsilon_{\text{est}}. \quad (42)$$

Since  $b_0^{(k)}$  is the feasible solution to Equation (40), it holds that  $J_{\varepsilon_{\text{est}}}^* \leq b_0^{(k)} = j^{(k+1)}$ . Accordingly, we have  $[i^{(k+1)}, j^{(k+1)}] \cap [J_{\varepsilon_{\text{est}}}^*, J^* + \varepsilon_{\text{est}} + \varepsilon_{\text{opt}}] \neq \emptyset$ . Therefore, the claim holds for any  $k \in \llbracket 0, K \rrbracket$ .  $\square$

We are now ready to prove Theorem 3.

*Proof of Theorem 3.* Note that  $j^{(k)} - i^{(k)} \leq (j^{(0)} - i^{(0)})2^{-k} = H2^{-k}$  due to the update rule of Equation (11). According to Lemma 20, we have  $J_{\varepsilon_{\text{est}}}^* \leq j^{(K)} \leq J^* + \varepsilon_{\text{est}} + \varepsilon_{\text{opt}} + H2^{-K}$ . Additionally, the returned policy  $\pi_{\text{ret}}$  satisfies

$$\Delta_{j^{(K)}}(\pi_{\text{ret}}) \stackrel{(a)}{\leq} \min_{\pi \in \Pi} \Delta_{j^{(K)}}(\pi) + \varepsilon_{\text{opt}} \stackrel{(b)}{\leq} \min_{\pi \in \Pi} \Delta_{J_{\varepsilon_{\text{est}}}^*}(\pi) + \varepsilon_{\text{opt}} \leq \varepsilon_{\text{est}} + \varepsilon_{\text{opt}},$$

where (a) uses Assumption 4 and (b) is due to  $J_{\varepsilon_{\text{est}}}^* \leq j^{(K)}$  and the fact that  $\min_{\pi} \Delta_{b_0}(\pi)$  is monotonically decreasing in  $b_0$ . Applying this to Lemma 19 with  $j^{(K)} \leq J^* + \varepsilon_{\text{est}} + \varepsilon_{\text{opt}} + H2^{-K}$  concludes the proof.  $\square$

## J.2 PROOF OF THEOREM 5

We prove the following restatement of Theorem 5 with concrete values.

**Theorem 6** (Restatement of Theorem 5). *Suppose Assumptions 2 and 5 hold. Suppose that Algorithm 2 is run with algorithms  $\hat{J}_n$  and  $\hat{J}^\partial$  that satisfy Assumption 1 and Assumption 3. Let*

$$C := \frac{\ell_{\text{LP}}}{2\ell_{\text{sm}}} + 2DH\sqrt{S} = \frac{1}{2\gamma H\sqrt{A}} + 2DH\sqrt{S} \underset{\text{when } \gamma \approx 1}{=} \tilde{\mathcal{O}}(DH\sqrt{S}),$$

where  $\ell_{\text{LP}}$  and  $\ell_{\text{sm}}$  are defined in Lemma 18 and  $D$  is defined in Theorem 4. Assume that the evaluators  $\hat{J}$  and  $\hat{J}^\partial$  are sufficiently accurate such that

$$\varepsilon_{\text{grd}} = C_\partial \varepsilon^2 \text{ where } C_\partial := \frac{1}{1024C^2\ell_{\text{sm}}\sqrt{S}} \text{ and } \varepsilon_{\text{est}} = C_J \varepsilon^2 \text{ where } C_J := \frac{1}{1024C^2\ell_{\text{sm}}}.$$

Set  $\alpha = C_\alpha \varepsilon^2$  and  $T = C_T \varepsilon^{-4}$  such that

$$C_\alpha := \frac{1}{64C^2\ell_{\text{sm}}(\ell_{\text{LP}}^2 + \varepsilon_{\text{grd}})} \text{ and } C_T := 4096C^4\ell_{\text{sm}}^2 S(\ell_{\text{LP}}^2 + \varepsilon_{\text{grd}}) = \tilde{\mathcal{O}}(D^4S^3A^3H^{14}).$$

Then, Algorithm 1 returns a policy  $\pi^{(t^*)}$  such that

$$\Delta_{b_0}(\pi^{(t^*)}) - \min_{\pi \in \Pi} \Delta_{b_0}(\pi) \leq \varepsilon.$$

We first introduce the following useful lemma.

**Lemma 21.** *Let  $(M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0}) : \pi \mapsto \min_{\pi' \in \Pi} \{ \Delta_{b_0}(\pi') + \ell_{\text{sm}} \|\pi - \pi'\|_2^2 \}$  be the Moreau envelope function of  $\Delta_{b_0}(\pi)$  with parameter  $1/2\ell_{\text{sm}}$ . For any policy  $\pi \in \Pi$ ,*

$$\Delta_{b_0}(\pi) - \min_{\pi' \in \Pi} \Delta_{b_0}(\pi') \leq C \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}} \circ \Delta_{b_0}} \right) (\pi) \right\|_2.$$

*Proof.* Define  $\bar{\pi} := \arg \min_{\pi' \in \Pi} \Delta_{b_0}(\pi') + \ell_{\text{sm}} \|\pi - \pi'\|_2^2$ . According to Lemma 14 with  $\tau = 1/2\ell_{\text{sm}}$ , there exists a subgradient  $g \in \partial \Delta_{b_0}(\bar{\pi})$  such that, for any  $\pi' \in \Pi$ ,

$$\begin{aligned} \langle \bar{\pi} - \pi', g \rangle &\leq \left\langle \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}} \circ \Delta_{b_0}} \right) (\pi), \bar{\pi} - \pi' \right\rangle \\ &\stackrel{(a)}{\leq} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}} \circ \Delta_{b_0}} \right) (\pi) \right\|_2 \|\bar{\pi} - \pi'\|_2 \\ &\stackrel{(b)}{\leq} 2\sqrt{S} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}} \circ \Delta_{b_0}} \right) (\pi) \right\|_2, \end{aligned} \tag{43}$$

where (a) is due to the Cauchy–Schwarz inequality and (b) uses that, for any  $\pi' \in \Pi$

$$\begin{aligned} \|\bar{\pi} - \pi'\|_2 &= \sqrt{\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} (\bar{\pi}(s, a) - \pi'(s, a))^2} \leq \sqrt{S} \max_{s \in \mathcal{S}} \sqrt{\sum_{a \in \mathcal{A}} (\bar{\pi}(s, a) - \pi'(s, a))^2} \\ &\leq \sqrt{S} \max_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} |\bar{\pi}(s, a) - \pi'(s, a)| \leq 2\sqrt{S}. \end{aligned} \tag{44}$$



Let  $\pi_{b_0}^* \in \arg \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$ . Inserting this result into Theorem 4, we have

$$\begin{aligned} \Delta_{b_0}(\bar{\pi}) - \Delta_{b_0}(\pi_{b_0}^*) &\leq DH \max_{\pi' \in \Pi} \langle \bar{\pi} - \pi', g \rangle \quad \forall g \in \partial \Delta_{b_0}(\bar{\pi}) \\ &\leq 2DH\sqrt{S} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi) \right\|_2. \end{aligned}$$

Therefore,

$$\begin{aligned} \Delta_{b_0}(\pi) - \Delta_{b_0}(\pi_{b_0}^*) &= \Delta_{b_0}(\pi) - \Delta_{b_0}(\bar{\pi}) + \Delta_{b_0}(\bar{\pi}) - \Delta_{b_0}(\pi_{b_0}^*) \\ &\stackrel{(a)}{\leq} \Delta_{b_0}(\pi) - \Delta_{b_0}(\bar{\pi}) + 2DH\sqrt{S} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi) \right\|_2 \\ &\stackrel{(b)}{\leq} \ell_{\text{LP}} \|\pi - \bar{\pi}\|_2 + 2DH\sqrt{S} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi) \right\|_2 \\ &\stackrel{(c)}{\leq} \frac{\ell_{\text{LP}}}{2\ell_{\text{sm}}} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi) \right\|_2 + 2DH\sqrt{S} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi) \right\|_2, \end{aligned}$$

where (a) is due to Theorem 4, (b) is due to the Lipschitz continuity by Lemma 18, and (c) uses Lemma 7. This concludes the proof.  $\square$

**Lemma 22.** *Under the settings of Theorem 6,*

$$\sum_{t=0}^{T-1} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) \right\|_2^2 \leq \frac{16\ell_{\text{sm}}S}{\alpha} + 4T \left( \alpha\ell_{\text{sm}}(\ell_{\text{LP}}^2 + \varepsilon_{\text{grd}}) + 4\ell_{\text{sm}}\varepsilon_{\text{grd}}\sqrt{S} + 2\ell_{\text{sm}}\varepsilon_{\text{est}} \right).$$

*Proof.* Recall that

$$\pi^{(t+1)} \in \arg \min_{\pi \in \Pi} \left\langle g^{(t)}, \pi - \pi^{(t)} \right\rangle + \frac{1}{2\alpha} \left\| \pi - \pi^{(t)} \right\|_2^2 = \text{Proj}_{\Pi} \left( \pi^{(t)} - \alpha g^{(t)} \right).$$

Define  $\bar{\pi}^{(t)} := \arg \min_{\pi'} \Delta_{b_0}(\pi') + \ell_{\text{sm}} \|\pi^{(t)} - \pi'\|_2^2$ . Then, we have

$$\begin{aligned} \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t+1)}) &= \min_{\pi \in \Pi} \Delta_{b_0}(\pi) + \ell_{\text{sm}} \left\| \pi^{(t+1)} - \pi \right\|_2^2 \\ &\leq \Delta_{b_0}(\bar{\pi}^{(t)}) + \ell_{\text{sm}} \left\| \pi^{(t+1)} - \bar{\pi}^{(t)} \right\|_2^2 \\ &= \Delta_{b_0}(\bar{\pi}^{(t)}) + \ell_{\text{sm}} \left\| \text{Proj}_{\Pi} \left( \pi^{(t)} - \alpha g^{(t)} \right) - \text{Proj}_{\Pi} \left( \bar{\pi}^{(t)} \right) \right\|_2^2 \\ &\leq \Delta_{b_0}(\bar{\pi}^{(t)}) + \ell_{\text{sm}} \left\| \pi^{(t)} - \alpha g^{(t)} - \bar{\pi}^{(t)} \right\|_2^2 \\ &= \underbrace{\Delta_{b_0}(\bar{\pi}^{(t)}) + \ell_{\text{sm}} \left\| \pi^{(t)} - \bar{\pi}^{(t)} \right\|_2^2}_{= \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)})} + \underbrace{2\ell_{\text{sm}}\alpha \left\langle g^{(t)}, \bar{\pi}^{(t)} - \pi^{(t)} \right\rangle}_{=: \textcircled{1}} + \underbrace{\alpha^2 \ell_{\text{sm}} \left\| g^{(t)} \right\|_2^2}_{=: \textcircled{2}}. \end{aligned} \tag{45}$$

We further upper bound  $\textcircled{1}$  and  $\textcircled{2}$ . Recall  $n^{(t)} \in \arg \max_{n \in \llbracket 0, N \rrbracket} \widehat{J}_n(\pi^{(t)}) - b_n$  and  $g^{(t)} = \widehat{J}_{n^{(t)}}^{\partial}(\pi^{(t)})$ . Due to Assumption 3, there exists an vector  $g' \in \mathbb{R}^{S \times A}$  that satisfies

$$g' \in \left\{ \nabla J_{c_{n^{(t)}}, P^{(t)}}(\pi^{(t)}) \mid P^{(t)} \in \arg \max_{P \in \mathcal{U}} J_{c_{n^{(t)}}, P}(\pi^{(t)}) \right\} \text{ and } \left\| g^{(t)} - g' \right\|_2^2 \leq \varepsilon_{\text{grd}}^2. \tag{46}$$

Thus, using Lemma 16, we have  $\textcircled{2} \leq \|g'\|_2^2 + \|g^{(t)} - g'\|_2^2 \leq \ell_{\text{LP}}^2 + \varepsilon_{\text{grd}}^2$ . Furthermore, we have

$$\begin{aligned} \textcircled{1} &= \left\langle g^{(t)}, \bar{\pi}^{(t)} - \pi^{(t)} \right\rangle = \left\langle g', \bar{\pi}^{(t)} - \pi^{(t)} \right\rangle + \left\langle g^{(t)} - g', \bar{\pi}^{(t)} - \pi^{(t)} \right\rangle \\ &\stackrel{(a)}{\leq} \left\langle g', \bar{\pi}^{(t)} - \pi^{(t)} \right\rangle + \left\| g^{(t)} - g' \right\|_2 \left\| \bar{\pi}^{(t)} - \pi^{(t)} \right\|_2 \stackrel{(b)}{\leq} \left\langle g', \bar{\pi}^{(t)} - \pi^{(t)} \right\rangle + 2\varepsilon_{\text{grd}}\sqrt{S} \end{aligned}$$

where (a) is due to the Cauchy–Schwarz inequality and (b) uses Equation (44).

By inserting the above inequalities to Equation (45), using  $g'$  defined in Equation (46), we have

$$\begin{aligned} \textcircled{3} &:= 2\ell_{\text{sm}}\alpha \left\langle g', \pi^{(t)} - \bar{\pi}^{(t)} \right\rangle \\ &\leq \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) - \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t+1)}) + \alpha^2 \ell_{\text{sm}} (\ell_{\text{LP}}^2 + \varepsilon_{\text{grd}}^2) + 4\ell_{\text{sm}}\alpha \varepsilon_{\text{grd}} \sqrt{S}. \end{aligned} \quad (47)$$

Next, we are going to derive the lower bound of  $\textcircled{3}$ . Define  $\Delta_{b_0}^{(t)}(\pi)$  such that

$$\Delta_{b_0}^{(t)}(\pi) := J_{c_{n^{(t)}}, \mathcal{U}}(\pi) - b_{n^{(t)}}. \quad (48)$$

Additionally, let  $\delta_n$  and  $\widehat{\delta}_n$  be shorthand such that  $\delta_n := J_{c_n, \mathcal{U}}(\pi) - b_n$  and  $\widehat{\delta}_n := \widehat{J}_n(\pi) - b_n$ . Then, Due to Assumption 1, for any  $\pi$ , we have

$$\begin{aligned} \left| \Delta_{b_0}^{(t)}(\pi) - \Delta_{b_0}(\pi) \right| &\stackrel{(a)}{\leq} \underbrace{\left| \Delta_{b_0}^{(t)}(\pi) - \widehat{\delta}_{n^{(t)}} \right|}_{\leq \varepsilon_{\text{est}} \text{ by Assumption 1}} + \left| \max_{n \in \llbracket 0, N \rrbracket} \widehat{\delta}_n - \Delta_{b_0}(\pi) \right| \\ &\stackrel{(b)}{\leq} \varepsilon_{\text{est}} + \max_{n \in \llbracket 0, N \rrbracket} \left| \widehat{\delta}_n - \delta_n \right| \leq 2\varepsilon_{\text{est}}. \end{aligned} \quad (49)$$

where (a) is due to the definition of  $n^{(t)}$  and (b) uses Lemma 11.

Due to the weak convexity of  $\Delta_{b_0}^{(t)}(\pi)$  with respect to  $\pi$  (Lemma 18) and since  $g' \in \partial \Delta_{b_0}^{(t)}(\pi^{(t)})$ ,

$$\begin{aligned} \textcircled{3} / 2\ell_{\text{sm}}\alpha &= \left\langle g', \pi^{(t)} - \bar{\pi}^{(t)} \right\rangle \\ &\geq \Delta_{b_0}^{(t)}(\pi^{(t)}) - \Delta_{b_0}^{(t)}(\bar{\pi}^{(t)}) - \frac{\ell_{\text{sm}}}{2} \left\| \bar{\pi}^{(t)} - \pi^{(t)} \right\|_2^2 \\ &\geq - \underbrace{\left| \Delta_{b_0}^{(t)}(\pi^{(t)}) - \Delta_{b_0}(\pi^{(t)}) \right|}_{\leq \varepsilon_{\text{est}} \text{ by Equation (49)}} - \underbrace{\left| \Delta_{b_0}(\bar{\pi}^{(t)}) - \Delta_{b_0}^{(t)}(\bar{\pi}^{(t)}) \right|}_{\leq \varepsilon_{\text{est}} \text{ by Equation (49)}} + \Delta_{b_0}(\pi^{(t)}) - \Delta_{b_0}(\bar{\pi}^{(t)}) - \frac{\ell_{\text{sm}}}{2} \left\| \bar{\pi}^{(t)} - \pi^{(t)} \right\|_2^2 \\ &= \Delta_{b_0}(\pi^{(t)}) + \ell_{\text{sm}} \left\| \pi^{(t)} - \pi^{(t)} \right\|_2^2 - \Delta_{b_0}(\bar{\pi}^{(t)}) - \ell_{\text{sm}} \left\| \bar{\pi}^{(t)} - \pi^{(t)} \right\|_2^2 + \frac{\ell_{\text{sm}}}{2} \left\| \bar{\pi}^{(t)} - \pi^{(t)} \right\|_2^2 - 2\varepsilon_{\text{est}} \\ &= \Delta_{b_0}(\pi^{(t)}) + \ell_{\text{sm}} \left\| \pi^{(t)} - \pi^{(t)} \right\|_2^2 - \min_{\pi' \in \Pi} \left( \Delta_{b_0}(\pi') + \ell_{\text{sm}} \left\| \pi' - \pi^{(t)} \right\|_2^2 \right) + \frac{\ell_{\text{sm}}}{2} \left\| \bar{\pi}^{(t)} - \pi^{(t)} \right\|_2^2 - 2\varepsilon_{\text{est}} \\ &\geq \frac{\ell_{\text{sm}}}{2} \left\| \bar{\pi}^{(t)} - \pi^{(t)} \right\|_2^2 - 2\varepsilon_{\text{est}} \\ &\stackrel{(a)}{=} \frac{\ell_{\text{sm}}}{2} \left\| \frac{1}{2\ell_{\text{sm}}} \nabla \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) \right\|_2^2 = \frac{1}{8\ell_{\text{sm}}} \left\| \nabla \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) \right\|_2^2 - 2\varepsilon_{\text{est}}, \end{aligned}$$

where (a) uses Lemma 7. By inserting this to Equation (47),

$$\begin{aligned} &\frac{\alpha}{4} \left\| \nabla \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) \right\|_2^2 - 8\alpha \ell_{\text{sm}} \varepsilon_{\text{est}} \\ &\leq \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) - \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t+1)}) + \alpha^2 \ell_{\text{sm}} (\ell_{\text{LP}}^2 + \varepsilon_{\text{grd}}^2) + 4\alpha \ell_{\text{sm}} \varepsilon_{\text{grd}} \sqrt{S}. \end{aligned}$$

By taking summation over  $\sum_{t=0}^{T-1}$ ,

$$\begin{aligned} \frac{\alpha}{4} \sum_{t=0}^{T-1} \left\| \nabla \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) \right\|_2^2 &\leq \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(0)}) - \left( \mathbb{M}_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(T)}) \\ &\quad + T \left( \alpha^2 \ell_{\text{sm}} (\ell_{\text{LP}}^2 + \varepsilon_{\text{grd}}^2) + 4\alpha \ell_{\text{sm}} \varepsilon_{\text{grd}} \sqrt{S} + 8\alpha \ell_{\text{sm}} \varepsilon_{\text{est}} \right). \end{aligned}$$

Note that

$$\begin{aligned}
& \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(0)}) - \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(T)}) \\
&= \min_{\pi \in \Pi} \left\{ \Delta_{b_0}(\pi) + \ell_{\text{sm}} \left\| \pi^{(0)} - \pi \right\|_2^2 \right\} - \min_{\pi \in \Pi} \left\{ \Delta_{b_0}(\pi) + \ell_{\text{sm}} \left\| \pi^{(T)} - \pi \right\|_2^2 \right\} \\
&= \Delta_{b_0}(\bar{\pi}^{(0)}) + \ell_{\text{sm}} \left\| \pi^{(0)} - \bar{\pi}^{(0)} \right\|_2^2 - \Delta_{b_0}(\bar{\pi}^{(T)}) - \ell_{\text{sm}} \left\| \pi^{(T)} - \bar{\pi}^{(T)} \right\|_2^2 \\
&\leq \Delta_{b_0}(\bar{\pi}^{(T)}) + \ell_{\text{sm}} \left\| \pi^{(0)} - \bar{\pi}^{(T)} \right\|_2^2 - \Delta_{b_0}(\bar{\pi}^{(T)}) - \ell_{\text{sm}} \left\| \pi^{(T)} - \bar{\pi}^{(T)} \right\|_2^2 \\
&\leq \ell_{\text{sm}} \left\| \pi^{(0)} - \bar{\pi}^{(T)} \right\|_2^2 \leq 4\ell_{\text{sm}}S,
\end{aligned}$$

where the last inequality uses Equation (44).

By combining all the results, we obtain

$$\sum_{t=0}^{T-1} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi^{(t)}) \right\|_2^2 \leq \frac{16\ell_{\text{sm}}S}{\alpha} + 4T \left( \alpha \ell_{\text{sm}} (\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2) + 4\ell_{\text{sm}}\varepsilon_{\text{grd}}\sqrt{S} + 2\ell_{\text{sm}}\varepsilon_{\text{est}} \right).$$

This concludes the proof.  $\square$

We are now ready to prove Theorem 6.

*Proof of Theorem 6.* Let  $\pi_{b_0}^* \in \arg \min_{\pi \in \Pi} \Delta_{b_0}(\pi)$ . Then,

$$\begin{aligned}
& \min_{t \in [0, T-1]} \Delta_{b_0}(\pi^{(t)}) - \Delta_{b_0}(\pi_{b_0}^*) \\
&\leq \frac{1}{T} \sum_{t=0}^{T-1} \Delta_{b_0}(\pi^{(t)}) - \Delta_{b_0}(\pi_{b_0}^*) \\
&\stackrel{(a)}{\leq} \frac{1}{T} C \sum_{t=0}^{T-1} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi) \right\|_2 \\
&\leq C \sqrt{\frac{1}{T} \sum_{t=0}^{T-1} \left\| \nabla \left( M_{\frac{1}{2\ell_{\text{sm}}}} \circ \Delta_{b_0} \right) (\pi) \right\|_2^2} \\
&\stackrel{(b)}{\leq} C \sqrt{\frac{16\ell_{\text{sm}}S}{T\alpha} + 4 \left( \alpha \ell_{\text{sm}} (\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2) + 4\ell_{\text{sm}}\varepsilon_{\text{grd}}\sqrt{S} + 2\ell_{\text{sm}}\varepsilon_{\text{est}} \right)} } \\
&\stackrel{(c)}{=} C \sqrt{\frac{16\ell_{\text{sm}}S}{\delta\sqrt{T}} + \frac{4\delta}{\sqrt{T}} \ell_{\text{sm}} (\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2) + 16\ell_{\text{sm}}\varepsilon_{\text{grd}}\sqrt{S} + 8\ell_{\text{sm}}\varepsilon_{\text{est}}} } \\
&\stackrel{(d)}{\leq} 4C\sqrt{\ell_{\text{sm}}S\delta^{-1}T^{-\frac{1}{4}}} + 2C\sqrt{\ell_{\text{sm}}(\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2)\delta T^{-\frac{1}{4}}} + 4C\sqrt{\ell_{\text{sm}} \left( S^{\frac{1}{4}}\sqrt{\varepsilon_{\text{grd}}} + \sqrt{\varepsilon_{\text{est}}} \right)} \\
&\stackrel{(e)}{=} 4C\sqrt{\ell_{\text{sm}}S^{\frac{1}{4}}(\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2)^{\frac{1}{4}}T^{-\frac{1}{4}}} + 4C\sqrt{\ell_{\text{sm}} \left( S^{\frac{1}{4}}\sqrt{\varepsilon_{\text{grd}}} + \sqrt{\varepsilon_{\text{est}}} \right)}.
\end{aligned}$$

where (a) uses Lemma 21, (b) uses Lemma 22, (c) replaces  $\alpha$  with  $\delta/\sqrt{T}$ , (d) uses  $\sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$ , and (e) sets  $\delta = \sqrt{S/(\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2)}$ .

Therefore, when  $\varepsilon_{\text{est}}$ ,  $\varepsilon_{\text{grd}}$ ,  $\alpha$ , and  $T$  satisfy:

$$\begin{aligned}
\varepsilon_{\text{grd}} &= \frac{\varepsilon^2}{1024C^2\ell_{\text{sm}}\sqrt{S}}, \quad \varepsilon_{\text{est}} = \frac{\varepsilon^2}{1024C^2\ell_{\text{sm}}}, \\
T &= 4096C^4\ell_{\text{sm}}^2S(\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2)\varepsilon^{-4}, \quad \text{and} \quad \alpha = \frac{\delta}{\sqrt{T}} = \frac{\varepsilon^2}{64C^2\ell_{\text{sm}}(\ell_{\text{Lp}}^2 + \varepsilon_{\text{grd}}^2)},
\end{aligned}$$

1890 we have

$$1891 \min_{t \in \llbracket 0, T-1 \rrbracket} \Delta_{b_0}(\pi^{(t)}) \leq \Delta_{b_0}(\pi_{b_0}^*) + \frac{3}{4}\varepsilon .$$

1892  
1893 Finally,  $t^* \in \arg \min_{t \in \llbracket 0, T-1 \rrbracket} \widehat{\Delta}^{(t)}$  satisfies that

$$1894 \begin{aligned} \Delta_{b_0}(\pi^{(t^*)}) &= \widehat{\Delta}^{(t^*)} + \Delta_{b_0}(\pi^{(t^*)}) - \widehat{\Delta}^{(t^*)} \\ 1895 &\leq \min_{t \in \llbracket 0, T-1 \rrbracket} \widehat{\Delta}^{(t)} + \varepsilon_{\text{est}} \\ 1896 &\leq \min_{t \in \llbracket 0, T-1 \rrbracket} \Delta_{b_0}(\pi^{(t)}) + \widehat{\Delta}^{(t)} - \Delta_{b_0}(\pi^{(t)}) + \varepsilon_{\text{est}} \\ 1897 &\leq \min_{t \in \llbracket 0, T-1 \rrbracket} \Delta_{b_0}(\pi^{(t)}) + 2\varepsilon_{\text{est}} \\ 1898 &\leq \Delta_{b_0}(\pi_{b_0}^*) + \frac{3}{4}\varepsilon + 2\varepsilon_{\text{est}} \\ 1899 &\leq \Delta_{b_0}(\pi_{b_0}^*) + \varepsilon . \end{aligned}$$

1900  
1901  
1902  
1903  
1904  
1905  
1906 This concludes the proof. □

1907  
1908  
1909  
1910  
1911  
1912  
1913  
1914  
1915  
1916  
1917  
1918  
1919  
1920  
1921  
1922  
1923  
1924  
1925  
1926  
1927  
1928  
1929  
1930  
1931  
1932  
1933  
1934  
1935  
1936  
1937  
1938  
1939  
1940  
1941  
1942  
1943