
A Dual Approach to Constrained Markov Decision Processes with Entropy Regularization

Donghao Ying

University of California, Berkeley
donghaoy@berkeley.edu

Yuhao Ding

University of California, Berkeley
yuhao_ding@berkeley.edu

Javad Lavaei

University of California, Berkeley
lavaei@berkeley.edu

Abstract

We study entropy-regularized constrained Markov decision processes (CMDPs) under the soft-max parameterization, in which an agent aims to maximize the entropy-regularized value function while satisfying constraints on the expected total utility. By leveraging the entropy regularization, our theoretical analysis shows that its Lagrangian dual function is smooth and the Lagrangian duality gap can be decomposed into the primal optimality gap and the constraint violation. Furthermore, we propose an accelerated dual-descent method for entropy-regularized CMDPs. We prove that our method achieves the global convergence rate $\tilde{O}(1/T)$ for both the optimality gap and the constraint violation for entropy-regularized CMDPs. A discussion about a linear convergence rate for CMDPs with a single constraint is also provided.

1 INTRODUCTION

In many sequential decision-making problems for safety-critical systems, e.g. autonomous driving (Fisac et al., 2018) and cyber-physical systems (Zhang et al., 2019), the optimality of an objective function by itself is not sufficient and a variety of constraints must be satisfied. This has naturally led to a generalization of the model of Markov Decision Processes (MDPs) to Constrained MDPs (CMDPs) (Altman, 1999), in which an agent aims to maximize the value function while satisfying given constraints on the expected total utility.

Direct policy search methods, including the policy gra-

dient and the natural policy gradient (NPG) methods, have had substantial empirical successes in solving CMDPs (Achiam et al., 2017; Chow et al., 2017; Bhatnagar and Lakshmanan, 2012; Borkar, 2005; Uchibe and Doya, 2007; Achiam et al., 2017). Recently, a major progress in understanding the theoretical non-asymptotic global convergence behavior of policy-based methods for CMDPs has also been achieved (Ding et al., 2020, 2021; Xu et al., 2021; Efroni et al., 2020; Chen et al., 2021).

For policy-based methods, entropy regularization is a popular technique for encouraging the exploration of an unknown environment and preventing a premature convergence (Williams and Peng, 1991; Mnih et al., 2016; Haarnoja et al., 2018; Zang et al., 2020). From a theoretical optimization perspective, it is shown in Mei et al. (2020) and Cen et al. (2021) that the entropy regularization can make the policy optimization landscape benign and achieve faster convergence rates even in the exact value evaluation setting. Nevertheless, most existing theoretical guarantees for the entropy-regularized policy optimization are restricted to unconstrained MDPs. The scope of the power of entropy regularization for CMDPs remains unknown even for the tabular setting with the exact value evaluation.

Inspired by the recent theoretical advances towards understanding entropy-regularized policy gradient methods (Mei et al., 2020; Cen et al., 2021) together with the global convergence of Lagrangian-based methods for CMDPs (Ding et al., 2020, 2021; Paternain et al., 2019; Xu et al., 2021), we investigate the optimization properties induced by the entropy regularization for CMDPs under the soft-max policy parameterization. We focus on the study of tabular CMDPs with the exact gradient evaluation. This is the setting commonly investigated in the literature since its understanding assists in demystifying the effectiveness of entropy-regularization in CMDPs with more complex settings.

1.1 Contributions

This work is the first one that certifies the effectiveness of entropy regularization in CMDPs from an optimization perspective. We summarize our contributions below:

- We first show that although the underlying problem is nonconcave, the Lagrangian dual function of CMDPs with the entropy regularization is smooth under the Slater condition and the exploratory initial distribution assumption. Under the same conditions, an $\mathcal{O}(\varepsilon)$ error bound for the dual optimality gap leads to an $\mathcal{O}(\sqrt{\varepsilon})$ error bound for the primal optimality gap and the constraint violation.
- To leverage the smoothness of the Lagrangian dual function, we propose a new accelerated dual-descent method for entropy-regularized CMDPs, which updates the dual variable via projected accelerated gradient descent and uses the natural policy gradient method in the inner loop.
- We prove that the proposed method achieves a global convergence with the rate $\tilde{\mathcal{O}}(1/T)$ for both the optimality gap and the constraint violation for entropy-regularized CMDPs.
- In the special case where CMDPs only have a single constraint, we show that a bisection-based dual approach can achieve a linear convergence rate.

1.2 Related Work

CMDPs Our work is related to policy-based CMDP algorithms (Altman, 1999; Borkar, 2005; Bhatnagar and Lakshmanan, 2012; Chow et al., 2017; Ding et al., 2020, 2021; Xu et al., 2021; Chen et al., 2021; Efroni et al., 2020). The papers Ding et al. (2020) and Xu et al. (2021) are closely related to our work. In Ding et al. (2020), the authors propose a natural policy gradient primal-dual method for CMDPs and prove that it achieves global convergence with the rate $\mathcal{O}(1/\sqrt{T})$ for both the optimality gap and the constraint violation under the soft-max policy parameterization. The work Xu et al. (2021) achieves a similar global convergence rate as Ding et al. (2020) using a primal-based approach. However, the entropy regularization, which is an effective technique for unconstrained MDPs, is not used in these algorithms.

Entropy-regularized RL Maximum entropy reinforcement learning optimizes policies to jointly maximize the expected return and the expected entropy of the policy. This framework has been used in many contexts. It has been shown that the maximum entropy

formulation provides a substantial improvement in exploration and robustness (Ziebart, 2010). It is robust in the face of model and estimation errors (Haarnoja et al., 2017) in both on-policy and off-policy settings (Haarnoja et al., 2018). More recently, the theoretical results in Mei et al. (2020) and Cen et al. (2021) have shown that the entropy regularization can help policy-based methods improve the convergence rate and the sample complexity compared with standard MDPs without the entropy regularization. However, despite the tremendous successes of the entropy regularization in unconstrained MDPs, the impact of the entropy regularization for CMDPs remains unknown.

1.3 Notations

Let $\Delta(\mathcal{S})$ denote the probability simplex over the set \mathcal{S} , and let $|\mathcal{S}|$ denote its cardinality. For a set $T \subset \mathbb{R}^p$, let $\text{cl}(T)$ denote the closure of T . When the variable s follows the distribution ρ , we write it as $s \sim \rho$. Let $\mathbb{E}[\cdot]$ and $\mathbb{E}[\cdot | \cdot]$, respectively, denote the expectation and conditional expectation of a random variable. Let \mathbb{R} denote the set of real numbers. For a number $a \in \mathbb{R}$, let $\text{sign}(a)$ denote the sign of a , i.e. $\text{sign}(a) = +1$ if $a \geq 0$ and $\text{sign}(a) = -1$ if $a < 0$. Let $[n]$ denote the set $\{1, 2, \dots, n\}$. For a vector x , we use x^T to denote the transpose of x , and use x_i or $(x)_i$ to denote the i -th entry of x . When applying a scalar function to x , e.g. $\log x$, the operation is understood as entry-wise. For vectors x and y , we use $x \geq y$ to denote an entry-wise inequality. We use the convention that $\|x\|_1 = \sum_i |x_i|$, $\|x\|_2 = \sqrt{\sum_i x_i^2}$, and $\|x\|_\infty = \max_i |x_i|$. For a matrix A , we use A_{ij} to denote its (i, j) -th entry, and let $\|A\|_F = \sqrt{\sum_{i,j} A_{ij}^2}$. Let I_n denote the $n \times n$ identity matrix. For square matrices A and B , we use $A \geq B$ to denote that $A - B$ is positive semi-definite. For a function $f(x)$, let $\nabla_x f(x)$ (resp. $\nabla_{xx}^2 f(x)$) denote its gradient (resp. Hessian) with respect to x , and we may omit x in the subscript when it is clear from the context. Let $\arg \min f(x)$ (resp. $\arg \max f(x)$) denote any arbitrary global minimum (resp. global maximum) of $f(x)$. We use boldface symbols for constraint-related vectors, e.g. λ .

2 PROBLEM FORMULATION

Markov Decision Processes An infinite-horizon Markov Decision Process $\text{MDP}(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ with a finite state-action space is specified by: a finite state space \mathcal{S} ; a finite action space \mathcal{A} ; a transition dynamics $P: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$, where $P(s'|s, a)$ is the probability of transition from state s to state s' when action a is taken; a reward function $r: \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, where $r(s, a)$ is the instantaneous reward when taking action a in state s ; a discount factor $\gamma \in [0, 1)$. A policy $\pi: \mathcal{S} \rightarrow \Delta(\mathcal{A})$

represents that the decision rule the agent uses, i.e. the agent takes action a with probability $\pi(a|s)$ in state s . We can also interpret a policy π as a vector in $\Delta(\mathcal{A})^{|\mathcal{S}|} \subset \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$.

Given a policy π , the value function $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ is defined to characterize the discounted sum of the rewards earned under π , i.e.

$$V^\pi(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| \pi, s_0 = s \right], \quad \forall s \in \mathcal{S} \quad (1)$$

where the expectation is taken over all possible trajectories, in which $a_t \sim \pi(\cdot|s_t)$ and $s_{t+1} \sim P(\cdot|s_t, a_t)$. When the initial state is sampled from some distribution ρ , we slightly abuse the notation and define the value function as

$$V^\pi(\rho) := \mathbb{E}_{s \sim \rho} [V^\pi(s)] \quad (2)$$

One classical property of the value function is that it is sufficiently smooth with respect to the policy if we view $V^\pi(\rho)$ as a function from the policy space $\Delta(\mathcal{A})^{|\mathcal{S}|}$ to \mathbb{R} . Especially, $V^\pi(\rho)$ has the following Lipschitz property.

Lemma 2.1 *For arbitrary policies π_1 and π_2 , it holds*

$$|V^{\pi_1}(\rho) - V^{\pi_2}(\rho)| \leq \ell_c \|\pi_1 - \pi_2\|_2, \quad (3)$$

where $\ell_c = \frac{\sqrt{|\mathcal{A}|}}{(1-\gamma)^2}$.

Lemma 2.1 follows from the bounded gradient of $V^\pi(\rho)$ with respect to π . We refer the reader to the supplement in Appendix A for more details.

The action-value function (or Q-function) $Q^\pi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ under policy π is defined as

$$Q^\pi(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| \pi, s_0 = s, a_0 = a \right] \quad (4)$$

which can be interpreted as the expected total reward with an initial state $s_0 = s$ and an initial action $a_0 = a$. Since $r(s, a) \in [0, 1]$ by assumption, we have that both $Q^\pi(s, a)$ and $V^\pi(\rho)$ are bounded between $[0, 1/(1-\gamma)]$ for any $(s, a) \in \mathcal{S} \times \mathcal{A}$ and initial distribution ρ .

For theoretical analysis, it is useful to define the so-called discounted state visitation distribution $d_{s_0}^\pi$ of a policy π :

$$d_{s_0}^\pi(s) := (1-\gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s | \pi, s_0), \quad \forall s \in \mathcal{S} \quad (5)$$

and we write $d_\rho^\pi(s) := \mathbb{E}_{s_0 \sim \rho} [d_{s_0}^\pi(s)]$ as the visitation distribution when the initial state follows ρ .

Soft-max Parameterization Parameterization is commonly deployed to model unknown policies to help

with the optimization process. One natural choice is the soft-max parameterization:

$$\pi_\theta(a|s) := \frac{\exp(\theta(s, a))}{\sum_{a' \in \mathcal{A}} \exp(\theta(s, a'))}, \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A} \quad (6)$$

where $\theta \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ is an unconstrained vector. We denote the class of all soft-max parameterized policies by Π . This policy class is complete in the sense that its closure $\text{cl}(\Pi)$ contains all stationary policies. In what follows, we will discard the subscript θ and just write $\pi \in \Pi$, whenever it is clear from the context.

Entropy Regularization To encourage exploration and accelerate convergence to the optimal policy, entropy regularization is widely used in solving MDPs. In the regularized setting, the agent seeks to optimize the entropy-regularized value function

$$V_\tau^\pi(\rho) := V^\pi(\rho) + \tau \cdot \mathcal{H}(\rho, \pi), \quad (7)$$

where $\tau \geq 0$ specifies the weight of regularization and $\mathcal{H}(\rho, \pi)$ is the discounted entropy defined by

$$\mathcal{H}(\rho, \pi) := \mathbb{E} \left[\sum_{t=0}^{\infty} -\gamma^t \log \pi(a_t | s_t) \middle| \pi, s_0 \sim \rho \right]. \quad (8)$$

We can also define the Q-function under regularization, which is referred to as the soft Q-function

$$Q_\tau^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a)} [V_\tau^\pi(s')]. \quad (9)$$

Constrained MDP With Entropy Regularization In a Constrained Markov Decision Process CMDP $(\mathcal{S}, \mathcal{A}, P, r, \mathbf{g}, \mathbf{b}, \gamma)$, besides the reward function r , we have a utility function $\mathbf{g} = (g_1, \dots, g_n) : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]^n$ and a threshold $\mathbf{b} \in [0, 1/(1-\gamma)]^n$. Under entropy regularization, the agent seeks to maximize the regularized value function $V_\tau^\pi(\rho)$ without violating the utility constraint $U_{\mathbf{g}}^\pi(\rho) \geq \mathbf{b}$, where the discounted utility $U_{\mathbf{g}}^\pi(\rho) := (U_{g_1}^\pi(\rho), \dots, U_{g_n}^\pi(\rho)) \in \mathbb{R}^n$ is defined by

$$U_{g_i}^\pi(\rho) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t g_i(s_t, a_t) \middle| \pi, s_0 \sim \rho \right]. \quad (10)$$

Equivalently, the agent solves the optimization problem

$$\max_{\pi \in \Pi} V_\tau^\pi(\rho) \quad \text{s.t.} \quad U_{\mathbf{g}}^\pi(\rho) \geq \mathbf{b}, \quad (11)$$

where $V_\tau^\pi(\rho) := V^\pi(\rho) + \tau \mathcal{H}(\rho, \pi)$. Consider the associated Lagrangian function $L(\pi, \boldsymbol{\lambda})$ and the dual function $D(\boldsymbol{\lambda})$ defined as:

$$L(\pi, \boldsymbol{\lambda}) := V_\tau^\pi(\rho) + \boldsymbol{\lambda}^T (U_{\mathbf{g}}^\pi(\rho) - \mathbf{b}), \quad (12a)$$

$$D(\boldsymbol{\lambda}) := \max_{\pi \in \Pi} L(\pi, \boldsymbol{\lambda}), \quad (12b)$$

where $\boldsymbol{\lambda} \in \mathbb{R}^n$ is the dual variable. For brevity, we omit the dependency of L and D on ρ and τ in the notations.

It can be seen from (12a) that $L(\pi, \boldsymbol{\lambda})$ can be viewed as an entropy-regularized value function with the reward $r_{\boldsymbol{\lambda}}(s, a) := r(s, a) + \boldsymbol{\lambda}^T \mathbf{g}(s, a)$ subtracted by the scalar $\boldsymbol{\lambda}^T \mathbf{b}$. Let $V_{\boldsymbol{\lambda}}^{\pi}(\rho)$ (resp. $Q_{\boldsymbol{\lambda}}^{\pi}(s, a)$) denote the value function (resp. Q-function) with the reward function $r_{\boldsymbol{\lambda}}$, i.e. $V_{\boldsymbol{\lambda}}^{\pi}(\rho) := V_{\tau}^{\pi}(\rho) + \boldsymbol{\lambda}^T U_{\mathbf{g}}^{\pi}(\rho)$.

It is worth mentioning that (11) is non-convex due to its non-concave objective function and non-convex constraints, thus making the problem challenging to solve. Henceforth, we slightly abuse the notation and still denote π_{τ}^* as an arbitrary optimal policy to the constrained problem (11). Let $\boldsymbol{\lambda}^*$ denote an optimal multiplier, i.e.

$$\boldsymbol{\lambda}^* := \arg \min_{\boldsymbol{\lambda} \geq 0} D(\boldsymbol{\lambda}), \quad (13)$$

and $\pi_{\boldsymbol{\lambda}}$ be the Lagrangian maximizer associated with the multiplier $\boldsymbol{\lambda}$, i.e.

$$\pi_{\boldsymbol{\lambda}} := \arg \max_{\pi \in \Pi} L(\pi, \boldsymbol{\lambda}). \quad (14)$$

We use the shorthand notations $V_{\tau}^* := V_{\tau}^{\pi_{\tau}^*}(\rho)$ and $D_{\tau}^* := \min_{\boldsymbol{\lambda} \geq 0} D(\boldsymbol{\lambda})$. As before, we hide the dependency of $\boldsymbol{\lambda}^*$, $\pi_{\boldsymbol{\lambda}}$ on ρ and τ , as well as the dependency of π_{τ}^* , V_{τ}^* , D_{τ}^* on ρ .

3 PROPERTIES OF CMDP WITH ENTROPY REGULARIZATION

Despite its non-convex nature, entropy-regularized CMDPs enjoy desirable properties, which we will discuss below. We refer the reader to the supplement in Appendix A for all the proofs in this section.

Assume that the Slater condition holds, i.e. there exists a strictly feasible policy.

Assumption 3.1 (Slater Condition) *There exist a policy $\bar{\pi} \in \Pi$ and $\boldsymbol{\xi} > 0$ such that $U_{\mathbf{g}}^{\bar{\pi}}(\rho) - \mathbf{b} \geq \boldsymbol{\xi}$.*

The Slater condition is standard in constrained optimization. It holds when the feasible region contains an interior point. In practical, such a point is often easy to find given prior knowledge of the problem. One direct consequence of the Slater condition is the strong duality (Altman, 1999; Paternain et al., 2019).

Lemma 3.2 (Strong Duality) *Under Assumption 3.1, there exist a primal-dual pair $(\pi_{\tau}^*, \boldsymbol{\lambda}^*)$ such that $V_{\tau}^* = D_{\tau}^* = L(\pi_{\tau}^*, \boldsymbol{\lambda}^*)$.*

In the remainder of the paper, we always assume that $(\pi_{\tau}^*, \boldsymbol{\lambda}^*)$ is a primal-dual pair. From the strong duality, we can derive an upper bound on $\boldsymbol{\lambda}^*$.

Lemma 3.3 *Under Assumption 3.1, it holds that*

$$0 \leq \lambda_i^* \leq \frac{V_{\tau}^* - V_{\tau}^{\bar{\pi}}(\rho)}{\xi_i}, \quad \forall i \in [n]. \quad (15)$$

Define

$$\Lambda := \left\{ \boldsymbol{\lambda} \mid 0 \leq \lambda_i \leq \frac{V_{\tau}^* - V_{\tau}^{\bar{\pi}}(\rho)}{\xi_i}, \text{ for all } i \in [n] \right\}. \quad (16)$$

Since the dual function $D(\boldsymbol{\lambda})$ is always convex, Lemmas 3.2 and 3.3 together imply that, instead of directly solving the non-convex primal problem (11), one can seek to solve the convex dual problem

$$\min_{\boldsymbol{\lambda}} D(\boldsymbol{\lambda}) \quad \text{s.t. } \boldsymbol{\lambda} \in \Lambda. \quad (17)$$

However, there are two open problems that need to be addressed. The first one is that although algorithms in convex optimization can be used to solve the dual problem, it is not clear how fast they will converge without discovering key properties of the dual function. The second problem is that optimizing the dual function gives a dual optimality bound, while our goal is find a primal solution and analyze the primal optimality gap together with the constraint violation. In the following sections, we will show that the entropy-regularized CMDP has special structures that can be leveraged to address the above issues.

3.1 Dual Smoothness

In optimization, smoothness plays an important role in establishing the convergence rate of an algorithm. Recall that a function $f : X \rightarrow \mathbb{R}$ is said to be ℓ -smooth if

$$\|\nabla f(x_1) - \nabla f(x_2)\|_2 \leq \ell \|x_1 - x_2\|_2, \quad (18)$$

for all $x_1, x_2 \in X$. In constrained optimization, however, smoothness is not always guaranteed, even when the primal problem is convex (Necoara et al., 2019). In addition, while the subgradient of the dual function exists in general, the dual function is not always differentiable due to the non-uniqueness of Lagrangian multipliers.

By leveraging the entropy regularization, we will show that the dual function $D(\boldsymbol{\lambda})$ in CMDPs is both differentiable and ℓ -smooth for some constant $\ell > 0$, under the following assumption on the discounted state visitation distribution.

Assumption 3.4 *The discounted state visitation distribution d_{ρ}^{π} is uniformly bounded away from 0 for all $\pi \in \Pi$, i.e. there exists $d > 0$, such that $d_{\rho}^{\pi}(s) \geq d$, $\forall s \in \mathcal{S}$, $\pi \in \Pi$.*

Assumption 3.4 ensures that the MDP *sufficiently explores* the state space. Since $d_{\rho}^{\pi}(s) \geq (1 - \gamma)\rho(s)$, it is satisfied when the initial distribution ρ lies in the interior of the probability simplex $\Delta(\mathcal{S})$. Similar assumptions are used in the prior literature (Agarwal et al., 2021; Mei et al., 2020, 2021).

The following proposition is crucial for the development of our main result.

Proposition 3.5 *For all policy π and $\lambda \geq 0$, it holds that*

$$L(\pi_\lambda, \lambda) - L(\pi, \lambda) \geq \frac{\tau d}{2(1-\gamma) \ln 2} \|\pi - \pi_\lambda\|_2^2. \quad (19)$$

Under Assumption 3.4, Proposition 3.5 implies that the Lagrangian function $L(\pi, \lambda)$ has a negative curvature at π_λ in all directions. The proof relies on the soft sub-optimality lemma (cf. Lemma D.4) and a lower bound on the KL divergence (Cover, 1999).

With the quadratic lower bound given by Proposition 3.5, we derive the following result.

Proposition 3.6 *Under Assumptions 3.1 and 3.4, the dual function $D(\lambda)$ satisfies the following properties:*

1. $D(\lambda)$ is differentiable and

$$\begin{aligned} \nabla D(\lambda) &= U_{\mathbf{g}}^{\pi_\lambda}(\lambda) - \mathbf{b} \\ &= (U_{g_1}^{\pi_\lambda}(\lambda) - b_1, \dots, U_{g_n}^{\pi_\lambda}(\lambda) - b_n). \end{aligned} \quad (20)$$

2. $D(\lambda)$ is ℓ -smooth on Λ , where

$$\ell = \frac{2 \times \ln 2 \times (n|\mathcal{A}| + (1-\gamma)^2 \sqrt{n|\mathcal{A}|})}{\tau(1-\gamma)^3 d}. \quad (21)$$

Proposition 3.6 asserts that the dual function $D(\lambda)$ is not only differentiable but also smooth on Λ . This is desirable since, along with the convexity, it establishes an improved convergence rate compared with the slow convergence rate of sub-gradient methods. We provide a short proof sketch for Proposition 3.6 below:

1. As subgradients of the dual function always exist for continuous problems, the differentiability follows from the uniqueness of the Lagrangian maximizer π_λ for every $\lambda \in \Lambda$ (Floudas, 1995)¹.
2. The smoothness of $D(\lambda)$ is the joint result of the Lipschitz continuity of $U_{\mathbf{g}}^{\pi}(\rho)$ with respect to π (cf. Lemma 2.1) and the Lipschitz continuity of π_λ with respect to λ , i.e. $\|\pi_{\lambda_1} - \pi_{\lambda_2}\|_2 \leq \ell_\Lambda \|\lambda_1 - \lambda_2\|_2$ for some $\ell_\Lambda > 0$. To prove the latter, the main idea is to use the quadratic lower bound given by Proposition 3.5 to conclude that π_λ is a second-order strict local maximum. After that, we apply a standard result from perturbation analysis which states that π_λ is Lipschitz stable at λ (Bonnans and Shapiro, 2013).

¹Although more than enough, Proposition 3.5 under Assumption 3.4 provides an intuitive way to think about the uniqueness of π_λ .

3.2 Optimality Gap And Constraint Violation

Given a candidate solution π to the CMDP problem in (11), our primary measures of the quality of the solution π are the primal optimality gap $|V_\tau^\pi(\rho) - V_\tau^*|$, and the constraint violation $\max_{i \in [n]} [b_i - U_{g_i}^\pi(\rho)]_+$, where $[x]_+ := \max\{x, 0\}$. However, dual-descent based methods could only guarantee a convergence bound in terms of the dual optimality gap $D(\lambda) - D_\tau^*$. In general, there is no guarantee that an ε -optimal dual solution λ , namely $D(\lambda) - D_\tau^* \leq \varepsilon$, would imply an $\mathcal{O}(\varepsilon^k)$ bound either on the primal optimality gap or on the constraint violation for the associated primal solution π_λ defined in (14), for some $k \in (0, 1]$.

However, in light of the entropy regularization, it is possible to show that an ε error bound for dual functions would yield an $\mathcal{O}(\sqrt{\varepsilon})$ error bound for the primal optimality gap and the constraint violation. We summarize the results in the following proposition.

Proposition 3.7 *Suppose that Assumptions 3.1 and 3.4 hold. If $\lambda \geq 0$ is an ε -optimal multiplier, i.e. $D(\lambda) - D_\tau^* \leq \varepsilon$, then there exist constants C_1 and C_2 such that the associated Lagrangian maximizer π_λ satisfies*

$$\|\pi_\lambda - \pi_\tau^*\|_2 \leq C_1 \sqrt{\varepsilon}, \quad (22a)$$

$$|V_\tau^{\pi_\lambda}(\rho) - V_\tau^*| \leq 2\varepsilon + \ell_c C_1 C_2 \sqrt{\varepsilon}, \quad (22b)$$

$$\max_{i \in [n]} [b_i - U_{g_i}^{\pi_\lambda}(\rho)]_+ \leq \ell_c C_1 \sqrt{\varepsilon}, \quad (22c)$$

where ℓ_c is the Lipschitz constant defined in (3).

The values of the problem-dependent constants C_1 and C_2 can be found in Appendix A.

In a nutshell, Proposition 3.7 enables the conversion of the dual optimality bound to primal metrics of interests, at the cost of enlarging the sub-optimality by a square root. This is a non-trivial result and it does not hold in a general setting without the entropy regularization. The proof of (22a) relies on the quadratic lower bound given in Proposition 3.5. Then, using the Lipschitz continuity of $U_{\mathbf{g}}^\pi(\rho)$ with respect to π (cf. Lemma 2.1), we can derive the bound (22c) on the constraint violation. Finally, (22b) can be obtained with some primal-dual properties.

4 FIRST-ORDER DUAL-DESCENT ALGORITHM

As shown in Section 3, the dual function $D(\lambda)$ for entropy-regularized CMDPs enjoys desirable properties, including the differentiability, the smoothness and the decomposition of the dual optimality gap. These favorable properties of entropy-regularized CMDPs motivate

us to use a dual-descent approach to solve the dual problem (17). In particular, we choose first-order methods, e.g. gradient projection method or Frank-Wolfe algorithm, while using the Natural Policy Gradient (NPG) algorithm as a subroutine for evaluating $D(\lambda)$ as well as $\nabla D(\lambda)$. To streamline the presentation, we mainly focus on the gradient projection method with the Nesterov acceleration as an example in the remainder of the paper. We begin with a brief introduction about the NPG algorithm.

4.1 Preliminary Tools

NPG Algorithm With Entropy Regularization

To optimize an unconstrained value function with respect to the policy, one commonly used first-order method is the Natural Policy Gradient algorithm (Kakade, 2001), which deploys a pre-conditioned gradient update and regularizes the descent direction by the Fisher-information matrix \mathcal{F}_ρ^θ (cf. Appendix B.1):

$$\theta \leftarrow \theta + \eta (\mathcal{F}_\rho^\theta)^\dagger \nabla_\theta V^{\pi_\theta}(\rho), \quad (23)$$

where η is the step-size and $(A)^\dagger$ denotes the Moore–Penrose inverse of a matrix A .

In the entropy regularized setting, the update scheme is obtained by replacing $\nabla_\theta V^{\pi_\theta}(\rho)$ in (23) with $\nabla_\theta V_\tau^{\pi_\theta}(\rho)$. Under the soft-max parameterization, the associated policy update has a fairly direct form, which is surprisingly independent from the initial distribution ρ :

$$\pi^{(t+1)}(a|s) \propto (\pi^{(t)}(a|s))^{1 - \frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta Q_\tau^{\pi^{(t)}}(s, a)}{1-\gamma}\right), \quad (24)$$

where we use the shorthand $\pi^{(t)}$ for the soft-max parameterized policy with respect to $\theta^{(t)}$, and Q_τ^π is the soft Q-function defined in (9). The right-hand side of (24) can be normalized by multiplying a factor $Z^{(t)}(s)$, defined as

$$Z^{(t)}(s) := \sum_{a \in \mathcal{A}} (\pi^{(t)}(a|s))^{1 - \frac{\eta\tau}{1-\gamma}} \exp\left(\frac{\eta Q_\tau^{\pi^{(t)}}(s, a)}{1-\gamma}\right), \quad (25)$$

to make $\pi^{(t+1)}$ be a valid distribution.

Cen et al. (2021) proved the global linear convergence of the entropy-regularized NPG method with a constant step-size. In particular, the error bound $\|\log \pi_\tau^* - \log \pi^{(t)}\|_\infty \leq \varepsilon$ can be achieved in

$$\frac{1}{1-\gamma} \log\left(\frac{2\|Q_\tau^* - Q_\tau^{\pi^{(0)}}\|_\infty}{\varepsilon\tau}\right), \quad (26)$$

iterations with the step-size $\eta = (1-\gamma)/\tau$, where π_τ^* is the optimal policy and $Q_\tau^*(s, a) := Q_\tau^{\pi_\tau^*}(s, a)$ is the associated optimal Q-function. Furthermore, they proved

that the convergence rate becomes quadratic around the optimum. We refer the reader to Appendix B.1 for more details.

Accelerated Gradient Projection Method with Inexact Gradient

The Nesterov acceleration is a momentum-based approach that can be used to modify a gradient descent-type method to improve its convergence (Nesterov, 1983, 2013). Consider the optimization problem

$$\min_x f(x) \quad \text{s.t.} \quad x \in X \quad (27)$$

where $f(x)$ is convex and differentiable, and X is a convex set. The accelerated gradient projection method takes the update rule

$$\begin{cases} x^{(k+1)} = \mathcal{P}_X(y^{(k)} - \alpha^k \nabla f(y^{(k)})) \\ y^{(k)} = x^{(k)} + \beta_k (x^{(k)} - x^{(k-1)}) \end{cases}, \quad k = 0, 1, \dots \quad (28)$$

where \mathcal{P}_X denotes the projection onto the set X , defined as $\mathcal{P}_X(y) := \arg \min_{x \in X} \|x - y\|_2$, and $\{\beta_k\}$ is chosen in a particular way to accelerate the convergence. The iteration (28) first computes an extrapolation point $y^{(k)}$ and then performs the gradient projection update on $y^{(k)}$ to find the next point $x^{(k+1)}$. It coincides with the standard gradient projection method when $\beta_k = 0$. For a convex and smooth function f , the accelerated gradient projection method (28) achieves an error bound of $\mathcal{O}(1/T^2)$ in T iterations (Nesterov, 2013).

When the gradient evaluation is inexact with a bounded error δ , i.e. we have access to some function $h : X \rightarrow \mathbb{R}^n$ such that $\|\nabla f(x) - h(x)\|_2 \leq \delta$ for all $x \in X$, Schmidt et al. (2011) proved that the accelerated gradient projection method still works with the slightly different error bound $\mathcal{O}(1/T^2 + T^2\delta^2 + \delta)$. Despite the $\mathcal{O}(1/T^2)$ shrinking term, there is an accumulated error incurred by the inexact gradient. We refer the reader to Proposition B.3 in Appendix B for a formal statement.

4.2 Accelerated Gradient Projection Method With NPG Subroutine

Before presenting our method, we first note that the feasible region Λ , as defined in (16), makes the dual problem (17) amenable to many constrained optimization methods. Especially, the projection operator $\mathcal{P}_\Lambda(\cdot)$ maps a point λ coordinate-wisely onto Λ such that

$$(\mathcal{P}_\Lambda(\lambda))_i = \text{Median}\left\{0, \frac{V_\tau^* - V_\tau^{\bar{\pi}}(\rho)}{\xi_i}, \lambda_i\right\} \quad (29)$$

where $\text{Median}\{\cdot, \cdot, \cdot\}$ returns the median of the input numbers.

Algorithm 1 Accelerated Gradient Projection Method with NPG Subroutine

- 1: **Input:** Initialization $\lambda^{(-1)}, \lambda^{(0)}, \tilde{\pi}_{\mu^{(-1)}}$; step-size $\{\alpha_k\}_{k \geq 0}$, η ; extrapolation weight $\{\beta_k\}_{k \geq 0}$; maximum number of iterations N_1, N_2, N_3 .
 - 2: **for** $t = 0, 1, 2, \dots, N_1 - 1$ **do**
 - 3: Compute the extrapolation point: $\mu^{(t)} = \lambda^{(t)} + \beta_k (\lambda^{(t)} - \lambda^{(t-1)})$.
 - 4: Estimate the optimal policy $\pi_{\mu^{(t)}}$ for problem (14) through the natural policy gradient subroutine: $\tilde{\pi}_{\mu^{(t)}} \leftarrow \text{NPG}_{Sub}(\mu^{(t)}, \tilde{\pi}_{\mu^{(t-1)}}, \eta, N_2)$.
 - 5: Compute the approximate gradient at $\mu^{(t)}$: $\tilde{\nabla} D(\mu^{(t)}) := U_{\tilde{\pi}_{\mu^{(t)}}}^{\mu^{(t)}}(\rho) - \mathbf{b}$.
 - 6: Take a gradient projection step at $\mu^{(t)}$: $\lambda^{(t+1)} \leftarrow \mathcal{P}_{\Lambda}(\mu^{(t)} - \alpha^k \tilde{\nabla} D(\mu^{(t)}))$, as defined by (29).
 - 7: **end for**
 - 8: Recover the policy from the dual variable: $\tilde{\pi}_{\lambda^{(N_1)}} \leftarrow \text{NPG}_{Sub}(\lambda^{(N_1)}, \tilde{\pi}_{\mu^{(N_1-1)}}, \eta, N_3)$.
-

The proposed method works in two loops. In the outer loop, we perform the accelerated gradient projection method on the dual function $D(\lambda)$, whereas we use the natural policy gradient method in the inner loop to evaluate $D(\lambda)$ by maximizing the Lagrangian $L(\pi, \lambda)$ with respect to π . We summarize the details of our method in Algorithm 1.

Specifically, in line 3, we compute the extrapolation point $\mu^{(t)}$. Then, in line 4, we estimate the corresponding Lagrangian maximizer $\pi_{\mu^{(t)}}$, defined in 14, using the natural policy gradient subroutine, which is displayed in Algorithm 2. With the estimated policy $\tilde{\pi}_{\mu^{(t)}}$, we evaluate the gradient $\nabla D(\mu^{(t)})$ by substituting the policy into the utility function in line 5. In line 6, we perform the gradient projection update at $\mu^{(t)}$ using the estimated gradient $\tilde{\nabla} D(\mu^{(t)})$. We remark that, as V_{τ}^* is generally unknown, the projection \mathcal{P}_{Λ} may not be precisely done in practical. Alternatively, one can perform the projection onto $\tilde{\Lambda} := \{\lambda \mid 0 \leq \lambda_i \leq (2 + 2\tau \log \mathcal{A}) / ((1 - \gamma)\xi_i), \text{ for all } i \in [n]\}$. Since the difference $V_{\tau}^* - V_{\tau}^{\tilde{\pi}}(\rho)$ is upper bounded by $(2 + 2\tau \log \mathcal{A}) / (1 - \gamma)$, it holds that $\lambda^* \in \Lambda \subseteq \tilde{\Lambda}$. This reduction would not influence the order of convergence. Finally, upon the termination of the outer loop, we recover the primal variable (policy) from the dual variable by running the NPG subroutine for N_3 iterations in line 8.

5 CONVERGENCE ANALYSIS

In this section, we analyze the convergence of Algorithm 1. The complete proofs of results in this section are postponed to Appendix B.3.

Algorithm 2 Natural Policy Gradient Subroutine (NPG_{Sub})

- 1: **Input:** Multiplier λ ; initial policy $\pi^{(0)}$; step-size η ; maximum number of iterations N .
 - 2: **for** $t = 0, 1, 2, \dots, N - 1$ **do**
 - 3: Compute the soft Q-function associated with the Lagrangian: $Q_{\lambda}^{\pi^{(t)}}$ of policy $\pi^{(t)}$.
 - 4: Update the policy with $Q_{\lambda}^{\pi^{(t)}}$ through (24).
 - 5: **end for**
-

A simple insight is that the NPG subroutine in Algorithm 1 computes the optimal policy at a linear rate in the inner loop and the accelerated gradient projection algorithm converges in $\mathcal{O}(1/\sqrt{\varepsilon})$ rate in the outer loop, leading to the overall convergence rate of $\tilde{\mathcal{O}}(1/\sqrt{\varepsilon})$. Then, we obtain the desired $\tilde{\mathcal{O}}(1/\varepsilon)$ rate in terms of the primal optimality gap and constraint violation by applying Proposition 3.7.

The above high-level technique requires subtle technicalities to be addressed here. We begin with the following proposition, which evaluates the accuracy of the gradient estimator defined in line 5 of Algorithm 1.

Proposition 5.1 *Suppose that π is an approximate solution to (14) such that $\|\log \pi - \log \pi_{\lambda}\|_{\infty} \leq \varepsilon$. The gradient estimator defined as $\tilde{\nabla} D(\lambda) := U_{\tilde{\pi}}^{\pi}(\rho) - \mathbf{b} = (U_{g_1}^{\pi}(\rho) - b_1, \dots, U_{g_n}^{\pi}(\rho) - b_n)$ satisfies the inequality $\|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_2 \leq \sqrt{n}|\mathcal{A}|\varepsilon / (1 - \gamma)^2$.*

This result can be deduced from the inequality $\|\pi - \pi_{\lambda}\|_{\infty} \leq \|\log \pi - \log \pi_{\lambda}\|_{\infty} \leq \varepsilon$ and the performance difference lemma of an unregularized value function (Lemma D.3). Recall that $L(\pi, \lambda) = V_{\lambda}^{\pi}(\rho) - \lambda^T \mathbf{b}$, where $V_{\lambda}^{\pi}(\rho) = V_{\tau}^{\pi}(\rho) + \lambda^T U_{\mathbf{g}}^{\pi}(\rho)$ is the value function with the reward $r_{\lambda}(s, a) := r(s, a) + \lambda^T \mathbf{g}(s, a)$ (cf. Section 2)). Therefore, Proposition 5.1 together with (26), implies that running Algorithm 2 with the step-size $\eta = (1 - \gamma)/\tau$ for

$$\frac{1}{1 - \gamma} \log \left(\frac{2\sqrt{n}|\mathcal{A}| \|Q_{\lambda}^{\pi} - Q_{\lambda}^{\pi^{(0)}}\|_{\infty}}{\delta(1 - \gamma)^2 \tau} \right) \quad (30)$$

iterations can guarantee a δ -accurate gradient estimation $\tilde{\nabla} D(\lambda)$, i.e. $\|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_2 \leq \delta$.

Below, we present our main convergence result of Algorithm 1.

Theorem 5.2 *Suppose that Assumptions 3.1 and 3.4 hold. For every $\varepsilon_1 > 0$, there exist some constants C_1 and $C_2 > 0$ such that Algorithm 1 with a random initialization and the parameters $\eta = (1 - \gamma)/\tau$, $\alpha_k = 1/\ell$, $\beta_k = (k - 1)/(k + 2)$, $N_1 = T$, $N_2 = \mathcal{O}(\log T)$ and*

$N_3 = \mathcal{O}(\log 1/\varepsilon_1)$ returns a solution pair (π, λ) such that

$$D(\lambda) - D_\tau^* \leq \varepsilon_0, \quad (31a)$$

$$\|\pi - \pi_\tau^*\|_2 \leq C_1 \sqrt{\varepsilon_0} + \varepsilon_1, \quad (31b)$$

$$|V_\tau^\pi(\rho) - V_\tau^*| \leq 2\varepsilon_0 + \ell_c C_1 C_2 \sqrt{\varepsilon_0} + \left(\ell_c C_2 + \frac{3\gamma}{2\tau\sqrt{n}} \right) \varepsilon_1, \quad (31c)$$

$$\max_{i \in [n]} [b_i - U_{g_i}^\pi(\rho)]_+ \leq \ell_c (C_1 \sqrt{\varepsilon_0} + \varepsilon_1), \quad (31d)$$

where

$$\varepsilon_0 = \frac{2\ell}{(T+1)^2} \left(\|\lambda^{(0)} - \lambda^*\|_2 + 1 \right)^2, \quad (32)$$

and where ℓ is the smoothness factor defined in (21) and ℓ_c is the Lipschitz constant defined in (3). The total iteration complexity is $N_1 \times N_2 + N_3 = \tilde{\mathcal{O}}(T)$ with primal error bounds $\mathcal{O}(1/T)$ given by (31b)-(31d), and a dual error bound $\mathcal{O}(1/T^2)$ given by (31a).

The values of the parameters N_2 , N_3 , and problem-dependent constants C_1 , C_2 can be found in Theorem B.5, Appendix B.3, where we restate the theorem.

Theorem 5.2 shows that Algorithm 1 achieves a global convergence with the rate $\tilde{\mathcal{O}}(1/\varepsilon)$. Specifically, it shows that with $\tilde{\mathcal{O}}(T)$ number of iterations in total, Algorithm 1 generates a solution with $\mathcal{O}(1/T)$ error bounds in terms of the policy (primal variable), primal optimality gap, and constraint violation, as well as $\mathcal{O}(1/T^2)$ error in terms of the dual optimality gap.

We briefly describe the intuition behind the proof of Theorem 5.2 below.

1. Firstly, the linear convergence of the natural policy gradient method (cf. Proposition B.1) and Proposition 5.1 imply that running the NPG subroutine for $N_2 = \mathcal{O}(\log T)$ iterations in the inner loop guarantees a sufficiently accurate estimation of $\nabla D(\lambda)$.
2. Then, we apply the convergence result by Schmidt et al. (2011) (refer to Section 4.1 and Appendix B.2), which implies the dual optimality gap $\mathcal{O}(1/T^2 + T^2\delta^2 + \delta)$, where δ is the gradient estimation error. Since the NPG subroutine converges linearly, we can suppress the constant in $N_2 = \mathcal{O}(\log T)$, and make δ sufficiently small such that the dual optimality gap equals $\mathcal{O}(1/T^2)$.
3. Let λ denote the dual variable returned by the for loop. Running the NPG subroutine in line 8 for additional $N_3 = \mathcal{O}(\log(1/\varepsilon_1))$ iterations guarantees an ε_1 -approximate solution π for the Lagrangian maximizer π_λ . Again, we can make ε_1 sufficiently small by suppressing the constant in N_3 .

4. Finally, by applying Proposition 3.7 and the triangular inequality, we prove (31b), stating that π is $\mathcal{O}(1/T)$ -optimal. We bound the constraint violation (31d) by using the Lipschitz continuity of $U_{\mathbf{g}}^\pi(\rho)$ with respect to π , and bound the primal optimality gap (31c) by using some primal-dual properties².

Remark 5.3 (Quadratic Convergence of NPG)

Our analysis of Algorithm 1 in this section is inspired by the global linear convergence of the entropy-regularized NPG method. However, as Cen et al. (2021) proved, the NPG method achieves a quadratic convergence around the optimum (cf. Proposition B.2). Therefore, it may be possible to improve the hidden constants in $N_2 = \mathcal{O}(\log T)$ and $N_3 = \mathcal{O}(\log(1/\varepsilon_1))$ under extra assumptions.

So far, we have only studied the entropy-regularized CMDP. However, adding entropy induces bias to the optimal solution of the standard unregularized CMDP. A standard way to deal with this mismatch issue is to choose the regularization parameter τ to be sufficiently small. The following corollary shows that we can compute a near-optimal policy with the rate $\tilde{\mathcal{O}}(1/\sqrt{T})$ for both the optimality gap and the constraint violation for the standard CMDP.

Corollary 5.4 *Suppose that Assumptions 3.1 and 3.4 hold. Then, Algorithm 1 with the choice $\tau = \mathcal{O}(\varepsilon)$ computes a solution π for the standard CMDP such that*

$$|V^{\pi^*}(\rho) - V^\pi(\rho)| = \mathcal{O}(\varepsilon), \quad (33a)$$

$$\max_{i \in [n]} [b_i - U_{g_i}^\pi(\rho)]_+ = \mathcal{O}(\varepsilon), \quad (33b)$$

in $\tilde{\mathcal{O}}(1/\varepsilon^2)$ iterations, where π^* is an optimal policy to the standard CMDP.

The proof of Corollary 5.4 relies on the following sandwich bound

$$V^{\pi_\tau^*}(\rho) \leq V^{\pi^*}(\rho) \leq V^{\pi_\tau^*}(\rho) + \frac{\tau}{1-\gamma} \log |\mathcal{A}|. \quad (34)$$

6 CMDPS WITH A SINGLE CONSTRAINT

Since we can convert the dual optimality bound to primal metrics of interests with little extra effort (cf. Proposition 3.7), the overall complexity relies on how fast we can solve the dual problem (17). For the special case where $n = 1$, the dual problem amounts to optimizing a convex function on an closed interval, which

²The techniques are similar to those in the proof of Proposition 3.7.

can be efficiently solved by the bisection method. Due to space restrictions, we refer the reader to Appendix C for more details about the algorithm (cf. Algorithm 3). We state the result in the following theorem.

Theorem 6.1 *Suppose that Assumptions 3.1 and 3.4 hold. When $n = 1$, for every $\varepsilon_0, \varepsilon_1 > 0$, Algorithm 3 returns a solution pair (π, λ) satisfying (31a)-(31d) in at most $\mathcal{O}(\log^2(1/\varepsilon_0) + \log(1/\varepsilon_1))$ iterations.*

By leveraging the linear convergence, we can derive a result analogous to Corollary 5.4 for the standard CMDP, but with a linear rate.

Corollary 6.2 *Suppose that Assumptions 3.1 and 3.4 hold. Then, Algorithm 3 with the choice $\tau = \mathcal{O}(\varepsilon)$ computes a solution π for the standard CMDP satisfying (33a)-(33b), in $\mathcal{O}(\log^2(1/\varepsilon))$ iterations.*

We refer the reader to Appendix C for the formal statements of Theorem 6.1 and Corollary 6.2 as well as their proofs.

7 CONCLUSION

In this paper, we showed that entropy regularization induces desirable properties to CMDPs from an optimization perspective. In particular, the Lagrangian dual function of CMDPs is smooth and an $\mathcal{O}(\varepsilon)$ error bound for the dual optimality gap yields an $\mathcal{O}(\sqrt{\varepsilon})$ error bound for the primal optimality gap and the constraint violation. In addition, we proposed a novel accelerated dual-descent algorithm and proved that it achieves a global convergence with the rate $\tilde{\mathcal{O}}(1/\varepsilon)$ for both the optimality gap and the constraint violation. It remains as an open question whether similar improved convergence results for the entropy-regularized CMDPs can be obtained with a sample-based policy gradient.

ACKNOWLEDGEMENTS

This work was supported by grants from AFOSR, ARO, ONR, NSF and C3.ai Digital Transformation Institute.

References

Achiam, J., Held, D., Tamar, A., and Abbeel, P. (2017). Constrained policy optimization. In *International Conference on Machine Learning*, pages 22–31. PMLR.

Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. (2021). On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *Journal of Machine Learning Research*, 22(98):1–76.

Altman, E. (1999). *Constrained Markov decision processes*, volume 7. CRC Press.

Bhatnagar, S. and Lakshmanan, K. (2012). An online actor-critic algorithm with function approximation for constrained markov decision processes. *Journal of Optimization Theory and Applications*, 153(3):688–708.

Bonnans, J. F. and Shapiro, A. (2013). *Perturbation analysis of optimization problems*. Springer Science & Business Media.

Borkar, V. S. (2005). An actor-critic algorithm for constrained markov decision processes. *Systems & control letters*, 54(3):207–213.

Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. (2021). Fast global convergence of natural policy gradient methods with entropy regularization.

Chen, Y., Dong, J., and Wang, Z. (2021). A primal-dual approach to constrained markov decision processes. *arXiv preprint arXiv:2101.10895*.

Chow, Y., Ghavamzadeh, M., Janson, L., and Pavone, M. (2017). Risk-constrained reinforcement learning with percentile risk criteria. *The Journal of Machine Learning Research*, 18(1):6070–6120.

Cover, T. M. (1999). *Elements of information theory*. John Wiley & Sons.

Ding, D., Wei, X., Yang, Z., Wang, Z., and Jovanovic, M. (2021). Provably efficient safe exploration via primal-dual policy optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 3304–3312. PMLR.

Ding, D., Zhang, K., Basar, T., and Jovanovic, M. R. (2020). Natural policy gradient primal-dual method for constrained markov decision processes. In *Conference on Neural Information Processing Systems*.

Efroni, Y., Mannor, S., and Pirotta, M. (2020). Exploration-exploitation in constrained mdps. *arXiv preprint arXiv:2003.02189*.

Fisac, J. F., Akametalu, A. K., Zeilinger, M. N., Kaynama, S., Gillula, J., and Tomlin, C. J. (2018). A general safety framework for learning-based control in uncertain robotic systems. *IEEE Transactions on Automatic Control*, 64(7):2737–2752.

Floudas, C. A. (1995). *Nonlinear and mixed-integer optimization: fundamentals and applications*. Oxford University Press.

- Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. (2017). Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pages 1352–1361. PMLR.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pages 1861–1870. PMLR.
- Kakade, S. M. (2001). A natural policy gradient. *Advances in neural information processing systems*, 14.
- Mei, J., Gao, Y., Dai, B., Szepesvari, C., and Schuurmans, D. (2021). Leveraging non-uniformity in first-order non-convex optimization. *arXiv preprint arXiv:2105.06072*.
- Mei, J., Xiao, C., Szepesvari, C., and Schuurmans, D. (2020). On the global convergence rates of softmax policy gradient methods. In *International Conference on Machine Learning*, pages 6820–6829. PMLR.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. (2016). Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pages 1928–1937. PMLR.
- Nachum, O., Norouzi, M., Xu, K., and Schuurmans, D. (2017). Bridging the gap between value and policy based reinforcement learning. *arXiv preprint arXiv:1702.08892*.
- Necoara, I., Patrascu, A., and Glineur, F. (2019). Complexity of first-order inexact lagrangian and penalty methods for conic convex programming. *Optimization Methods and Software*, 34(2):305–335.
- Nesterov, Y. (2013). Gradient methods for minimizing composite functions. *Mathematical Programming*, 140(1):125–161.
- Nesterov, Y. E. (1983). A method for solving the convex programming problem with convergence rate $o(1/k^2)$. In *Dokl. akad. nauk Sssr*, volume 269, pages 543–547.
- Paternain, S., Calvo-Fullana, M., Chamon, L. F., and Ribeiro, A. (2019). Safe policies for reinforcement learning via primal-dual methods. *arXiv preprint arXiv:1911.09101*.
- Schmidt, M., Roux, N. L., and Bach, F. (2011). Convergence rates of inexact proximal-gradient methods for convex optimization. *arXiv preprint arXiv:1109.2415*.
- Uchibe, E. and Doya, K. (2007). Constrained reinforcement learning from intrinsic and extrinsic rewards. In *2007 IEEE 6th International Conference on Development and Learning*, pages 163–168. IEEE.
- Williams, R. J. and Peng, J. (1991). Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268.
- Xu, T., Liang, Y., and Lan, G. (2021). Crpo: A new approach for safe reinforcement learning with convergence guarantee. In *International Conference on Machine Learning*, pages 11480–11491. PMLR.
- Zang, H., Li, X., Zhang, L., Zhao, P., and Wang, M. (2020). Teac: Integrating trust region and max entropy actor critic for continuous control.
- Zhang, X., Zhang, K., Miehling, E., and Basar, T. (2019). Non-cooperative inverse reinforcement learning. *Advances in Neural Information Processing Systems*, 32.
- Ziebart, B. D. (2010). *Modeling purposeful adaptive behavior with the principle of maximum causal entropy*. Carnegie Mellon University.

Supplementary Material: A Dual Approach to Constrained Markov Decision Processes with Entropy Regularization

A Proofs of Results in Sections 2 and 3

Lemma A.1 (Restatement of Lemma 2.1) *For an unregularized value function $V^\pi(\rho)$ with the reward function $r(s, a) \in [0, 1]$, it holds that*

$$|V^{\pi_1}(\rho) - V^{\pi_2}(\rho)| \leq \ell_c \|\pi_1 - \pi_2\|_2, \quad (35)$$

for arbitrary policies π_1 and π_2 , where $\ell_c = \sqrt{|\mathcal{A}|}/(1-\gamma)^2$.

Proof. It follows from the policy gradient for the direct parameterization (Lemma D.1) that

$$\frac{\partial V^\pi(\rho)}{\partial \pi(a|s)} = \frac{1}{1-\gamma} d_\rho^\pi(s) Q^\pi(s, a). \quad (36)$$

Thus, we can bound $\nabla_\pi V^\pi(\rho)$ as

$$\begin{aligned} \|\nabla_\pi V^\pi(\rho)\|_2 &= \frac{1}{1-\gamma} \sqrt{\sum_{s \in \mathcal{S}, a \in \mathcal{A}} (d_\rho^\pi(s) Q^\pi(s, a))^2} \\ &\leq \frac{\max_{s \in \mathcal{S}, a \in \mathcal{A}} Q^\pi(s, a)}{1-\gamma} \sqrt{\sum_{s \in \mathcal{S}, a \in \mathcal{A}} (d_\rho^\pi(s))^2} \\ &\stackrel{(i)}{\leq} \frac{\sqrt{|\mathcal{A}|}}{(1-\gamma)^2} \|d_\rho^\pi(\cdot)\|_2 \\ &\stackrel{(ii)}{\leq} \frac{\sqrt{|\mathcal{A}|}}{(1-\gamma)^2} =: \ell_c, \end{aligned} \quad (37)$$

where (i) uses $Q^\pi(s, a) \leq 1/(1-\gamma)$ and (ii) is because of $\|d_\rho^\pi(\cdot)\|_2 \leq \|d_\rho^\pi(\cdot)\|_1 = 1$. Then, (35) follows from

$$|V^{\pi_1}(\rho) - V^{\pi_2}(\rho)| \leq \sup_\pi \{\|\nabla_\pi V^\pi(\rho)\|_2\} \|\pi_1 - \pi_2\|_2 \leq \ell_c \|\pi_1 - \pi_2\|_2. \quad (38)$$

This completes the proof. \square

Lemma A.2 (Restatement of Lemma 3.2) *Under Assumption 3.1, there exist a primal-dual pair (π_τ^*, λ^*) such that $V_\tau^* = D_\tau^* = L(\pi_\tau^*, \lambda^*)$.*

Proof. We refer the reader to (Paternain et al., 2019) for a proof of the strong duality. \square

Lemma A.3 (Restatement of Lemma 3.3) *Under Assumption 3.1, it holds that*

$$0 \leq \lambda_i^* \leq \frac{V_\tau^* - V_\tau^{\bar{\pi}}(\rho)}{\xi_i}, \quad \forall i \in [n]. \quad (39)$$

Proof. Let $C \in \mathbb{R}$. For every $\lambda \geq 0$ such that $D(\lambda) \leq C$, one can write

$$\begin{aligned} C \geq D(\lambda) &\stackrel{(i)}{\geq} V_\tau^{\bar{\pi}}(\rho) + \sum_{i=1}^n \lambda_i (U_{g_i}^{\bar{\pi}}(\rho) - b_i) \\ &\stackrel{(ii)}{\geq} V_\tau^{\bar{\pi}}(\rho) + \sum_{i=1}^n \lambda_i \xi_i, \end{aligned} \quad (40)$$

where (i) follows from the definition of $D(\boldsymbol{\lambda})$ and (ii) is due to Assumption 3.1. Since $\boldsymbol{\xi} > 0$ and $\boldsymbol{\lambda} \geq 0$, (40) gives rise to the bound $0 \leq \lambda_i \leq (C - V_\tau^\pi(\rho)) / \xi_i$, for $i = 1, 2, \dots, n$. Now, by letting $C = V_\tau^*$, it results from the strong duality that $\{\boldsymbol{\lambda} \geq 0 \mid D(\boldsymbol{\lambda}) \leq C\}$ becomes the set of optimal dual variables. This completes the proof. \square

Proposition A.4 (Restatement of Proposition 3.5) *For all policy π and $\boldsymbol{\lambda} \geq 0$, it holds that*

$$L(\pi_{\boldsymbol{\lambda}}, \boldsymbol{\lambda}) - L(\pi, \boldsymbol{\lambda}) \geq \frac{\tau d}{2(1-\gamma) \ln 2} \|\pi - \pi_{\boldsymbol{\lambda}}\|_2^2. \quad (41)$$

Proof. Recall that $L(\pi, \boldsymbol{\lambda}) = V_{\boldsymbol{\lambda}}^\pi(\rho) - \boldsymbol{\lambda}^T \mathbf{b}$, where $V_{\boldsymbol{\lambda}}^\pi(\rho) = V_\tau^\pi(\rho) + \boldsymbol{\lambda}^T U_{\mathbf{g}}^\pi(\rho)$ is the value function with the reward $r_{\boldsymbol{\lambda}}(s, a) := r(s, a) + \boldsymbol{\lambda}^T \mathbf{g}(s, a)$ (cf. Section 2). The soft sub-optimality gap (Lemma D.4) then gives

$$L(\pi_{\boldsymbol{\lambda}}, \boldsymbol{\lambda}) - L(\pi, \boldsymbol{\lambda}) = V_{\boldsymbol{\lambda}}^{\pi_{\boldsymbol{\lambda}}}(\rho) - V_{\boldsymbol{\lambda}}^\pi(\rho) = \frac{\tau}{1-\gamma} \sum_{s \in \mathcal{S}} d_\rho^\pi(s) D_{\text{KL}}[\pi(\cdot|s) \mid \pi_{\boldsymbol{\lambda}}(\cdot|s)], \quad (42)$$

where $D_{\text{KL}}[P(\cdot) \mid Q(\cdot)] := \sum_x P(x) (\log P(x) - \log Q(x))$ is the KL divergence between probability distributions $P(\cdot)$ and $Q(\cdot)$. Now, we use a well-known bound relating the KL divergence to the vector 1-norm (Cover, 1999):

$$D_{\text{KL}}[P(\cdot) \mid Q(\cdot)] \geq \frac{1}{2 \ln 2} \|P(\cdot) - Q(\cdot)\|_1^2. \quad (43)$$

Combine (42) and (43) yields

$$L(\pi_{\boldsymbol{\lambda}}, \boldsymbol{\lambda}) - L(\pi, \boldsymbol{\lambda}) \geq \frac{\tau}{2(1-\gamma) \ln 2} \sum_{s \in \mathcal{S}} d_\rho^\pi(s) \|\pi(\cdot|s) - \pi_{\boldsymbol{\lambda}}(\cdot|s)\|_1^2 \stackrel{(i)}{\geq} \frac{\tau d}{2(1-\gamma) \ln 2} \|\pi - \pi_{\boldsymbol{\lambda}}\|_2^2, \quad (44)$$

where (i) is due to $\|\cdot\|_1 \geq \|\cdot\|_2$ and Assumption 3.4. This completes the proof. \square

Proposition A.5 (Restatement of Proposition 3.6) *Under Assumptions 3.1 and 3.4, the dual function $D(\boldsymbol{\lambda})$ satisfies the following properties:*

1. $D(\boldsymbol{\lambda})$ is differentiable and

$$\nabla D(\boldsymbol{\lambda}) = U_{\mathbf{g}}^{\pi_{\boldsymbol{\lambda}}}(\boldsymbol{\lambda}) - \mathbf{b} = (U_{g_1}^{\pi_{\boldsymbol{\lambda}}}(\boldsymbol{\lambda}) - b_1, \dots, U_{g_n}^{\pi_{\boldsymbol{\lambda}}}(\boldsymbol{\lambda}) - b_n). \quad (45)$$

2. $D(\boldsymbol{\lambda})$ is ℓ -smooth on Λ , where

$$\ell = \frac{2 \times \ln 2 \times (n|\mathcal{A}| + (1-\gamma)^2 \sqrt{n|\mathcal{A}|})}{\tau(1-\gamma)^3 d}. \quad (46)$$

The proof of Proposition A.5 relies on the following result by Bonnans and Shapiro (2013).

Proposition A.6 (Lipschitz stability of parametric local maximizers) *Given a set $T \subset \mathbb{R}^p$, consider a parametric optimization problem $P(t)$ with $t \in T$, stated as*

$$\max_{x \in F} f(x, t) \quad \text{s.t. } F = \{x \in \mathbb{R}^n \mid h_j(x) \leq 0, j = 1, 2, \dots, m\}, \quad (47)$$

where f, h_j are twice continuously differentiable and $F \neq \emptyset$. For every $\bar{t} \in T$, if $\bar{x} = x(\bar{t})$ is a strict local maximizer of $P(\bar{t})$ of order 2, i.e. $\nabla_{xx}^2 f(\bar{x}, \bar{t}) \geq w_1 I_n$ for some $w_1 > 0$, then there exist $\varepsilon, \delta, L > 0$ such that for all $t \in B_\varepsilon(\bar{t}) := \{t \mid \|t - \bar{t}\|_2 < \varepsilon\}$, there exists at least one local maximizer $x(t) \in B_\delta(\bar{x})$ of $P(t)$ and for each such local maximizer we have

$$\|x(t) - \bar{x}\|_2 \leq L \|t - \bar{t}\|_2. \quad (48)$$

Especially, taking $L = w_2/w_1$ fulfills the requirement, with

$$w_2 = \max_{z \in \text{cl}(B_\delta(\bar{x}))} [\|\nabla_{xt}^2 f(z, \bar{t})\|_F + 1]. \quad (49)$$

Proof of Proposition A.5. We first prove the differentiability of $D(\boldsymbol{\lambda})$. For a fixed $\boldsymbol{\lambda}$, solving for $\pi_{\boldsymbol{\lambda}} = \arg \max_{\pi \in \Pi} L(\pi, \boldsymbol{\lambda})$ is equivalent to solving an unconstrained MDP with entropy regularization (cf. Section 2):

$$\max_{\pi \in \Pi} V_{\boldsymbol{\lambda}}^{\pi}(\rho). \quad (50)$$

As shown in (Nachum et al., 2017), $\pi_{\boldsymbol{\lambda}}$ can be uniquely characterized as

$$\pi_{\boldsymbol{\lambda}}(a|s) \propto \exp\left(\frac{Q_{\boldsymbol{\lambda}}^{\pi_{\boldsymbol{\lambda}}}(s, a) - V_{\boldsymbol{\lambda}}^{\pi_{\boldsymbol{\lambda}}}(s)}{\tau}\right), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}. \quad (51)$$

Therefore, a standard result in the duality theory (Floudas, 1995) implies that $D(\boldsymbol{\lambda})$ is differentiable with the gradient

$$\nabla D(\boldsymbol{\lambda}) = U_{\mathbf{g}}^{\pi_{\boldsymbol{\lambda}}}(\boldsymbol{\lambda}) - \mathbf{b} = (U_{g_1}^{\pi_{\boldsymbol{\lambda}}}(\boldsymbol{\lambda}) - b_1, \dots, U_{g_n}^{\pi_{\boldsymbol{\lambda}}}(\boldsymbol{\lambda}) - b_n). \quad (52)$$

Next, we show that $\nabla D(\boldsymbol{\lambda})$ is Lipschitz continuous on Λ , which implies smoothness. Consider the statements:

1. $(U_{g_i}^{\pi}(\rho) - b_i)$ is ℓ_c -Lipschitz continuous with respect to π , which is already proved in Lemma 2.1.
2. $\pi_{\boldsymbol{\lambda}}$ is ℓ_{Λ} -Lipschitz continuous with respect to $\boldsymbol{\lambda}$ for some $\ell_{\Lambda} > 0$.

If these statements hold true, it follows that

$$|(U_{g_i}^{\pi_{\boldsymbol{\lambda}_1}}(\rho) - b_i) - (U_{g_i}^{\pi_{\boldsymbol{\lambda}_2}}(\rho) - b_i)| \leq \ell_c \|\pi_{\boldsymbol{\lambda}_1} - \pi_{\boldsymbol{\lambda}_2}\|_2 \leq \ell_c \ell_{\Lambda} \|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2\|_2, \quad (53)$$

which leads to

$$\|\nabla D(\boldsymbol{\lambda}_1) - \nabla D(\boldsymbol{\lambda}_2)\|_2 \leq \sqrt{n} \ell_c \ell_{\Lambda} \|\boldsymbol{\lambda}_1 - \boldsymbol{\lambda}_2\|_2, \quad (54)$$

i.e. $D(\boldsymbol{\lambda})$ is $\sqrt{n} \ell_c \ell_{\Lambda}$ -strongly smooth on Λ .

To prove Statement 2, consider the Lagrangian $L(\pi, \boldsymbol{\lambda})$, which is twice continuously differentiable on $(0, 1)^{|\mathcal{S}| \times |\mathcal{A}|} \times \Lambda$. The hidden constraint for the maximization problem (12b) is linear and has the form

$$\sum_{a \in \mathcal{A}} \pi(a|s) = 1, \quad \forall s \in \mathcal{S}. \quad (55)$$

By Proposition 3.5, it holds that

$$\nabla_{\pi\pi} L(\pi_{\boldsymbol{\lambda}}, \boldsymbol{\lambda}) \geq \frac{\tau d}{2(1-\gamma) \ln 2} I_{|\mathcal{S}| \times |\mathcal{A}|}, \quad (56)$$

which implies that $\pi_{\boldsymbol{\lambda}}$ is a strict global maximizer of order 2 under Assumption 3.4.

Consider $\nabla_{\pi\boldsymbol{\lambda}}^2 L(\pi, \boldsymbol{\lambda})$, which is a matrix of dimension $|\mathcal{S}| |\mathcal{A}| \times n$. Specifically, it holds

$$\frac{\partial^2 L(\pi, \boldsymbol{\lambda})}{\partial \pi(a|s) \partial \lambda_i} \stackrel{(i)}{=} \frac{\partial}{\partial \pi(a|s)} (U_{g_i}^{\pi}(\rho) - b_i) \stackrel{(ii)}{=} \frac{1}{1-\gamma} d_{\rho}^{\pi}(s) Q_{g_i}^{\pi}(s, a), \quad (57)$$

where (i) follows from definition (12a) and (ii) is due to the policy gradient (cf. Lemma D.1).

Following the same argument as in the proof of Proposition A.1, we have

$$\|\nabla_{\pi\boldsymbol{\lambda}}^2 L(\pi, \boldsymbol{\lambda})\|_F \leq \frac{\sqrt{n|\mathcal{A}|}}{(1-\gamma)^2}. \quad (58)$$

Therefore, applying Proposition A.6 with

$$w_1 = \frac{\tau d}{2(1-\gamma) \ln 2}, \quad w_2 = \frac{\sqrt{n|\mathcal{A}|}}{(1-\gamma)^2} + 1, \quad (59)$$

we conclude that $\pi_{\boldsymbol{\lambda}}$ is locally ℓ_{Λ} -Lipschitz continuous with respect to $\boldsymbol{\lambda}$ for all $\boldsymbol{\lambda} \in \Lambda$, where $\ell_{\Lambda} = w_2/w_1$. Since ℓ_{Λ} is universal and does not depend on $\boldsymbol{\lambda}$, the local Lipschitz property is ready to be extended to Λ . The proof is completed by setting $\ell = \sqrt{n} \ell_c \ell_{\Lambda}$. \square

Proposition A.7 (Restatement of Proposition 3.7) *Suppose that Assumptions 3.1 and 3.4 hold. If $\lambda \geq 0$ is an ε -optimal multiplier, i.e. $D(\lambda) - D_\tau^* \leq \varepsilon$, then the associated Lagrangian maximizer π_λ satisfies*

$$\|\pi_\lambda - \pi_\tau^*\|_2 \leq C_1 \sqrt{\varepsilon}, \quad (60a)$$

$$|V_\tau^{\pi_\lambda}(\rho) - V_\tau^*| \leq 2\varepsilon + \ell_c C_1 C_2 \sqrt{\varepsilon}, \quad (60b)$$

$$\max_{i \in [n]} [b_i - U_{g_i}^{\pi_\lambda}(\rho)]_+ \leq \ell_c C_1 \sqrt{\varepsilon}, \quad (60c)$$

where

$$\ell_c = \frac{\sqrt{|\mathcal{A}|}}{(1-\gamma)^2}, \quad C_1 = \sqrt{\frac{2(1-\gamma)\ln 2}{\tau d}}, \quad C_2 = (V_\tau^* - V_\tau^{\bar{\pi}}(\rho)) \left(\sum_{i=1}^n \frac{1}{\xi_i} \right). \quad (61)$$

Proof. We can write $D(\lambda) = L(\pi_\lambda, \lambda)$ and $D_\tau^* = L(\pi_\tau^*, \lambda^*)$, where (π_τ^*, λ^*) is any primal-dual pair. Then, by the strong duality (Lemma 3.2), we have

$$L(\pi_\tau^*, \lambda^*) = \min_{\mu \geq 0} L(\pi_\tau^*, \mu) \leq L(\pi_\tau^*, \lambda). \quad (62)$$

Therefore,

$$\begin{aligned} \varepsilon &\geq L(\pi_\lambda, \lambda) - L(\pi_\tau^*, \lambda^*) = L(\pi_\lambda, \lambda) - \min_{\mu \geq 0} L(\pi_\tau^*, \mu) \\ &\geq L(\pi_\lambda, \lambda) - L(\pi_\tau^*, \lambda) \\ &\stackrel{(i)}{\geq} \frac{\tau d}{2(1-\gamma)\ln 2} \|\pi_\lambda - \pi_\tau^*\|_2^2, \end{aligned} \quad (63)$$

where (i) results from the quadratic lower bound given by Proposition 3.5. Then, (60a) is obtained after rearranging the terms in (63).

Next, we can use the Lipschitz continuity of the utility function (cf. Lemma 2.1) to bound the constraint violation. For every $i = 1, 2, \dots, n$, it holds that

$$\left| U_{g_i}^{\pi_\lambda}(\rho) - U_{g_i}^{\pi_\tau^*}(\rho) \right| \leq \ell_c \|\pi_\lambda - \pi_\tau^*\|_2 \leq \ell_c C_1 \sqrt{\varepsilon}. \quad (64)$$

As the optimal policy π_τ^* must be feasible to (11), i.e. $U_{\mathbf{g}}^{\pi_\tau^*}(\rho) \geq \mathbf{b}$, we can bound the constraint violation as

$$\max_{i \in [n]} [b_i - U_{g_i}^{\pi_\lambda}(\rho)]_+ \leq \max_{i \in [n]} \left\{ [b_i - U_{g_i}^{\pi_\tau^*}(\rho)]_+ + \left| U_{g_i}^{\pi_\lambda}(\rho) - U_{g_i}^{\pi_\tau^*}(\rho) \right| \right\} \leq \ell_c C_1 \sqrt{\varepsilon}. \quad (65)$$

Finally, to bound the primal optimality gap, we note that

$$0 \stackrel{(i)}{\leq} L(\pi_\tau^*, \lambda) - L(\pi_\tau^*, \lambda^*) \stackrel{(ii)}{\leq} L(\pi_\lambda, \lambda) - L(\pi_\tau^*, \lambda^*) = D(\lambda) - D(\lambda^*) \leq \varepsilon, \quad (66)$$

where (i) follows from the strong duality and (ii) is due to the definition of π_λ (cf. (14)). Thus, by expanding the Lagrangian as

$$\begin{aligned} L(\pi_\tau^*, \lambda) - L(\pi_\tau^*, \lambda^*) &= V_\tau^{\pi_\tau^*}(\rho) + \lambda^T \left(U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - \mathbf{b} \right) - V_\tau^{\pi_\tau^*}(\rho) - (\lambda^*)^T \left(U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - \mathbf{b} \right) \\ &= (\lambda - \lambda^*)^T \left(U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - \mathbf{b} \right), \end{aligned} \quad (67)$$

and applying the complementary slackness $(\lambda^*)^T \left(U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - \mathbf{b} \right) = 0$, we obtain the bound

$$0 \leq (\lambda)^T \left(U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - \mathbf{b} \right) \leq \varepsilon. \quad (68)$$

Therefore,

$$\begin{aligned} \left| (\lambda)^T \left(U_{\mathbf{g}}^{\pi_\lambda}(\rho) - \mathbf{b} \right) \right| &\stackrel{(i)}{\leq} \left| (\lambda)^T \left(U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - \mathbf{b} \right) \right| + \left| (\lambda)^T \left(U_{\mathbf{g}}^{\pi_\lambda}(\rho) - U_{\mathbf{g}}^{\pi_\tau^*}(\rho) \right) \right| \\ &\stackrel{(ii)}{\leq} \varepsilon + \ell_c C_1 (V_\tau^* - V_\tau^{\bar{\pi}}(\rho)) \left(\sum_{i=1}^n \frac{1}{\xi_i} \right) \sqrt{\varepsilon}, \end{aligned} \quad (69)$$

where (i) is due to the triangular inequality and (ii) uses the bound (64) and the boundedness of Λ (cf. Lemma 3.3), i.e. $0 \leq \lambda_i \leq (V_\tau^* - V_\tau^{\bar{\pi}}(\rho))/\xi_i$ for all $i \in [n]$ and $\boldsymbol{\lambda} \in \Lambda$. Thus, we can bound the primal optimality gap as

$$\begin{aligned}
 |V_\tau^{\pi_\lambda}(\rho) - V_\tau^*| &= |V_\tau^{\pi_\lambda}(\rho) - V_\tau^{\pi_\tau^*}(\rho)| \\
 &\stackrel{(i)}{=} \left| \left[V_\tau^{\pi_\lambda}(\rho) + (\boldsymbol{\lambda})^T (U_{\mathbf{g}}^{\pi_\lambda}(\rho) - b) \right] - (\boldsymbol{\lambda})^T (U_{\mathbf{g}}^{\pi_\lambda}(\rho) - b) \right. \\
 &\quad \left. - \left[V_\tau^{\pi_\tau^*}(\rho) + (\boldsymbol{\lambda}^*)^T (U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - b) \right] \right| \\
 &\stackrel{(ii)}{\leq} |L(\pi_\lambda, \boldsymbol{\lambda}) - L(\pi_\tau^*, \boldsymbol{\lambda}^*)| + |(\boldsymbol{\lambda})^T (U_{\mathbf{g}}^{\pi_\lambda}(\rho) - b)| \\
 &\stackrel{(iii)}{\leq} \varepsilon + \left(\varepsilon + \ell_c C_1 (V_\tau^* - V_\tau^{\bar{\pi}}(\rho)) \left(\sum_{i=1}^n \frac{1}{\xi_i} \right) \sqrt{\varepsilon} \right) \\
 &= 2\varepsilon + \ell_c C_1 (V_\tau^* - V_\tau^{\bar{\pi}}(\rho)) \left(\sum_{i=1}^n \frac{1}{\xi_i} \right) \sqrt{\varepsilon} \\
 &= 2\varepsilon + \ell_c C_1 C_2 \sqrt{\varepsilon},
 \end{aligned} \tag{70}$$

where (i) uses the complementary slackness $(\boldsymbol{\lambda}^*)^T (U_{\mathbf{g}}^{\pi_\tau^*}(\rho) - b) = 0$ and (ii) uses the triangular inequality and the definition of Lagrangian (12a). In (iii), we use the assumption

$$D(\boldsymbol{\lambda}) - D(\boldsymbol{\lambda}^*) = L(\pi_\lambda, \boldsymbol{\lambda}) - L(\pi_\tau^*, \boldsymbol{\lambda}^*) \leq \varepsilon, \tag{71}$$

and the inequality (69). This completes the proof. \square

B Supplementary Materials for Sections 4 and 5

B.1 Entropy-regularized NPG

For entropy-regularized MDPs, the natural policy gradient update rule can be written as

$$\theta \leftarrow \theta + \eta (\mathcal{F}_\rho^\theta)^\dagger \nabla_\theta V_\tau^{\pi_\theta}(\rho), \tag{72}$$

where \mathcal{F}_ρ^θ is the Fisher information matrix, defined as

$$\mathcal{F}_\rho^\theta := \mathbb{E}_{s \sim d_\rho^{\pi_\theta}, a \sim \pi_\theta(\cdot|s)} \left[(\nabla_\theta \log \pi_\theta(a|s)) (\nabla_\theta \log \pi_\theta(a|s))^\top \right]. \tag{73}$$

Under the soft-max parameterization, the associated policy update has a fairly direct form (cf. (24) and (25)). We refer the reader to (Cen et al., 2021) for a detailed derivation.

Cen et al. (2021) proved that the entropy-regularized NPG method enjoys a global linear convergence and a local quadratic convergence. We summarize the two results in Propositions B.1 and B.2, where we abuse the notations and denote the optimal unconstrained value function with entropy regularization, the corresponding Q-function, and the associated optimal policy with V_τ^* , Q_τ^* , and π_τ^* respectively. Let μ_τ^* denote the stationary distribution over \mathcal{S} of the MDP under policy π_τ^* ³.

Proposition B.1 (Global linear convergence) *If the step-size $\eta = (1 - \gamma)/\tau$ is used, the entropy-regularized NPG algorithm (24) satisfies the error bounds:*

$$\|V_\tau^* - V_\tau^{\pi^{(t+1)}}\|_\infty \leq 3 \|Q_\tau^* - Q_\tau^{\pi^{(0)}}\|_\infty \gamma^{t+1}, \tag{74a}$$

$$\|\log \pi_\tau^* - \log \pi^{(t+1)}\|_\infty \leq 2 \|Q_\tau^* - Q_\tau^{\pi^{(0)}}\|_\infty \tau^{-1} \gamma^t, \tag{74b}$$

³It is straightforward to verify that $d_{\mu_\tau^*}^{\pi_\tau^*} = \mu_\tau^*$.

for all $t > 0$, where

$$\|V_\tau^* - V_\tau^{\pi^{(t+1)}}\|_\infty := \max_{s \in \mathcal{S}} |V_\tau^*(s) - V_\tau^{\pi^{(t+1)}}(s)|, \quad (75a)$$

$$\|Q_\tau^* - Q_\tau^{\pi^{(0)}}\|_\infty := \max_{s \in \mathcal{S}, a \in \mathcal{A}} |Q_\tau^*(s, a) - Q_\tau^{\pi^{(0)}}(s, a)|, \quad (75b)$$

$$\|\log \pi_\tau^* - \log \pi^{\pi^{(t+1)}}\|_\infty := \max_{s \in \mathcal{S}, a \in \mathcal{A}} |\log \pi_\tau^*(a|s) - \log \pi^{\pi^{(t+1)}}(a|s)|. \quad (75c)$$

We note that Cen et al. (2021) proved a more general result for all step-sizes $\eta \in [0, (1 - \gamma)/\tau]$, whereas the fastest convergence is achieved with the maximum step-size $\eta = (1 - \gamma)/\tau$.

Proposition B.2 (Local quadratic convergence) *Suppose that the entropy-regularized NPG algorithm (24) with the step-size $\eta = (1 - \gamma)/\tau$ satisfies*

$$\|\log \pi^{(t)} - \log \pi_\tau^*\|_\infty \leq 1, \quad (76)$$

for all $t \geq 0$. There exist problem-dependent constants K_1 and K_2 such that

$$V_\tau^*(\rho) - V_\tau^{(t)}(\rho) \leq K_1 \left(K_2 \left(V_\tau^*(\mu_\tau^*) - V_\tau^{\pi^{(0)}}(\mu_\tau^*) \right) \right)^{2^t}. \quad (77)$$

In our work, $V_\lambda^\pi(\rho)$ is the entropy-regularized value function associated with the Lagrangian $L(\pi, \lambda)$, which has the reward function $r_\lambda(s, a) = r(s, a) + \lambda^T \mathbf{g}(s, a)$. Therefore, Proposition B.1 implies that, with step-size $\eta = (1 - \gamma)/\tau$, the error bound $\|\log \pi_\lambda - \log \pi^{(t)}\|_\infty \leq \varepsilon$ can be achieved in

$$\frac{1}{1 - \gamma} \log \left(\frac{2 \|Q_\lambda^{\pi_\lambda} - Q_\lambda^{\pi^{(0)}}\|_\infty}{\varepsilon \tau} \right), \quad (78)$$

iterations (cf. (26)). Furthermore, since $\lambda \in \Lambda = \{\lambda \mid 0 \leq \lambda_i \leq (V_\tau^* - V_\tau^{\bar{\pi}}(\rho))/\xi_i, \text{ for all } i \in [n]\}$, we have that $r_\lambda(s, a) \in [0, 1 + C_2]$, where $C_2 = (V_\tau^* - V_\tau^{\bar{\pi}}(\rho)) (\sum_{i=1}^n 1/\xi_i)$. Together with the elementary entropy bound $\mathcal{H}(\rho, \pi) \in [0, \log |\mathcal{A}|/(1 - \gamma)]$, it holds that

$$Q_\lambda^\pi(s, a) \in \left[0, \frac{1 + C_2 + \tau \log |\mathcal{A}|}{1 - \gamma} \right], \quad (79)$$

for all $\lambda \in \Lambda$. Thus, we can drop the dependency of Q-function in (78) to obtain the following bound on the number of iterations:

$$\frac{1}{1 - \gamma} \log \left(\frac{2(1 + C_2 + \tau \log |\mathcal{A}|)}{\varepsilon \tau (1 - \gamma)} \right). \quad (80)$$

B.2 Accelerated Gradient Projection Method with Inexact Gradient

Gradient projection method is a feasible direction method for solving constrained optimization problems of the form:

$$\begin{aligned} \min_x & f(x) \\ \text{s.t. } & x \in X \end{aligned} \quad (81)$$

where $f(x)$ is convex and differentiable and X is convex. The general update scheme is

$$x^{(k+1)} = \mathcal{P}_X \left(x^{(k)} - \alpha^k \nabla f(x^{(k)}) \right). \quad (82)$$

When the gradient is inexact, Schmidt et al. (2011) proved the following bound for the general update (82) and the accelerated update (28).

Proposition B.3 (Convergence of inexact gradient projection method) *Assume that $f(x)$ is convex and L -smooth on X , and that we have access to a gradient oracle $h(x)$ such that $\|\nabla f(x) - h(x)\|_2 \leq \delta$ for*

all $x \in X$. Let $x^* = \arg \min_{x \in X} f(x)$. By selecting $\alpha_k = 1/L$ and $\beta_k = (k-1)/(k+2)$, then the iterates of algorithm (82) satisfy

$$f\left(\frac{1}{k} \sum_{i=1}^k x^{(i)}\right) - f(x^*) \leq \frac{L}{2k} \left(\|x^{(0)} - x^*\|_2 + \frac{2k\delta}{L} \right)^2 = O\left(\frac{1}{k}\right) + O(k^2\delta^2) + O(k\delta). \quad (83)$$

Moreover, for the accelerated version (28), it holds that

$$f(x^{(k)}) - f(x^*) \leq \frac{2L}{(k+1)^2} \left(\|x^{(0)} - x^*\|_2 + \frac{(k+1)k\delta}{L} \right)^2 = O\left(\frac{1}{k^2}\right) + O(k^2\delta^2) + O(\delta). \quad (84)$$

B.3 Proofs of Results in Section 5

Proposition B.4 (Restatement of Proposition 5.1) Suppose that π is an approximate solution to (14) such that $\|\log \pi - \log \pi_\lambda\|_\infty \leq \varepsilon$. The gradient estimator defined by

$$\tilde{\nabla} D(\lambda) := U_{\mathbf{g}}^\pi(\rho) - \mathbf{b} = (U_{g_1}^\pi(\rho) - b_1, \dots, U_{g_n}^\pi(\rho) - b_n), \quad (85)$$

satisfies

$$\|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_2 \leq \frac{\sqrt{n|\mathcal{A}|}}{(1-\gamma)^2} \varepsilon. \quad (86)$$

Proof. Since $(\log x)' = 1/x \geq 1$ for all $x \in (0, 1]$, it holds that $\|\pi - \pi_\lambda\|_\infty \leq \|\log \pi - \log \pi_\lambda\|_\infty \leq \varepsilon$. To bound the quantity $\|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_2 = \|U_{\mathbf{g}}^\pi(\rho) - U_{\mathbf{g}}^{\pi_\lambda}(\rho)\|_2$, we can either use the Lipschitz continuity of $U_{\mathbf{g}}^\pi(\rho)$ (cf. Lemma 2.1) or use the performance difference lemma (cf. Lemma D.3).

With the Lipschitz continuity, we have

$$|U_{g_i}^\pi(\rho) - U_{g_i}^{\pi_\lambda}(\rho)| \leq \ell_c \|\pi - \pi_\lambda\|_2 \leq \ell_c \sqrt{n} \|\pi - \pi_\lambda\|_\infty = \frac{\sqrt{n|\mathcal{A}|}}{(1-\gamma)^2} \varepsilon, \quad (87)$$

where we have used the inequality $\|\cdot\|_2 \leq \sqrt{n} \|\cdot\|_\infty$. Therefore,

$$\|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_2 \leq \sqrt{n} \|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_\infty = \sqrt{n} \|U_{\mathbf{g}}^\pi(\rho) - U_{\mathbf{g}}^{\pi_\lambda}(\rho)\|_\infty \leq \frac{n\sqrt{|\mathcal{A}|}}{(1-\gamma)^2} \varepsilon. \quad (88)$$

On the other hand, we can use the performance difference lemma to obtain

$$\begin{aligned} |U_{g_i}^\pi(\rho) - U_{g_i}^{\pi_\lambda}(\rho)| &= \left| \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} d_\rho^{\pi_\lambda}(s) \sum_{a \in \mathcal{A}} (\pi(a|s) - \pi_\lambda(a|s)) Q_{g_i}^\pi(s, a) \right| \\ &\leq \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} d_\rho^{\pi_\lambda}(s) \sum_{a \in \mathcal{A}} |\pi(a|s) - \pi_\lambda(a|s)| Q_{g_i}^\pi(s, a) \\ &\stackrel{(i)}{\leq} \frac{\varepsilon}{1-\gamma} \sum_{s \in \mathcal{S}} d_\rho^{\pi_\lambda}(s) \sum_{a \in \mathcal{A}} Q_{g_i}^\pi(s, a) \\ &\stackrel{(ii)}{\leq} \frac{|\mathcal{A}|}{(1-\gamma)^2} \varepsilon, \end{aligned} \quad (89)$$

where (i) is based on the bound $\|\pi - \pi_\lambda\|_\infty \leq \varepsilon$ and (ii) is due to $Q_{g_i}^\pi(s, a) \leq 1/(1-\gamma)$ and the fact that $d_\rho^{\pi_\lambda}(\cdot)$ is a probability distribution. Repeating (88) with the bound (89) yields that

$$\|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_2 \leq \sqrt{n} \|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_\infty = \sqrt{n} \|U_{\mathbf{g}}^\pi(\rho) - U_{\mathbf{g}}^{\pi_\lambda}(\rho)\|_\infty \leq \frac{\sqrt{n|\mathcal{A}|}}{(1-\gamma)^2} \varepsilon. \quad (90)$$

Equations (88) and (90) give two upper bounds on the quantity $\|\tilde{\nabla} D(\lambda) - \nabla D(\lambda)\|_2$. In this work, we use the bound (90), as it has a weaker dependence on the number of constraints n . This completes the proof. \square

Proposition 5.1 implies that running Algorithm 2 with the step-size $\eta = (1-\gamma)/\tau$ for

$$\frac{1}{1-\gamma} \log \left(\frac{2\sqrt{n|\mathcal{A}|}(1+C_2+\tau \log |\mathcal{A}|)}{\delta(1-\gamma)^3 \tau} \right). \quad (91)$$

iterations, where $C_2 = (V_\tau^* - V_\tau^{\bar{\pi}}(\rho)) (\sum_{i=1}^n 1/\xi_i)$, guarantees a δ -accurate gradient estimation $\tilde{\nabla}D(\boldsymbol{\lambda})$, i.e. $\|\tilde{\nabla}D(\boldsymbol{\lambda}) - \nabla D(\boldsymbol{\lambda})\|_2 \leq \delta$ (cf. (30)).

Theorem B.5 (Restatement of Theorem 5.2) *Suppose that Assumptions 3.1 and 3.4 hold. For every $\varepsilon_1 > 0$, Algorithm 1 with a random initialization and the parameters $\eta = (1 - \gamma)/\tau$, $\alpha_k = 1/\ell$, $\beta_k = (k - 1)/(k + 2)$, and*

$$N_1 = T, N_2 = \frac{1}{1 - \gamma} \log \left(\frac{2\sqrt{n}|\mathcal{A}|T(T + 1)(1 + C_2 + \tau \log |\mathcal{A}|)}{(1 - \gamma)^3 \tau \ell} \right), N_3 = \frac{1}{1 - \gamma} \log \left(\frac{2\sqrt{n}(1 + C_2 + \tau \log |\mathcal{A}|)}{\varepsilon_1 \tau (1 - \gamma)} \right), \quad (92)$$

returns a solution pair $(\pi, \boldsymbol{\lambda})$ such that

$$D(\boldsymbol{\lambda}) - D_\tau^* \leq \varepsilon_0, \quad (93a)$$

$$\|\pi - \pi_\tau^*\|_2 \leq C_1 \sqrt{\varepsilon_0} + \varepsilon_1, \quad (93b)$$

$$|V_\tau^\pi(\rho) - V_\tau^*| \leq 2\varepsilon_0 + \ell_c C_1 C_2 \sqrt{\varepsilon_0} + \left(\ell_c C_2 + \frac{3\gamma}{2\tau\sqrt{n}} \right) \varepsilon_1, \quad (93c)$$

$$\max_{i \in [n]} [b_i - U_{g_i}^\pi(\rho)]_+ \leq \ell_c (C_1 \sqrt{\varepsilon_0} + \varepsilon_1), \quad (93d)$$

where

$$\varepsilon_0 = \frac{2\ell}{(T + 1)^2} \left(\|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|_2 + 1 \right)^2, \quad \ell = \frac{2 \ln 2 \left(n|\mathcal{A}| + (1 - \gamma)^2 \sqrt{n|\mathcal{A}|} \right)}{\tau(1 - \gamma)^3 d}, \quad (94)$$

and

$$\ell_c = \frac{\sqrt{|\mathcal{A}|}}{(1 - \gamma)^2}, \quad C_1 = \sqrt{\frac{2(1 - \gamma) \ln 2}{\tau d}}, \quad C_2 = (V_\tau^* - V_\tau^{\bar{\pi}}(\rho)) \left(\sum_{i=1}^n \frac{1}{\xi_i} \right). \quad (95)$$

The total iteration complexity is $N_1 \times N_2 + N_3 = \tilde{\mathcal{O}}(T)$ with primal error bounds $\mathcal{O}(1/T)$ given by (93b)-(93d) and a dual error bound $\mathcal{O}(1/T^2)$ given by (93a).

Proof. Under Assumptions 3.1 and 3.4, it follows from Proposition 3.6 that $D(\boldsymbol{\lambda})$ is convex, differentiable, and ℓ -smooth on Λ . Now, we fix the gradient accuracy as

$$\delta = \frac{\ell}{T(T + 1)}. \quad (96)$$

Then, it follows from Proposition 5.1 and (91) that running the NPG subroutine in line 4 of Algorithm 1 for

$$N_2 = \frac{1}{1 - \gamma} \log \left(\frac{2\sqrt{n}|\mathcal{A}|T(T + 1)(1 + C_2 + \tau \log |\mathcal{A}|)}{(1 - \gamma)^3 \tau \ell} \right), \quad (97)$$

iterations guarantees obtaining an estimation $\tilde{\nabla}D(\boldsymbol{\lambda})$ such that

$$\|\tilde{\nabla}D(\boldsymbol{\lambda}) - \nabla D(\boldsymbol{\lambda})\|_2 \leq \frac{\ell}{T(T + 1)}, \quad (98)$$

where we have use the bound $\|Q_{\boldsymbol{\lambda}}^{\pi_{\boldsymbol{\lambda}}} - Q_{\boldsymbol{\lambda}}^{(0)}\|_\infty \leq (1 + C_2 + \log |\mathcal{A}|)/(1 - \gamma)$ for all $\boldsymbol{\lambda} \in \Lambda$ (cf. (79)). Therefore, by Proposition B.3, running the outer loop in Algorithm 1 for $N_1 = T$ iterations generates a solution $\boldsymbol{\lambda}^{(T)}$ such that

$$D(\boldsymbol{\lambda}^{(T)}) - D_\tau^* \leq \frac{2\ell}{(T + 1)^2} \left(\|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|_2 + \frac{(T + 1)T\delta}{\ell} \right)^2 = \frac{2\ell}{(T + 1)^2} \left(\|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|_2 + 1 \right)^2, \quad (99)$$

which satisfies (93a).

Below, we adopt the proof of Proposition 3.7 (cf. Appendix A). We first apply Proposition 3.7 with

$$\varepsilon_0 = \frac{2\ell}{(T + 1)^2} \left(\|\boldsymbol{\lambda}^{(0)} - \boldsymbol{\lambda}^*\|_2 + 1 \right)^2 \quad (100)$$

to obtain $\|\pi_{\lambda^{(T)}} - \pi_{\tau}^*\|_2 \leq C_1\sqrt{\varepsilon_0}$, where π_{τ}^* is an optimal policy. By Propositions B.1, we can compute an approximate Lagrangian maximizer $\tilde{\pi}_{\lambda^{(T)}}$ to (14) such that $\|\log \tilde{\pi}_{\lambda^{(T)}} - \log \pi_{\lambda^{(T)}}\|_{\infty} \leq \varepsilon_1/\sqrt{n}$, by running Algorithm 2 for

$$N_3 = \frac{1}{1-\gamma} \log \left(\frac{2\sqrt{n}(1+C_2+\tau \log |\mathcal{A}|)}{\varepsilon_1\tau(1-\gamma)} \right), \quad (101)$$

iterations (cf. (78)). Now, we show that $\tilde{\pi}_{\lambda^{(T)}}$ is a solution to (11) satisfying (93b)-(93c). Firstly, we have

$$\|\log \tilde{\pi}_{\lambda^{(T)}} - \log \pi_{\lambda^{(T)}}\|_2 \leq \sqrt{n} \|\log \tilde{\pi}_{\lambda^{(T)}} - \log \pi_{\lambda^{(T)}}\|_{\infty} \leq \varepsilon_1. \quad (102)$$

By applying the triangular inequality and using the strong concavity of the logarithm function on $(0, 1]$, it holds that

$$\|\tilde{\pi}_{\lambda^{(T)}} - \pi_{\tau}^*\|_2 \leq \|\tilde{\pi}_{\lambda^{(T)}} - \pi_{\lambda^{(T)}}\|_2 + \|\pi_{\lambda^{(T)}} - \pi_{\tau}^*\|_2 \leq C_1\sqrt{\varepsilon_0} + \varepsilon_1. \quad (103)$$

Then, we bound the constraint violation. It follows from the Lipschitz continuity of the utility function (cf. (3)) that

$$\left| U_{g_i}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - U_{g_i}^{\pi_{\tau}^*}(\rho) \right| \leq \ell_c \|\tilde{\pi}_{\lambda^{(T)}} - \pi_{\tau}^*\|_2 \leq \ell_c (C_1\sqrt{\varepsilon_0} + \varepsilon_1), \quad \forall i = 1, 2, \dots, n. \quad (104)$$

As the optimal policy π_{τ}^* must be a feasible solution to (11), i.e. $U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho) \geq \mathbf{b}$, the constraint violation is bounded as

$$\max_{i \in [n]} \left[b_i - U_{g_i}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) \right]_+ \leq \max_{i \in [n]} \left\{ \left[b_i - U_{g_i}^{\pi_{\tau}^*}(\rho) \right]_+ + \left| U_{g_i}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - U_{g_i}^{\pi_{\tau}^*}(\rho) \right| \right\} \leq \ell_c (C_1\sqrt{\varepsilon_0} + \varepsilon_1). \quad (105)$$

Finally, to bound the primal optimality gap, we note that

$$0 \stackrel{(i)}{\leq} L(\pi_{\tau}^*, \lambda^{(T)}) - L(\pi_{\tau}^*, \lambda^*) \stackrel{(ii)}{\leq} L(\pi_{\lambda^{(T)}}, \lambda^{(T)}) - L(\pi_{\tau}^*, \lambda^*) = D(\lambda^{(T)}) - D_{\tau}^* \leq \varepsilon_0, \quad (106)$$

where (i) follows from the strong duality and (ii) is due to the definition of $\pi_{\lambda^{(T)}}$. Thus, by expanding the Lagrangian as

$$\begin{aligned} L(\pi_{\tau}^*, \lambda^{(T)}) - L(\pi_{\tau}^*, \lambda^*) &= V_{\tau}^{\pi_{\tau}^*}(\rho) + (\lambda^{(T)})^T (U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho) - \mathbf{b}) - V_{\tau}^{\pi_{\tau}^*}(\rho) - (\lambda^*)^T (U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho) - \mathbf{b}) \\ &= (\lambda^{(T)} - \lambda^*)^T (U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho) - \mathbf{b}), \end{aligned} \quad (107)$$

and applying the complementary slackness $(\lambda^*)^T (U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho) - \mathbf{b}) = 0$, we obtain the bound

$$0 \leq (\lambda^{(T)})^T (U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho) - \mathbf{b}) \leq \varepsilon_0. \quad (108)$$

Therefore,

$$\begin{aligned} \left| (\lambda^{(T)})^T (U_{\mathbf{g}}^{\pi_{\lambda^{(T)}}}(\rho) - \mathbf{b}) \right| &\stackrel{(i)}{\leq} \left| (\lambda^{(T)})^T (U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho) - \mathbf{b}) \right| + \left| (\lambda^{(T)})^T (U_{\mathbf{g}}^{\pi_{\lambda^{(T)}}}(\rho) - U_{\mathbf{g}}^{\pi_{\tau}^*}(\rho)) \right| \\ &\stackrel{(ii)}{\leq} \varepsilon_0 + \ell_c (V_{\tau}^* - V_{\tau}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho)) \left(\sum_{i=1}^n \frac{1}{\xi_i} \right) (C_1\sqrt{\varepsilon_0} + \varepsilon_1) \\ &= \varepsilon_0 + \ell_c C_2 (C_1\sqrt{\varepsilon_0} + \varepsilon_1), \end{aligned} \quad (109)$$

where (i) is based on the triangular inequality, and (ii) uses the bound (104) and the boundedness of Λ , i.e. $0 \leq \lambda_i \leq (V_{\tau}^* - V_{\tau}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho))/\xi_i$ for all $i \in [n]$ and $\lambda \in \Lambda$ (cf. Lemma 3.3). Thus, we can bound the primal optimality

gap as

$$\begin{aligned}
\left| V_{\tau}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - V_{\tau}^{\star} \right| &= \left| V_{\tau}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - V_{\tau}^{\pi_{\tau}^{\star}}(\rho) \right| \\
&\stackrel{(i)}{=} \left| \left[V_{\tau}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) + (\boldsymbol{\lambda}^{(T)})^T \left(U_{\mathbf{g}}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - \mathbf{b} \right) \right] - (\boldsymbol{\lambda}^{(T)})^T \left(U_{\mathbf{g}}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - \mathbf{b} \right) \right. \\
&\quad \left. - \left[V_{\tau}^{\pi_{\tau}^{\star}}(\rho) + (\boldsymbol{\lambda}^{\star})^T \left(U_{\mathbf{g}}^{\pi_{\tau}^{\star}}(\rho) - \mathbf{b} \right) \right] \right| \\
&\stackrel{(ii)}{\leq} \left| L(\tilde{\pi}_{\lambda^{(T)}}, \boldsymbol{\lambda}^{(T)}) - L(\pi_{\tau}^{\star}, \boldsymbol{\lambda}^{\star}) \right| + \left| (\boldsymbol{\lambda}^{(T)})^T \left(U_{\mathbf{g}}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - \mathbf{b} \right) \right| \\
&\stackrel{(iii)}{\leq} \left| L(\tilde{\pi}_{\lambda^{(T)}}, \boldsymbol{\lambda}^{(T)}) - L(\pi_{\lambda^{(T)}}, \boldsymbol{\lambda}^{(T)}) \right| + \left| L(\pi_{\lambda^{(T)}}, \boldsymbol{\lambda}^{(T)}) - L(\pi_{\tau}^{\star}, \boldsymbol{\lambda}^{\star}) \right| \\
&\quad + \left| (\boldsymbol{\lambda}^{(T)})^T \left(U_{\mathbf{g}}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - \mathbf{b} \right) \right| \\
&\stackrel{(iv)}{\leq} \left(\frac{3\gamma}{2\tau\sqrt{n}} \varepsilon_1 \right) + (\varepsilon_0) + (\varepsilon_0 + \ell_c C_2 (C_1 \sqrt{\varepsilon_0} + \varepsilon_1)) \\
&= 2\varepsilon_0 + \ell_c C_1 C_2 \sqrt{\varepsilon_0} + \left(\ell_c C_2 + \frac{3\gamma}{2\tau\sqrt{n}} \right) \varepsilon_1,
\end{aligned} \tag{110}$$

where (i) uses the complementary slackness $(\boldsymbol{\lambda}^{\star})^T (U_{\mathbf{g}}^{\pi_{\tau}^{\star}}(\rho) - \mathbf{b}) = 0$, (ii) uses the triangular inequality and the definition of Lagrangian (12a), and (iii) uses the triangular inequality again. The inequality (iv) contains three parts where the first part uses Proposition B.1 as

$$\left| L(\tilde{\pi}_{\lambda^{(T)}}, \boldsymbol{\lambda}^{(T)}) - L(\pi_{\lambda^{(T)}}, \boldsymbol{\lambda}^{(T)}) \right| = \left| V_{\lambda^{(T)}}^{\tilde{\pi}_{\lambda^{(T)}}}(\rho) - V_{\lambda^{(T)}}^{\pi_{\lambda^{(T)}}}(\rho) \right| \leq \frac{3\gamma}{2\tau} \|\log \tilde{\pi}_{\lambda^{(T)}} - \log \pi_{\lambda^{(T)}}\|_{\infty} = \frac{3\gamma}{2\tau\sqrt{n}} \varepsilon_1, \tag{111}$$

the second part uses the assumption

$$D(\boldsymbol{\lambda}^{(T)}) - D_{\tau}^{\star} = L(\pi_{\lambda^{(T)}}, \boldsymbol{\lambda}^{(T)}) - L(\pi_{\tau}^{\star}, \boldsymbol{\lambda}^{\star}) \leq \varepsilon_0, \tag{112}$$

and the third part uses the inequality (109). This completes the proof. \square

Corollary B.6 (Restatement of Corollary 5.4) *Suppose that Assumptions 3.1 and 3.4 hold. Let*

$$\tau = \frac{(1-\gamma)\varepsilon}{4\log|\mathcal{A}|}. \tag{113}$$

Then, Algorithm 1 computes a solution π for the standard CMDP such that

$$\left| V^{\pi^{\star}}(\rho) - V^{\pi}(\rho) \right| = \mathcal{O}(\varepsilon), \tag{114a}$$

$$\max_{i \in [n]} [b_i - U_{g_i}^{\pi}(\rho)]_+ = \mathcal{O}(\varepsilon), \tag{114b}$$

in $\tilde{\mathcal{O}}(1/\varepsilon^2)$ iterations, where π^{\star} is an optimal policy to the standard CMDP.

Proof. Since the total iteration complexity of Algorithm 1 is dominated by $N_1 \times N_2$ and the error bounds (31b)-(31d) are dominated by $\sqrt{\varepsilon_0}$ (cf. (94)), we ignore the effect of ε_0 , ε_1 and only focus on $\sqrt{\varepsilon_0}$ terms in (31b)-(31d) through the analysis below.

Firstly, by invoking the optimality of π_{τ}^{\star} with respect to the entropy-regularized CMDP and the elementary entropy bound $0 \leq \mathcal{H}(\rho, \pi) \leq \log|\mathcal{A}|/(1-\gamma)$, we obtain

$$V^{\pi_{\tau}^{\star}}(\rho) + \frac{\tau}{1-\gamma} \log|\mathcal{A}| \geq V^{\pi_{\tau}^{\star}}(\rho) + \tau \mathcal{H}(\rho, \pi_{\tau}^{\star}) = V_{\tau}^{\star}(\rho) \geq V_{\tau}^{\pi^{\star}}(\rho) \geq V^{\pi^{\star}}(\rho), \tag{115}$$

which implies the sandwich bound (34). Now, we choose T in such a way that $\ell_c C_1 C_2 \sqrt{\varepsilon_0} = \varepsilon/2$, where ℓ_c , C_1 , and C_2 are specified in Theorem B.5. Then, Theorem B.5 implies that running Algorithm 1 with

$$N_1 = T, \quad N_2 = \frac{1}{1-\gamma} \log \left(\frac{2\sqrt{n}|\mathcal{A}|T(T+1)(1+C_2+\tau\log|\mathcal{A}|)}{(1-\gamma)^{3\tau\ell}} \right) \tag{116}$$

returns a solution π such that $|V_\tau^* - V_\tau^\pi(\rho)| = \mathcal{O}(\varepsilon/2)$ (cf. (93c)). It then follows that

$$\left| V^{\pi^*}(\rho) - V^\pi(\rho) \right| \stackrel{(i)}{\leq} |V^*(\rho) - V_\tau^*(\rho)| + |V_\tau^*(\rho) - V_\tau^\pi(\rho)| + |V_\tau^\pi(\rho) - V^\pi(\rho)| \stackrel{(ii)}{\leq} \frac{2\tau \log |\mathcal{A}|}{1-\gamma} + \mathcal{O}\left(\frac{\varepsilon}{2}\right) = \mathcal{O}(\varepsilon), \quad (117)$$

where (i) is due to the triangular inequality and (ii) uses the bound (115) (cf. (34)). The $\mathcal{O}(\varepsilon)$ -constraint violation follows directly from Theorem 5.2, since it enjoys the same order of convergence as the primal optimality gap.

We note that $\ell_c C_1 C_2 \sqrt{\varepsilon_0}$ can be written as $C/(\tau(T+1))$, where

$$C = \ell_c (C_1 \cdot \sqrt{\tau}) C_2 \left(\sqrt{2\ell \cdot \tau} \right) \left(\|\lambda^{(0)} - \lambda^*\|_2 + 1 \right) \quad (118)$$

is a constant that does not depend on T and ε . Thus, the choice $\tau = [(1-\gamma)\varepsilon]/(4\log |\mathcal{A}|)$ implies that

$$\frac{C}{T+1} = \frac{(1-\gamma)\varepsilon}{4\log |\mathcal{A}|} \times \frac{\varepsilon}{2}, \quad (119)$$

i.e. $T = \mathcal{O}(1/\varepsilon^2)$. Since the total iteration complexity is $N_1 \times N_2 = T \times \mathcal{O}(\log T)$, we obtain the $\tilde{\mathcal{O}}(1/\sqrt{T})$ error bound for the primal optimality gap and the constraint violation. This completes the proof. \square

C Supplementary Materials for Section 6

In this subsection, we consider the special situation where there is a single constraint ($n = 1$). In particular, we use the non-bold notations to emphasize that the associated notations denote numbers instead of vectors used in the previous sections, e.g. multiplier λ , constraint $U_g^\pi(\rho) \geq b$, Slater condition $V_\tau^\pi(\rho) - b \geq \xi$.

Since the feasible region $\Lambda = [0, C_2]$, where $C_2 = (V_\tau^* - V_\tau^\pi(\rho))/\xi$, is bounded and $D(\lambda)$ is convex, an approximate stationary point is also an approximate optimal solution. Specifically, if $|\nabla D(\lambda)| < \varepsilon$, then

$$D(\lambda) - D_\tau^* \leq |\nabla D(\lambda)| \times |\lambda - \lambda^*| < C_2 \varepsilon. \quad (120)$$

Additionally, if $\text{sign}(\nabla D(0)) = \text{sign}(\nabla D(C_2)) = 1$, then $D(\lambda)$ attains the optimum at $\lambda = 0$ due to the convexity. Similarly, $D(\lambda)$ attains the optimum at $\lambda = C_2$ if $\text{sign}(\nabla D(0)) = \text{sign}(\nabla D(C_2)) = -1$. Therefore, it only remains to consider the case where $D(0) < 0$ and $D(C_2) > 0$.

The proposed method aims to find an approximate stationary point with the bisection scheme. For a given search interval, it computes the gradient at the middle point. If the gradient is greater than ε , it shrinks the search interval to the left-half interval; if the gradient is smaller than $-\varepsilon$, it shrinks the search interval to the right-half interval. The iterates terminate when it finds a point λ such that $|\tilde{\nabla} D(\lambda)| < \varepsilon$, where $\tilde{\nabla} D(\lambda)$ is the approximate gradient. We summarize the proposed method in Algorithm 3, where we separately define the gradient estimator (cf. lines 4 and 5 in Algorithm 1) as a new subroutine Grad_{Sub} (cf. Algorithm 4) for the ease of presentation. We also assume that the initialization is non-trivial in the sense that $\nabla D(0) < 0$ and $\nabla D(C_2) > 0$, since otherwise we can return the optimal solution λ^* as 0 or C_2 .

Below, we give a formal statement for the convergence result of Algorithm 3 (cf. Theorem 6.1).

Theorem C.1 (Restatement of Theorem 6.1) *Suppose that Assumptions 3.1 and 3.4 hold. When $n = 1$, for every $\varepsilon, \varepsilon_1 > 0$, Algorithm 3 with the parameters*

$$\eta = \frac{(1-\gamma)}{\tau}, \quad N_1 = \frac{1}{1-\gamma} \log \left(\frac{4|\mathcal{A}|(1+C_2+\tau \log |\mathcal{A}|)}{(1-\gamma)^3 \tau \varepsilon} \right), \quad N_2 = \frac{1}{1-\gamma} \log \left(\frac{2(1+C_2+\tau \log |\mathcal{A}|)}{\varepsilon_1 \tau (1-\gamma)} \right), \quad (121)$$

returns a solution (π, λ) in at most $\log_2(\ell C_2/\varepsilon)$ outer loops, such that

$$D(\lambda) - D_\tau^* \leq \varepsilon_0, \quad (122a)$$

$$\|\pi - \pi_\tau^*\|_2 \leq C_1 \sqrt{\varepsilon_0} + \varepsilon_1, \quad (122b)$$

$$|V_\tau^\pi(\rho) - V_\tau^*| \leq 2\varepsilon_0 + \ell_c C_1 C_2 \sqrt{\varepsilon_0} + \left(\ell_c C_2 + \frac{3\gamma}{2\tau} \right) \varepsilon_1, \quad (122c)$$

$$\left[b - U_g^\pi(\rho) \right]_+ \leq \ell_c (C_1 \sqrt{\varepsilon} + \varepsilon_1), \quad (122d)$$

Algorithm 3 Bisection Method with NPG Subroutine

```

1: Input: Initialization  $\pi$ ,  $p_0 = 0$ ,  $q_0 = C_2$ ; step-size  $\eta$ ; maximum number of iterations  $N_1$ ,  $N_2$ ; threshold  $\varepsilon$ .
2: for  $t = 0, 1, 2, \dots$  do
3:   Let  $\lambda = (p_t + q_t)/2$ .
4:   if  $|\text{Grad}_{Sub}(\lambda, \pi, \eta, N_1)| < \varepsilon$  then
5:     break
6:   else
7:     if  $\text{Grad}_{Sub}(\lambda, \pi, \eta, N_1) \geq \varepsilon$  then
8:       Let  $p_{t+1} \leftarrow p_t$  and  $q_{t+1} \leftarrow \lambda$ .
9:     else
10:      Let  $p_{t+1} \leftarrow \lambda$  and  $q_{t+1} \leftarrow q_t$ .
11:    end if
12:  end if
13: end for
14: Recover the policy from the dual variable:  $\tilde{\pi}_\lambda \leftarrow \text{NPG}_{Sub}(\lambda, \pi, \eta, N_2)$ .

```

Algorithm 4 Gradient Estimator (Grad_{Sub})

```

1: Input: Target point  $\lambda$ , initialization  $\pi$ , step-size  $\eta$ , maximum number of iterations  $N$ .
2: Estimate the optimal policy  $\pi_\lambda$  for problem (14) through the natural policy gradient subroutine:  $\tilde{\pi}_\lambda \leftarrow \text{NPG}_{Sub}(\lambda, \pi, \eta, N)$ .
3: Compute and output the approximate gradient at  $\lambda$ :  $\tilde{\nabla}D(\lambda) := U_g^{\tilde{\pi}_\lambda}(\rho) - b$ .

```

where

$$\varepsilon_0 = \frac{3C_2}{2}\varepsilon, \quad \ell = \frac{2 \ln 2 \left(|\mathcal{A}| + (1-\gamma)^2 \sqrt{|\mathcal{A}|} \right)}{\tau(1-\gamma)^3 d}, \quad \ell_c = \frac{\sqrt{|\mathcal{A}|}}{(1-\gamma)^2}, \quad C_1 = \sqrt{\frac{2(1-\gamma) \ln 2}{\tau d}}, \quad C_2 = \frac{V_\tau^* - V_\tau^{\tilde{\pi}}(\rho)}{\xi}. \quad (123)$$

The total iteration complexity is $\log_2(\ell C_2/\varepsilon) \times N_2 + N_3 = \mathcal{O}(\log^2(1/\varepsilon) + \log(1/\varepsilon_1))$ with primal error bounds $\mathcal{O}(\sqrt{\varepsilon} + \varepsilon_1)$ given by (122b) - (122d) and a dual error bound $\mathcal{O}(\varepsilon)$ given by (122a).

We remark that the constants ℓ , ℓ_c , C_1 , and C_2 used in Theorem C.1 coincide with those used in previous sections, except that they correspond to the 1-dimensional situation. The maximum number of iterations N_1 and N_2 in Theorem C.1, respectively, correspond to the N_2 and N_3 in Theorem B.5.

Proof. Under the assumption that $\nabla D(0) < 0$ and $\nabla D(\xi) > 0$, the optimal dual variable λ^* must belong to the interval $(0, C_2)$ and $\nabla D(\lambda^*) = 0$. Denote $\tilde{\nabla}D(\lambda) = \text{Grad}_{Sub}(\lambda, \pi, \eta, N_1)$, i.e. the output of the gradient estimator. For a given threshold ε , when

$$N_1 = \frac{1}{1-\gamma} \log \left(\frac{4|\mathcal{A}|(1+C_2+\tau \log |\mathcal{A}|)}{(1-\gamma)^3 \tau \varepsilon} \right), \quad (124)$$

it follows from the proof of Theorem 5.2 (cf. (98)) that

$$|\tilde{\nabla}D(\lambda) - \nabla D(\lambda)| \leq \frac{\varepsilon}{2}. \quad (125)$$

Thus, if $\tilde{\nabla}D(\lambda) \geq \varepsilon$, we have that $\nabla D(\lambda) \geq \varepsilon/2$. Similarly, if $\tilde{\nabla}D(\lambda) \leq -\varepsilon$, we have that $\nabla D(\lambda) \leq -\varepsilon/2$. Therefore, the lines 7-11 in Algorithm 3 shrink by a factor of 2 the search region that contains the optimal solution λ^* .

By leveraging the triangular inequality, we have that

$$|\tilde{\nabla}D(\lambda)| \leq |\tilde{\nabla}D(\lambda) - \nabla D(\lambda)| + |\nabla D(\lambda) - \nabla D(\lambda^*)| \leq \frac{\varepsilon}{2} + \ell |\lambda - \lambda^*|. \quad (126)$$

where we apply the smoothness of the dual function (cf. Proposition 3.6). Thus, Algorithm 3 terminates in at most $t = \log_2(\ell C_2/\varepsilon)$ iterations with an ε -optimal stationary point due to

$$|\tilde{\nabla}D(\lambda)| \leq \frac{\varepsilon}{2} + \ell |\lambda - \lambda^*| \leq \frac{\varepsilon}{2} + \ell \left(\frac{1}{2} \right)^t C_2 \leq \varepsilon, \quad (127)$$

where λ denotes the midpoint $(p_t + q_t)/2$ generated in the t -th iteration in line 3.

Now, suppose that λ is the output solution with $|\tilde{\nabla}D(\lambda)| \leq \varepsilon$. It holds that

$$D(\lambda) - D_\tau^* \leq |\nabla D(\lambda)| \times |\lambda - \lambda^*| \leq \left(|\tilde{\nabla}D(\lambda)| + \frac{\varepsilon}{2} \right) C_2 = \frac{3C_2}{2} \varepsilon, \quad (128)$$

where we have used (120) and (127). By substituting (128) into the proof of Theorem 5.2 and letting $\varepsilon_0 = 3C_2\varepsilon/2$, the desired bounds follow. The convergence is linear and the total iteration complexity is upper-bounded by $\log_2(\ell C_2/\varepsilon) \times N_2 + N_3 = \mathcal{O}(\log^2(1/\varepsilon) + \log(1/\varepsilon_1))$. \square

Corollary C.2 (Restatement of Corollary 6.2) *Suppose that Assumptions 3.1 and 3.4 hold. Let*

$$\tau = \frac{(1-\gamma)\varepsilon}{4\log|\mathcal{A}|}. \quad (129)$$

Then, Algorithm 3 computes a solution π for the standard CMDP such that

$$\left| V^{\pi^*}(\rho) - V^\pi(\rho) \right| = \mathcal{O}(\varepsilon), \quad (130a)$$

$$\left[b - U_g^\pi(\rho) \right]_+ = \mathcal{O}(\varepsilon), \quad (130b)$$

in $\mathcal{O}(\log^2(1/\varepsilon))$ iterations, where π^ is an optimal policy to the standard CMDP.*

Proof. The proof can be fully adopted from that of Corollary B.6. The main difference lies in the outer-loop complexity. To have $\ell_c C_1 C_2 \sqrt{\varepsilon_0} = \varepsilon/2$, compared to the previous $\tilde{\mathcal{O}}(1/\varepsilon^2)$ total iterations, it only requires $\mathcal{O}(\log^2(1/\varepsilon))$ total iterations for Algorithm 3 when there is a single constraint. \square

D Supporting Lemmas

The followings are standard results about unregularized and entropy-regularized MDPs. We refer the reader to (Agarwal et al., 2021; Mei et al., 2020) for the proofs.

Lemma D.1 (Policy gradient for direct parameterization) *Suppose that $V^\pi(\rho)$ is an unregularized value function. For the direct policy parameterization where $\theta(s, a) = \pi_\theta(a|s)$, the gradient is*

$$\frac{\partial V^\pi(\rho)}{\partial \pi(a|s)} = \frac{1}{1-\gamma} d_\rho^\pi(s) Q^\pi(s, a). \quad (131)$$

Lemma D.2 (Policy gradient for soft-max parameterization) *Suppose that $V^\pi(\rho)$ is an unregularized value function. For the soft-max policy parameterization (cf. (6)), the gradient is*

$$\frac{\partial V^{\pi_\theta}(\rho)}{\partial \theta(s, a)} = \frac{1}{1-\gamma} d_\rho^{\pi_\theta}(s) \pi_\theta(a|s) A^{\pi_\theta}(s, a). \quad (132)$$

Lemma D.3 (Performance difference) *Suppose that $V^\pi(\rho)$ is an unregularized value function. For all policies π and π' , it holds that*

$$V^{\pi'}(\rho) - V^\pi(\rho) = \frac{1}{1-\gamma} \sum_{s \in \mathcal{S}} d_\rho^\pi(s) \sum_{a \in \mathcal{A}} (\pi'(a|s) - \pi(a|s)) \cdot Q^{\pi'}(s, a). \quad (133)$$

Lemma D.4 (Soft sub-optimality) *Suppose that $V_\tau^\pi(\rho)$ is an entropy-regularized value function and π_τ^* is the optimal policy. For every policy π , it holds that*

$$V_\tau^{\pi^*}(\rho) - V_\tau^\pi(\rho) = \frac{\tau}{1-\gamma} \sum_{s \in \mathcal{S}} d_\rho^\pi(s) D_{\text{KL}}[\pi(\cdot|s) | \pi_\tau^*(\cdot|s)], \quad (134)$$

where $D_{\text{KL}}[P(\cdot) | Q(\cdot)] := \sum_x P(x) (\log P(x) - \log Q(x))$ is the KL divergence between probability distributions $P(\cdot)$ and $Q(\cdot)$.