
Anchor-Changing Regularized Natural Policy Gradient for Multi-Objective Reinforcement Learning

Ruida Zhou*
Texas A&M University
ruida@tamu.edu

Tao Liu*
Texas A&M University
tliu@tamu.edu

Dileep Kalathil
Texas A&M University
dileep.kalathil@tamu.edu

P. R. Kumar
Texas A&M University
prk@tamu.edu

Chao Tian
Texas A&M University
chao.tian@tamu.edu

Abstract

We study policy optimization for Markov decision processes (MDPs) with multiple reward value functions, which are to be jointly optimized according to given criteria such as proportional fairness (smooth concave scalarization), hard constraints (constrained MDP), and max-min trade-off. We propose an Anchor-changing Regularized Natural Policy Gradient (ARNPG) framework, which can systematically incorporate ideas from well-performing first-order methods into the design of policy optimization algorithms for multi-objective MDP problems. Theoretically, the designed algorithms based on the ARNPG framework achieve $\tilde{O}(1/T)$ global convergence with exact gradients. Empirically, the ARNPG-guided algorithms also demonstrate superior performance compared to some existing policy gradient-based approaches in both exact gradients and sample-based scenarios.

1 Introduction

In many sequential decision-making scenarios, agents usually face multiple objectives simultaneously. This motivates the study of reinforcement learning (RL) with multiple reward values $V_{1:m}^\pi(\rho)$.² Given the achievable region $\mathcal{V} = \{V_{1:m}^\pi(\rho)\}_{\pi \in \Pi}$ consisting of value vectors achieved by policies in policy class Π , the agent employs certain criteria to reflect the system requirement. For example,

1. Proportional fairness [13]: Given $a_{1:m} > 0$, find $v \in \mathcal{V}$ that $\sum_{i=1}^m a_i \frac{v'_i - v_i}{v_i} \leq 0, \forall v' \in \mathcal{V}$.
2. Hard constraints [4]: Given $b_{2:m}$, maximize $v \in \mathcal{V}$ v_1 , subject to $v_i \geq b_i, \forall i = 2, \dots, m$.
3. Max-min trade-off [8]: Given $c_{1:m} > 0$, maximize $v \in \mathcal{V}$ $\min_{i \in [m]} (v_i/c_i)$.

We study policy gradient-based approaches that optimize over parameterized policies $\Pi = \{\pi_\theta : \theta \in \Theta\}$ through policy gradient. In general, the optimization problems above may not be convex in terms of θ , not even for single-objective MDPs with direct parameterization by $\theta_{s,a} = \pi_\theta(a|s)$ [2]. Due to the non-convexity, $O(1/T)$ global convergence of policy gradient-based methods was only established very recently for single-objective MDPs with exact gradients [2, 20]. These breakthrough results have motivated the study of policy optimization for multi-objective MDPs, e.g., smooth concave scalarization [5], constrained MDPs (CMDPs) [11, 31].

However, under the exact gradients scenario, the previous approaches for multi-objective MDPs, either suffer from slow provable $O(1/\sqrt{T})$ global convergence [11], or require extra assumptions

*The first two authors contributed equally.

²The notations are formally defined in Section 2.

[37, 33, 18]. The compactness of Θ is assumed in [37], but this assumption forbids a very common softmax parameterization, where $\Theta = \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$. The NPG-based methods have been analyzed in [33, 18] under an ergodicity assumption, but such an assumption is not required for NPG in single-objective MDPs [2], and therefore appears artificial.

The above criteria for multi-objective MDPs could be viewed as convex optimization problems w.r.t. a value vector $v \in \mathcal{V}$, for which there are a wide array of well-performing first-order methods for convex optimization problems in general. It is desirable to take full advantage of such efficient first-order methods in a unified and flexible manner when designing policy gradient-based algorithms for multi-objective MDPs.

Main contributions

1. We propose an anchor-changing regularized natural policy gradient (ARNPG) framework in Section 3 that can exploit and integrate first-order methods for the design of policy gradient-based algorithms for multi-objective MDPs.
2. We demonstrate the strength of the ARNPG framework by designing algorithms for three general criteria: smooth concave scalarization (Section 4.1), constrained MDPs (Section 4.2), and max-min trade-off (Section 4.3).
3. Under softmax parameterization with exact gradients, the proposed algorithms inherit the advantages of the integrated first-order methods, and are guaranteed to have $\tilde{O}(1/T)$ global convergence without further assumptions on the underlying MDP.
4. In addition to the theoretical advantages, we provide the results of extensive experimentation in Section 5 and Appendices A and B which demonstrate that the ARNPG-guided algorithms provide superior performance in exact gradient and sample-based tabular scenarios, as well as actor-critic deep RL scenarios, compared to several existing policy gradient-based approaches.

1.1 Related works

Policy gradient (PG)-based methods have drawn much attention recently [1, 20, 10, 14] due to their simplicity as well as the potential to generalize to large scale problems. Despite their non-convex nature, PG-based methods have been shown to converge globally for single-objective MDPs [1, 20]. Their convergence may be further accelerated with appropriate regularization [10, 17], e.g., entropy regularization, but the algorithms only converge to the optimum of the regularized problem instead of the desired (unregularized) problem.

This paper considers *single-policy* multi-objective MDPs, including CMDPs where constraints are specified on some objectives. Global convergence of PG-based approaches in the multi-objective MDPs has been previously studied. For smooth concave scalarization, Bai et al. [5] showed an $O(1/\epsilon^4)$ sample complexity (to achieve ϵ -optimal in expectation) of the policy-gradient method under sample-based scenarios. However, with exact gradients, we are unaware of works with fast $\tilde{O}(1/T)$ convergence. For CMDPs, Ding et al. [11] have studied a primal-dual NPG algorithm achieving $O(1/\sqrt{T})$ global convergence for both the optimality gap and the constraint violation. Xu et al. [31] have proposed a primal approach that reduces constraint violations with a higher priority than optimizing objective, and enjoys the same $O(1/\sqrt{T})$ global convergence. In work conducted concurrently with ours, [33] and [18] have proposed algorithms that achieve $\tilde{O}(1/T)$ convergence but with extra ergodicity assumptions.

A general setting of optimizing a concave function of the state-action visitation distribution has been considered in [37]. Though the problem is more general, its gradient estimation is more complicated than the canonical policy gradient estimate. Zhang et al. [37] showed that the gradient ascent achieves $O(1/T)$ global convergence for smooth scalarization with exact gradients, under several assumptions such as convexity and compactness of the parameter set Θ . Directly viewing the state-action visitation as the decision variables and imposing equality constraints for their feasibility, a smooth concave scalarization has been studied in [36] and later generalized to the constrained setting in [6]. These two works focus on sample-based scenarios, but due to their primal-dual approach with equality constraints, the convergence rate is only $O(1/\sqrt{T})$ even with exact gradients. Moreover, the state-action visitation parameterization is difficult to generalize to larger scale deep RL scenarios.

A more thorough discussion on related works is given in Appendix F.

2 Preliminaries

System model A Markov decision process (MDP) is represented by a tuple $(\mathcal{S}, \mathcal{A}, P, \rho, \gamma, r)$, where \mathcal{S} is the state space, \mathcal{A} the action space, $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ the transition kernel, $\rho \in \Delta(\mathcal{S})$ the initial state distribution, $\gamma \in (0, 1)$ the discount factor, and $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ the reward function. Given any policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ and any reward function $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$, we define the state value function $V_r^\pi : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$, and the state-action value function $Q_r^\pi : \mathcal{S} \times \mathcal{A} \rightarrow [0, \frac{1}{1-\gamma}]$, as

$$V_r^\pi(s) := \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, \pi], \quad Q_r^\pi(s, a) := \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a, \pi],$$

where expectation \mathbb{E} is taken over the random trajectory of the Markov chain induced by the policy π and the transition kernel P . With a slight abuse of notation, we denote $V_r^\pi(\rho) := \mathbb{E}_{s \sim \rho}[V_r^\pi(s)]$. Define the discounted state-action visitation distribution (state-action visitation for short) of policy π with initial state distribution ρ by $d_\rho^\pi(s, a) := (1 - \gamma) \mathbb{E}_{s_0 \sim \rho}[\sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s, a_t = a \mid s_0, \pi)]$. It then follows that $V_r^\pi(\rho) = \frac{1}{1-\gamma} \langle d_\rho^\pi, r \rangle$ by viewing d_ρ^π and r as $|\mathcal{S}| |\mathcal{A}|$ -dimensional vectors indexed by $(s, a) \in \mathcal{S} \times \mathcal{A}$. When it is clear from the context, we denote the state visitation distribution by $d_\rho^\pi(s) := \mathbb{E}_{s_0 \sim \rho}[(1 - \gamma) \sum_{t=0}^{\infty} \gamma^t \mathbb{P}(s_t = s \mid s_0)]$, which is the marginal distribution of the state-action visitation $d_\rho^\pi(s, a)$, i.e., $d_\rho^\pi(s) = \sum_{a \in \mathcal{A}} d_\rho^\pi(s, a)$.

We study an MDP with m objectives represented by $(\mathcal{S}, \mathcal{A}, P, \rho, \gamma, r_{1:m})$, where $r_i : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ is the i -th reward function for each $i \in [m]$. For simplicity, denote $V_i^\pi(\cdot) := V_{r_i}^\pi(\cdot)$ and $V_{1:m}^\pi(\cdot) := (V_1^\pi(\cdot), \dots, V_m^\pi(\cdot))$. We consider parameterized policies in $\Pi = \{\pi_\theta : \theta \in \Theta\}$, where $\Theta \subset \mathbb{R}^n$ is the parameter space. For example, the softmax policy is $\pi_\theta(a|s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$ with $\Theta = \mathbb{R}^{|\mathcal{S}| |\mathcal{A}|}$; and neural softmax policy is $\pi_\theta(a|s) = \frac{\exp(\text{NN}_\theta(s, a))}{\sum_{a'} \exp(\text{NN}_\theta(s, a'))}$, where NN_θ is some neural network parameterized θ . Define $\mathcal{V} := \{V_{1:m}^{\pi_\theta}(\rho) : \theta \in \Theta\}$ as the achievable region of value vectors. The agent wishes to optimize the policy in Π for a given specific multi-objective criterion on value vectors in \mathcal{V} .

Mirror ascent As one of the most well-known iterative optimization methods, mirror descent (actually ascent in the context of our formulation as a maximization problem) [21, 7] is a general class that encompasses many first-order methods in convex optimization. Given a variable x in a compact convex set $\mathcal{X} \subset \mathbb{R}^n$ and an ascent direction $g \in \mathbb{R}^n$, the variational representation of the mirror ascent update is

$$x' \in \arg \max_{y \in \mathcal{X}} \{ \langle g, y \rangle - \alpha B_h(y \| x) \}, \quad (1)$$

where $B_h(x \| y) := h(x) - h(y) - \langle \nabla h(y), x - y \rangle$ is some Bregman divergence generated by a differentiable convex function $h : \mathcal{X} \rightarrow \mathbb{R}$. When analyzing the convergence of first-order methods, certain fundamental inequalities are usually established to facilitate the proof. One such inequality is

$$\langle g, x' \rangle - \alpha B_h(x' \| x) \geq \langle g, y \rangle - \alpha B_h(y \| x) + \alpha B_h(y \| x'), \quad \forall y \in \mathcal{X}, \quad (2)$$

which is a critical step in many previous works, e.g., [22, 27, 16].

It is desirable to construct a similar fundamental inequality for multi-objective MDPs that can facilitate the analysis of convergence. As we will show in the next section, such an inequality can indeed be established in a new framework, which we refer to as the Anchor-Changing Regularized Natural Policy Gradient (ARNPG).

Notations Denote KL-divergence between two n -dimensional probability vectors x, y by $D(x \| y) := \sum_{i=1}^n x_i \log(x_i / y_i)$, which is a widely-used Bregman divergence. For any policies π, π' and state visitation distribution d , define $D_d(\pi \| \pi') := \sum_{s \in \mathcal{S}} d(s) D(\pi(\cdot | s) \| \pi'(\cdot | s))$. A *uniform policy* is one which chooses actions uniformly at random.

3 Anchor-changing regularized natural policy gradient

Let us consider a hypothetical mirror ascent update on decision value vector $v_k \in \mathcal{V}$ according to (1). Given an ascent direction \tilde{G}_k along which to improve v_k , the updated value vector is

$$v' \in \arg \max_{v \in \mathcal{V}} \{ \langle \tilde{G}_k, v \rangle - \alpha B_h(v \| v_k) \}. \quad (3)$$

Suppose the value vector v_k is achieved by a policy π_{θ_k} , i.e., $v_k = V_{1:m}^{\pi_{\theta_k}}(\rho)$. Denote the reward function in the ascent direction as $\tilde{r}_k(s, a) = \langle \tilde{G}_k, r_{1:m}(s, a) \rangle$. It follows that $\langle \tilde{G}_k, v_k \rangle = V_{\tilde{r}_k}^{\pi_{\theta_k}}(\rho)$. Note that $B_h(v||v_k)$ in (3) serves the role of a soft constraint on v by keeping v within a vicinity of v_k . Replacing $B(v||v_k)$ by $\frac{D_{d_\rho}^{\pi_\theta}(\pi_\theta||\pi_{\theta_k})}{1-\gamma}$ will induce a similar soft constraint that prefers the vicinity of the ‘‘anchor’’ policy π_{θ_k} . Therefore we consider replacing the variational update in (3) by

$$\theta' \in \arg \max_{\theta \in \Theta} \left\{ \tilde{V}_{k,\alpha}^{\pi_\theta}(\rho) \right\}, \quad \text{where} \quad \tilde{V}_{k,\alpha}^{\pi_\theta}(\rho) := V_{\tilde{r}_k}^{\pi_\theta}(\rho) - \alpha \frac{D_{d_\rho}^{\pi_\theta}(\pi_\theta||\pi_{\theta_k})}{1-\gamma}. \quad (4)$$

ARNPG Motivated by the intuition above, we propose the Anchor-Changing Regularized Natural Policy Gradient (ARNPG) framework. At (macro) step k , the ARNPG framework determines the reward function in the ascent direction \tilde{r}_k and the anchor policy π_{θ_k} , which can exploit well-performed first-order methods in convex optimization literature utilizing the features of the specific criteria in use. With \tilde{r}_k and π_{θ_k} , we wish to solve for (4) to improve the value vector. However the optimal solution θ' of (4) is generally not determinable explicitly. ARNPG therefore approaches the optimal solution via a subroutine that executes a natural policy gradient (NPG) algorithm w.r.t. the KL-regularized value function $\tilde{V}_{k,\alpha}^{\pi_\theta}(\rho)$. We refer to this subroutine, given in Algorithm 1, as `InnerLoop`($\tilde{r}_k, \pi_{\theta_k}, \alpha, \eta, t_k$). It iteratively updates the parameter $\theta_k^{(t)}$ for t_k (micro) steps according to the NPG update rule as in (5), where $\mathcal{F}_\rho(\theta)^\dagger$ is the Moore-Penrose inverse of the Fisher information matrix $\mathcal{F}_\rho(\theta) := \mathbb{E}_{(s,a) \sim d_\rho^{\pi_\theta}} \left[\nabla_\theta \log \pi_\theta(a|s) (\nabla_\theta \log \pi_\theta(a|s))^\top \right]$.

Algorithm 1: `InnerLoop`($\tilde{r}_k, \pi_{\theta_k}, \alpha, \eta, t_k$)

Initialize $\theta_k^{(0)} = \theta_k$
for $t = 0, 1, \dots, t_k - 1$ **do**
 $\theta_k^{(t+1)} \leftarrow \theta_k^{(t)} + \eta \mathcal{F}_\rho(\theta_k^{(t)})^\dagger \nabla_\theta \tilde{V}_{k,\alpha}^{\pi_{\theta_k^{(t)}}}(\rho)$
Return $\theta_k^{(t_k)}$

The choice of the number of iterations in `InnerLoop` (i.e., t_k) involves a trade-off between the variational update precision and the overall efficiency. On the one hand, a larger t_k leads to a more accurate approximation of the optimal solution θ' to (4), but it may cause the algorithm to spend unnecessary computational resources on the regularized objective $\tilde{V}_{k,\alpha}^{\pi_\theta}(\rho)$, instead of on the true optimization problem. On the other hand, a smaller t_k saves inner loop iterations but the update follows less closely to the underlying mirror-ascent update in improving the value vector. In our experiments, we choose t_k within 10 to strike a balance and empirically observe $t_k > 1$ has better performance.

We note that when $t_k = 1$, the gradient $\nabla_\theta \tilde{V}_{k,\alpha}^{\pi_{\theta_k}}(\rho) = \nabla_\theta V_{\tilde{r}_k}^{\pi_{\theta_k}}(\rho)$, since $D_{d_\rho}^{\pi_\theta}(\pi_\theta||\pi_{\theta_k})$ has zero gradient at $\theta = \theta_k$. The update in (5) reduces to an NPG update on the unregularized value function $V_{\tilde{r}_k}^{\pi_\theta}(\rho)$. For single-objective MDPs, it reduces to the canonical NPG method.

3.1 Theoretical guarantee of ARNPG

We now present the main theoretical tool for the analysis of the ARNPG framework. Recall the discussion of the fundamental inequality after (2). Proposition 1 establishes such a fundamental inequality with controllable approximation error under the softmax policy parameterization, i.e., $\pi_\theta(a|s) = \frac{\exp(\theta_{s,a})}{\sum_{a'} \exp(\theta_{s,a'})}$. In the rest of the paper, we omit θ in π_θ when it is clear from the context, but it should be noted that all updates of policies are performed on the parameters.

Proposition 1. *Under the softmax parameterization, given $\epsilon_k > 0$, for any $\tilde{r}_k, t_k \geq \frac{1}{1-\gamma} \log(\frac{5||\tilde{r}_k||_\infty}{(1-\gamma)^2 \epsilon_k}) + 1$, $\alpha > 0$ and $\eta = \frac{1-\gamma}{\alpha}$, the update $\pi_{k+1} \leftarrow \text{InnerLoop}(\pi_k, \tilde{r}_k, \alpha, \eta, t_k)$ satisfies*

$$V_{\tilde{r}_k}^{\pi_{k+1}}(\rho) - \alpha \frac{D_{d_\rho}^{\pi_{k+1}}(\pi_{k+1}||\pi_k)}{1-\gamma} \geq V_{\tilde{r}_k}^{\pi_k}(\rho) - \alpha \frac{D_{d_\rho}^{\pi_k}(\pi_k||\pi_k) - D_{d_\rho}^{\pi_k}(\pi_k||\pi_{k+1})}{1-\gamma} - \epsilon_k, \quad \forall \pi. \quad (6)$$

The inequality (6) is critical to the convergence proof. Its right hand side allows telescoping, which by summing over k can iteratively cancel the terms $D_{d_\rho^\pi}(\pi||\pi_k)$. Since $t_k = \Theta(\log(1/\epsilon_k))$ it suffices to use very few iterations in InnerLoop for maintaining precision.

Remark. It has been shown that for the entropy-regularized MDP, i.e., KL-regularized with the uniform policy as the anchor policy, NPG converges linearly (i.e., geometrically fast) to the regularized optimal policy [10]. It is natural to anticipate that for the KL-regularized MDP $\tilde{V}_{k,\alpha}^{\pi_k}(\rho)$ with anchor π_k , NPG would similarly converge linearly (i.e., $\tilde{V}_{k,\alpha}^{\pi_k} \geq \tilde{V}_{k,\alpha}^{\pi_k^*} - \epsilon$ for $t_k = \Theta(\log(1/\epsilon))$) to a corresponding optimal policy, denoted as π_k^* . In contrast, the right hand side of inequality (6) has a *positive drift* $\alpha \frac{D_{d_\rho^\pi}(\pi||\pi_{k+1})}{1-\gamma}$ for any policy π , which is considerably stronger.

Proof sketch of Proposition 1. We can show that InnerLoop approximately solves the variational update in (4) with linear convergence as anticipated. However to establish (6), the difficulty lies in the introduction of positive drift, since $V_{\tilde{r}_k}^{\pi_\theta}(\rho)$ is not concave w.r.t. θ and $D_{d_\rho^{\pi_\theta}}(\pi_\theta||\pi_{\theta_k})$ may not be a Bregman divergence. We tackle this difficulty by showing that optimizing π_θ in InnerLoop implicitly performs a mirror ascent update for state action visitation $d_\rho^{\pi_\theta}$. \square

As demonstrated in the next section, Proposition 1 ensures that the convergence rate of the algorithms derived from the ARNPG framework is of the same rate as the underlying first-order methods with only extra logarithmic factors.

4 Theoretical applications

In this section, we apply the ARNPG framework to several important multi-objective MDP scenarios and obtain new policy optimization algorithms by integrating first-order methods in convex optimization. All the theoretical results presented in this section are under the softmax parameterization with exact gradients. However, the obtained algorithms can be implemented in more general settings such as neural softmax and sample-based scenarios, as in the next section. We theoretically establish $\tilde{O}(1/T)$ convergence of these algorithms by leveraging the fundamental inequality in Proposition 1.

4.1 Smooth concave scalarization function

We start by considering the following optimization problem

$$\max_{\theta} F(V_{1:m}^{\pi_\theta}(\rho)), \quad (7)$$

where F is a concave function, and β -smooth w.r.t. $\|\cdot\|_\infty$ norm, i.e., $\|\nabla F(v) - \nabla F(v')\|_1 \leq \beta\|v - v'\|_\infty$. Since the set of achievable values $\mathcal{V} \subseteq \left[0, \frac{1}{1-\gamma}\right]^m$, it can be verified that $\|\nabla F(v)\|_1 \leq L$ for some factor $L > 0$.

The proportional fair criterion discussed in Section 1 can be approximated by $F(v) := \sum_{i=1}^m a_i \log(\delta + v_i)$, where $\delta > 0$ is some constant introduced to circumvent the pathological case $v_i = 0$ for some $i \in [m]$. Under this criterion, $\beta = \sum_{i=1}^m a_i/\delta^2$ and $L = \sum_{i=1}^m a_i/\delta$.

When v is viewed as the decision variable, at macro step k with value vector $V_{1:m}^{\pi_k}(\rho)$, the ascent direction in a typical gradient ascent step is the gradient $\tilde{G}_k = \nabla_v F(V_{1:m}^{\pi_k}(\rho))$. This naturally determines the reward in the ascent direction as $\tilde{r}_k(s, a) = \langle \tilde{G}_k, r_{1:m}(s, a) \rangle$. Adapting the ARNPG framework to this specific context, we present the algorithm for solving the program (7) in Algorithm 2. We refer to it as “implicit mirror descent” because the algorithm implicitly employs mirror descent.

Algorithm 2: ARNPG Implicit Mirror Descent (ARNPG-IMD)

Input $\pi_0, \alpha, \eta, t_{0:K-1}, K$
for $k = 0, 1, \dots, K-1$ **do**
 Update $\pi_{k+1} \leftarrow \text{InnerLoop}(\pi_k, \tilde{r}_k, \alpha, \eta, t_k)$
Return the policy in $\{\pi_k\}_{k=1}^K$ with the largest $F(V_{1:m}^{\pi_k}(\rho))$

Let π^* be the optimal policy for (7). Based on Proposition 1, we present the following theorem which guarantees the convergence of ARNPG-IMD with appropriately selected parameters π_0, α, η, t_k .

Theorem 1. For any $K \geq 1$, take uniform policy π_0 , $\alpha \geq \frac{\beta}{(1-\gamma)^3}$, $\eta = \frac{1-\gamma}{\alpha}$, and $t_k = \lceil \frac{1}{1-\gamma} \log(\frac{5LK}{\beta \log(|\mathcal{A}|)}) + 1 \rceil$. The optimality gap of ARNPG-IMD (Algorithm 2) satisfies

$$F(V_{1:m}^{\pi^*}(\rho)) - \max_{k \in [1:K]} F(V_{1:m}^{\pi_k}(\rho)) \leq F(V_{1:m}^{\pi^*}(\rho)) - \frac{1}{K} \sum_{k=1}^K F(V_{1:m}^{\pi_k}(\rho)) \leq \frac{2\alpha \log(|\mathcal{A}|)}{(1-\gamma)K}. \quad (8)$$

There are a total of K macro steps, and the total number of iterations is $T = \sum_{k=0}^{K-1} t_k = \Theta(\frac{K}{1-\gamma} \log(K))$. The following corollary provides the convergence rate in terms of T .

Corollary 1. Under the same conditions as in Theorem 1, the ARNPG-IMD algorithm satisfies $F(V_{1:m}^{\pi^*}(\rho)) - \frac{1}{K} \sum_{k=1}^K F(V_{1:m}^{\pi_k}(\rho)) = O\left(\frac{\beta \log(T)}{(1-\gamma)^5 T}\right)$.

Remark. In the absence of knowledge of K , we can select time-varying numbers of InnerLoop iterations, such as $t_k = \Theta(\log(k))$, and ARNPG-IMD will still have the same $\tilde{O}(1/T)$ convergence.

4.2 Constrained Markov decision process

Another way of trading off the objectives is to optimize one while setting hard constraints on the others. This can be formulated as the following constrained MDP (CMDP) problem:

$$\max_{\theta} V_1^{\pi_{\theta}}(\rho), \quad \text{s.t. } V_i^{\pi_{\theta}}(\rho) \geq b_i, \quad \forall i \in [2:m], \quad (9)$$

where $b_{2:m} \in [0, \frac{1}{1-\gamma}]^{m-1}$. Let $\pi^* = \pi_{\theta^*}$ be the optimal policy of the CMDP problem in (9).

Define the Lagrangian of the CMDP problem as $\mathcal{L}(\pi_{\theta}, \lambda) = V_1^{\pi_{\theta}}(\rho) + \sum_{i=2}^m \lambda_i (V_i^{\pi_{\theta}}(\rho) - b_i)$, where λ_i is the Lagrange multiplier (dual variable) corresponding to the constraint $V_i^{\pi_{\theta}} \geq b_i$, for each $i \in [2:m]$. The Lagrange dual function $\max_{\pi} \mathcal{L}(\pi, \cdot)$ is a convex function of dual variables $\lambda \geq 0$. Denote by λ^* the optimal dual variables that minimize the Lagrange dual function. We assume λ^* is finite, which is guaranteed by Slater's condition, i.e., there is some π_{θ} and $\xi > 0$ with $V_i^{\pi_{\theta}}(\rho) - b_i \geq \xi$ for any $i \in [2:m]$. Note (π^*, λ^*) is a saddle point of the Lagrangian $\mathcal{L}(\pi, \lambda)$. This motivates the primal-dual approach, which iteratively performs gradient ascent for π_{θ} and gradient descent for λ . This is suitable for the CMDP setting, since for any fixed λ , the Lagrangian $\mathcal{L}(\pi, \lambda)$ corresponds to an MDP for which policy gradient can be employed.

The canonical primal-dual gradient ascent-descent method for constrained convex optimization can only guarantee $O(1/\sqrt{T})$ convergence, and consequently the primal-dual policy gradient-based approach for CMDPs [11] has the same convergence. Recently, Yu et al. [35] have proposed a primal-dual-based method with $O(1/T)$ convergence under the Euclidean setting, i.e., $B_h(x||y) = \frac{1}{2} \|x-y\|_2^2$. Adopting ideas from [35], we next propose the ARNPG with Extra Primal-Dual (ARNPG-EPD) algorithm (Algorithm 3). To the best of our knowledge, this new primal-dual update appears in the CMDP-related literature for the first time.

Note that $b_i - V_i^{\pi}(\rho)$ is the amount of constraint violation. There are two key ideas we adopt from [35]. The first is the design of the reward in the ascent direction

$$\tilde{r}_k(s, a) := r_1(s, a) + \sum_{i=2}^m (\lambda_{k,i} + \eta' (b_i - V_i^{\pi_k}(\rho))) r_i(s, a),$$

where an extra constraint violation term is added to the dual variables. The second idea is that the update of dual variables should not fall below the negative constraint violation (the first term in (10)), and it can alleviate the overshooting of dual variables. The extra constraint violation terms in \tilde{r}_k and the dual update work jointly to ensure the $\tilde{O}(1/T)$ convergence.

Theorem 2. For any $K \geq 1$ and $\eta' \in (0, 1]$, take uniform policy π_0 , $\alpha \geq \frac{2\eta' m}{(1-\gamma)^3}$, $\eta = \frac{1-\gamma}{\alpha}$, and choose $t_k = \lceil \frac{1}{1-\gamma} \log(\frac{5L_k K}{2\eta' m \log(|\mathcal{A}|)}) + 1 \rceil$ with $L_k = 1 + \frac{\eta' (m-1)}{1-\gamma} + \sum_{i=2}^m \lambda_{k,i}$. The average optimality gap and the average constraint violation of ARNPG-EPD (Algorithm 3) satisfy

$$V_1^{\pi^*}(\rho) - \frac{1}{K} \sum_{k=1}^K V_1^{\pi_k}(\rho) \leq \frac{3\alpha \log(|\mathcal{A}|)}{(1-\gamma)K}, \quad (11)$$

$$b_i - \frac{1}{K} \sum_{k=1}^K V_i^{\pi_k}(\rho) \leq \frac{1}{K} \left(\frac{2\|\lambda^*\|_2}{\eta'} + 3\sqrt{\frac{\alpha \log(|\mathcal{A}|)}{(1-\gamma)\eta'}} \right) \quad \forall i \in [2:m]. \quad (12)$$

Algorithm 3: ARNPG with Extra Primal Dual (ARNPG-EPD)

Input $\pi_0, \eta', \alpha, \eta, t_{0:K-1}, K$ **Initialize** $\lambda_{0,i} = \max\{\eta'(V_i^{\pi_0}(\rho) - b_i), 0\}, \forall i \in [2 : m]$ **for** $k = 0, 1, \dots, K-1$ **do** Update $\pi_{k+1} \leftarrow \text{InnerLoop}(\pi_k, \tilde{r}_k, \alpha, \eta, t_k)$ Update $\lambda_{k+1,i} = \max\{\eta'(V_i^{\pi_{k+1}}(\rho) - b_i), \lambda_{k,i} + \eta'(b_i - V_i^{\pi_{k+1}}(\rho))\}, \forall i \in [2 : m]$ (10)**Return:** a policy randomly chosen from $\{\pi_k\}_{k=1}^K$

Note that the number of micro steps t_k is chosen according to the dual variables λ_k in the previous theorem. Denote by $T := \sum_{k=0}^{K-1} t_k$ the total number of iterations.

Corollary 2. *Under the same conditions as in Theorem 2, the ARNPG-EPD algorithm satisfies $V_1^{\pi^*}(\rho) - \frac{1}{K} \sum_{k=1}^K V_1^{\pi_k}(\rho) = O(\frac{m \log(T)}{(1-\gamma)^{5T}})$, and $b_i - \frac{1}{K} \sum_{k=1}^K V_i^{\pi_k}(\rho) = O(\frac{\sqrt{m} \log(T)}{(1-\gamma)^{2.5T}})$.*

The theorem and corollary establish convergence of the average optimality gap and the average constraint violation, in the same manner as many previous works [11, 31, 12, 19] on CMDPs. However, a guarantee on the last iterate is more preferable. This drawback is inherited from the primal-dual algorithm for convex optimization, where the primal-dual algorithm with sublinear convergence can only be guaranteed on the average solution, as of our knowledge. Last iterate convergence is still an on-going open research topic.

4.3 Max-min trade-off criteria

Finally, we consider the max-min trade-off criterion defined as

$$\max_{\theta} \min_{\lambda \in \Lambda} \Phi(V_{1:m}^{\pi_{\theta}}(\rho), \lambda), \quad (13)$$

where Λ is a subset of the m -dimensional probability simplex $\Delta([m])$. We assume $\Phi(\cdot, \lambda)$ is concave and $\Phi(v, \cdot)$ is convex. We also assume Φ is β -smooth w.r.t. the norm $\Psi(v, \lambda) = \|v\|_{\infty} + \|\lambda\|_1$.

The max-min criterion mentioned in Section 1 can be represented by $\Phi(v, \lambda) = \sum_{i=1}^m v_i \lambda_i / c_i$ and $\Lambda = \Delta([m])$. Φ satisfies the concave-convex assumption and is β -smooth w.r.t. the norm Ψ with $\beta = O(m)$.

Denote $F(v) := \min_{\lambda \in \Lambda} \Phi(v, \lambda)$, which is concave but not necessarily smooth. Thus we cannot apply the ARNPG-IMD algorithm (Algorithm 2) due to the non-smoothness of F , and the subgradient-based method can only guarantee $O(1/\sqrt{T})$ convergence.

We next integrate the optimistic mirror descent ascent (OMDA) method [27] for solving minimax optimization in the ARNPG framework. Denote the gradients $\tilde{G}_k^{\lambda} = \nabla_{\lambda} \Phi(V_{1:m}^{\pi_k}(\rho), \tilde{\lambda}_k)$ and $\tilde{G}_k^v = \nabla_v \Phi(V_{1:m}^{\pi_k}(\rho), \tilde{\lambda}_k)$. It can be verified that $\|\tilde{G}_k^v\|_1 \leq L$ for some L due to the smoothness of Φ . OMDA performs gradient ascent along the direction \tilde{G}_k^v w.r.t. the value vector, and therefore we construct the reward in the ascent direction as $\tilde{r}_k(s, a) = \langle \tilde{G}_k^v, r_{1:m}(s, a) \rangle$. OMDA performs mirror descent along direction \tilde{G}_k^{λ} w.r.t. the dual vector λ . A key ingredient of OMDA is that it updates twice in each macro step. ARNPG-OMD adopts this idea and update (π, λ) from the same anchor points (π_k, λ_k) , first with ascent direction $(\tilde{r}_k, -\tilde{G}_k^{\lambda}) \in \mathbb{R}^{2m}$ and then a further step with direction $(\tilde{r}_{k+1}, -\tilde{G}_{k+1}^{\lambda}) \in \mathbb{R}^{2m}$.

We present ARNPG-OMDA in Algorithm 4, and establish the following performance guarantees:

Theorem 3. *For any $K \geq 1$, take uniform policy $\pi_0, \eta' \leq \frac{1}{6\beta}, \alpha \geq \frac{6\beta}{(1-\gamma)^3}, \eta = \frac{1-\gamma}{\alpha}$, and $t_k = \lceil \frac{1}{1-\gamma} \log(\frac{5LK}{6\beta \log(|\mathcal{A}|)}) + 1 \rceil$. The ARNPG-OMDA algorithm (Algorithm 4) satisfies*

$$F(V_{1:m}^{\pi^*}(\rho)) - F\left(\frac{1}{K} \sum_{k=1}^K V_{1:m}^{\pi_k}(\rho)\right) \leq \frac{3\alpha \log(|\mathcal{A}|)}{(1-\gamma)K} + \frac{\log(m)}{\eta' K}. \quad (14)$$

Algorithm 4: ARNPG with Optimistic Mirror Descent Ascent (ARNPG-OMDA)

Input $\pi_0, \lambda_0, \eta', \alpha, \eta, t_{0:K-1}, K$

Initialize $\tilde{\pi}_0 = \pi_0$ and $\lambda_0, \tilde{\lambda}_0$ as uniform distribution on $[m]$

for $k = 0, 1, \dots, K-1$ **do**

 Update $\tilde{\pi}_{k+1} \leftarrow \text{InnerLoop}(\pi_k, \tilde{r}_k, \alpha, \eta, t_k), \tilde{\lambda}_{k+1} \leftarrow \arg \min_{\lambda \in \Lambda} \{ \langle \tilde{G}_k^\lambda, \lambda \rangle + \frac{D(\lambda || \lambda_k)}{\eta'} \}$

 Update $\pi_{k+1} \leftarrow \text{InnerLoop}(\pi_k, \tilde{r}_{k+1}, \alpha, \eta, t_k), \lambda_{k+1} \leftarrow \arg \min_{\lambda \in \Lambda} \{ \langle \tilde{G}_{k+1}^\lambda, \lambda \rangle + \frac{D(\lambda || \lambda_k)}{\eta'} \}$

Return: a policy randomly chosen from $\{\tilde{\pi}_k\}_{k=1}^K$

Similar to the discussion after Corollary 2, Theorem 3 provides a performance guarantee on the average value vector $F(\frac{1}{K} \sum_{k=1}^K V_{1:m}^{\tilde{\pi}_k}(\rho))$, which is inherited from the OMDA methods. Denote the total number of iterations by $T := \sum_{k=0}^{K-1} 2t_k$.

Corollary 3. *Under the same conditions as in Theorem 3, ARNPG-OMDA satisfies $F(V_{1:m}^{\pi^*}(\rho)) - F(\frac{1}{K} \sum_{k=1}^K V_{1:m}^{\tilde{\pi}_k}(\rho)) = O\left(\frac{\beta \log(T)}{(1-\gamma)^5 T}\right)$.*

5 Empirical evaluation and application

In this section, we present the experimental results on CMDP. We compare the performance of the proposed ARNPG-EPD algorithm (Algorithm 3) with two benchmarks: NPG-PD [11] and CRPO [31]. Experimental details on CMDP are postponed to Appendix A and further experiments on smooth concave scalarization and max-min trade-off are presented in Appendix B. We provide code at <https://github.com/tliu1997/ARNPG-MORL>.

5.1 Tabular CMDP with exact gradients

Recall that under softmax policy with exact gradients, Corollary 2 (Theorem 2) guarantees $\tilde{O}(1/T)$ convergence of both performance measures: average optimality gap and average constraint violation. We compare the proposed ARNPG-EPD with the benchmarks NPG-PD and CRPO under both performance measures on a randomly generated CMDP with a single constraint, which are illustrated in Figure 1. The horizontal axis is the total number of iterations, i.e., including the micro steps in InnerLoop of ARNPG-EPD.

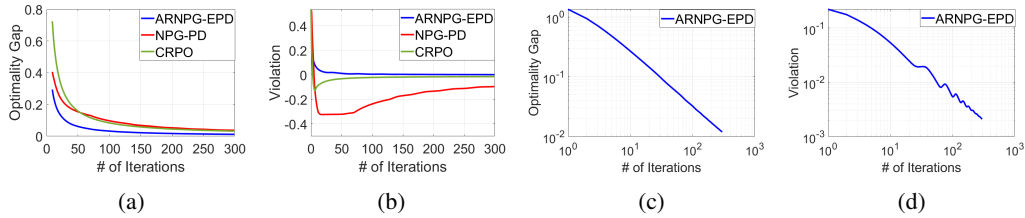


Figure 1: The average optimality gap and the average constraint violation versus the total number of iterations, for ARNPG-EPD, NPG-PD, and CRPO on a randomly generated CMDP.

Figures 1(a) and 1(b) show that both the average optimality gap and the average constraint violation of the ARNPG-EPD algorithm converge faster than those of NPG-PD. Since the CRPO focuses on the violated constraint, the policy becomes feasible quickly, though at the cost of an initially slower convergence for the optimality gap. As illustrated in Figures 1(c) and 1(d), the slopes of both the optimality gap and the constraint violation of the ARNPG-EPD algorithm in the log-log plots are approximately between -0.9 and -1, indicating a converge rate of $\tilde{O}(1/T)$.

5.2 Sample-based tabular CMDP

We next consider the same tabular CMDP described in Section 5.1 without exact policy gradients. Instead, policy gradients are estimated by samples from a generative model that can generate

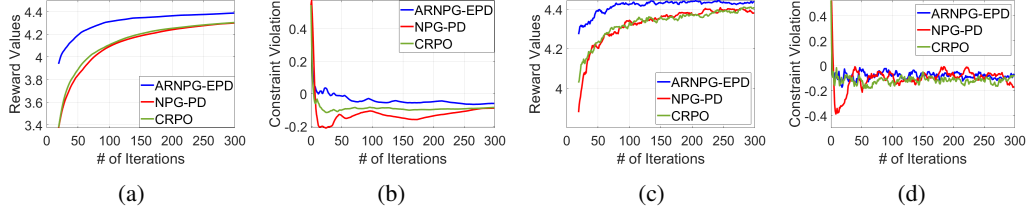


Figure 2: The reward values and the constraint violation with respect to the total number of iterations, for sample-based ARNPG-EPD, NPG-PD, and CRPO on a randomly generated CMDP.

independent trajectories starting from any state and action pair. The assumption of such a generative model is common [17, 11, 31].

The performances of CRPO, NPG-PD, and ARNPG-EPD in the sample-based scenario are shown in Figure 2. Figures 2(a) and 2(b) display the averaged performance, while Figures 2(c) and 2(d) display the performance of the current iterate (a.k.a. last-iterate in optimization literature). It shows that in this sample-based scenario, ARNPG-EPD achieves higher reward values with faster convergence, while all three algorithms satisfy the constraint after a few iterations.

5.3 Acrobot-v1

To demonstrate the efficacy of ARNPG-EPD on complex tasks, we have conducted experiments on the Acrobot-v1 environment from OpenAI Gym [9]. We follow the same experiment setup in [31], where there is a reward value to maximize, and two cost values to be constrained below some thresholds. The superior performance of ARNPG-EPD is shown in Figure 3.

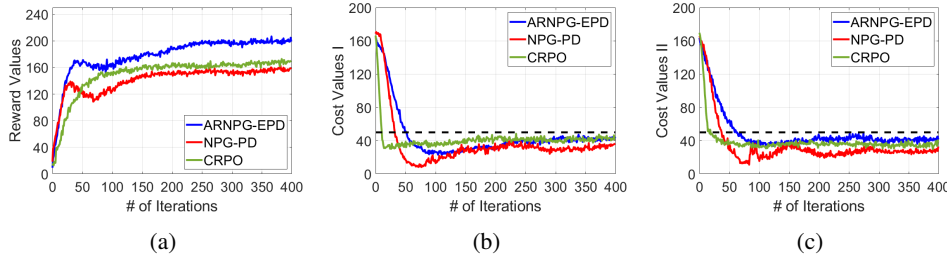


Figure 3: Last-iterate performance for sample-based ARNPG-EPD, NPG-PD, CRPO averaged over 10 random seeds. The black dashed lines in (b) and (c) represent given thresholds.

Figure 3(a) shows that ARNPG-EPD achieves a higher reward value compared to NPG-PD and CRPO, while Figures 3(b) and 3(c) demonstrate that the cost values of all three algorithms are below the thresholds after a few initial iterations. We believe the superiority is due to the new primal-dual design inspired by [34] (discussed in Section 4.2) and the flexibility of choosing t_k in the InnerLoop in the framework. More experiments with different t_k are presented in Appendix A.

6 Conclusion and future works

We propose an ARNPG framework to systematically integrate well-performing first-order methods into the design of policy gradient-based algorithms for multi-objective MDPs. The designed algorithms achieve a global $\mathcal{O}(1/T)$ convergence rate under the softmax parameterization with exact gradients, and empirically have satisfactory performance beyond tabular and exact gradient settings. We believe that ARNPG has potential applications in other scenarios, since the general and flexible framework allows integration with more advanced first-order methods, currently and in the future.

A natural future direction is to extend the theoretical results to more general settings such as function approximation and sample-based scenarios. Viewing ARNPG as a heuristic, the anchor-changing ideas can also be applied to policy optimization for multi-agent RL and meta RL.

7 Acknowledgement

P. R. Kumar’s work is based upon work partially supported by the US Army Contracting Command under W911NF-22-1-0151, US Office of Naval Research under N00014-21-1-2385, 4/21-22 DARES: Army Research Office W911NF-21-20064, US National Science Foundation under CMMI-2038625. The views expressed herein and conclusions contained in this document are those of the authors and should not be interpreted as representing the views or official policies, either expressed or implied, of the U.S. Army Contracting Command, ONR, ARO, NSF, or the United States Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein.

Dileep Kalathil gratefully acknowledges funding from the U.S. National Science Foundation (NSF) grants NSF-CRII-CPS-1850206 and NSF-CAREER-EPCN-2045783.

We thank Dongsheng Ding and Tengyu Xu for generously sharing their code in [11, 31] as baselines.

References

- [1] Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. Optimality and approximation with policy gradient methods in markov decision processes. In Conference on Learning Theory, pages 64–66. PMLR, 2020.
- [2] Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. Journal of Machine Learning Research, 22(98):1–76, 2021.
- [3] Mridul Agarwal and Vaneet Aggarwal. Reinforcement learning for joint optimization of multiple rewards. arXiv preprint arXiv:1909.02940, 2019.
- [4] Eitan Altman. Constrained Markov decision processes, volume 7. CRC Press, 1999.
- [5] Qinbo Bai, Mridul Agarwal, and Vaneet Aggarwal. Joint optimization of multi-objective reinforcement learning with policy gradient based algorithm. arXiv preprint arXiv:2105.14125, 2021.
- [6] Qinbo Bai, Amrit Singh Bedi, Mridul Agarwal, Alec Koppel, and Vaneet Aggarwal. Achieving zero constraint violation for constrained reinforcement learning via primal-dual approach. arXiv preprint arXiv:2109.06332, 2021.
- [7] Amir Beck and Marc Teboulle. Mirror descent and nonlinear projected subgradient methods for convex optimization. Operations Research Letters, 31(3):167–175, 2003.
- [8] Dimitri Bertsekas and Robert Gallager. Data networks. Athena Scientific, 2021.
- [9] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.
- [10] Shicong Cen, Chen Cheng, Yuxin Chen, Yuting Wei, and Yuejie Chi. Fast global convergence of natural policy gradient methods with entropy regularization. Operations Research, 2021.
- [11] Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo R Jovanovic. Natural policy gradient primal-dual method for constrained markov decision processes. In Advances in Neural Information Processing Systems, 2020.
- [12] Arushi Jain, Sharan Vaswani, Reza Babanezhad, Csaba Szepesvari, and Doina Precup. Towards painless policy optimization for constrained mdp. arXiv preprint arXiv:2204.05176, 2022.
- [13] Frank P Kelly, Aman K Maulloo, and David Kim Hong Tan. Rate control for communication networks: shadow prices, proportional fairness and stability. Journal of the Operational Research society, 49(3):237–252, 1998.
- [14] Sajad Khodadadian, Prakirt Raj Jhunjhunwala, Sushil Mahavir Varma, and Siva Theja Maguluri. On the linear convergence of natural policy gradient algorithm. arXiv preprint arXiv:2105.01424, 2021.

- [15] Panagiotis Kyriakis, Jyotirmoy Deshmukh, and Paul Bogdan. Pareto policy adaptation. In International Conference on Learning Representations, 2022.
- [16] Guanghui Lan. First-order and Stochastic Optimization Methods for Machine Learning. Springer Nature, 2020.
- [17] Guanghui Lan. Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. arXiv preprint arXiv:2102.00135, 2021.
- [18] Tianjiao Li, Ziwei Guan, Shaofeng Zou, Tengyu Xu, Yingbin Liang, and Guanghui Lan. Faster algorithm and sharper analysis for constrained markov decision process. arXiv preprint arXiv:2110.10351, 2021.
- [19] Tao Liu, Ruida Zhou, Dileep Kalathil, Panganamala Kumar, and Chao Tian. Learning policies with zero or bounded constraint violation for constrained mdps. Advances in Neural Information Processing Systems, 34, 2021.
- [20] Jincheng Mei, Chenjun Xiao, Csaba Szepesvari, and Dale Schuurmans. On the global convergence rates of softmax policy gradient methods. In International Conference on Machine Learning, pages 6820–6829. PMLR, 2020.
- [21] Arkadi Nemirovski. Prox-method with rate of convergence $o(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. SIAM Journal on Optimization, 15(1):229–251, 2004.
- [22] Sasha Rakhlin and Karthik Sridharan. Optimization, learning, and games with predictable sequences. Advances in Neural Information Processing Systems, 26, 2013.
- [23] Diederik M Roijers, Peter Vamplew, Shimon Whiteson, and Richard Dazeley. A survey of multi-objective sequential decision-making. Journal of Artificial Intelligence Research, 48:67–113, 2013.
- [24] Sandhya Saisubramanian, Ece Kamar, and Shlomo Zilberstein. A multi-objective approach to mitigate negative side effects. In Proceedings of the Twenty-Ninth International Conference on Artificial Intelligence, pages 354–361, 2021.
- [25] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based control. In 2012 IEEE/RSJ international conference on intelligent robots and systems, pages 5026–5033. IEEE, 2012.
- [26] Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Scalarized multi-objective reinforcement learning: Novel design techniques. In 2013 IEEE Symposium on Adaptive Dynamic Programming and Reinforcement Learning (ADPRL), pages 191–199. IEEE, 2013.
- [27] Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Linear last-iterate convergence in constrained saddle-point optimization. In International Conference on Learning Representations, 2020.
- [28] Xiaohan Wei, Hao Yu, and Michael J Neely. Online primal-dual mirror descent under stochastic constraints. Proceedings of the ACM on Measurement and Analysis of Computing Systems, 4(2):1–36, 2020.
- [29] Kyle Hollins Wray, Shlomo Zilberstein, and Abdel-Ilhah Mouaddib. Multi-objective mdps with conditional lexicographic reward preferences. In Twenty-ninth AAAI conference on artificial intelligence, 2015.
- [30] Jingfeng Wu, Vladimir Braverman, and Lin F Yang. Accommodating picky customers: Regret bound and exploration complexity for multi-objective reinforcement learning. arXiv preprint arXiv:2011.13034, 2020.
- [31] Tengyu Xu, Yingbin Liang, and Guanghui Lan. CRPO: A new approach for safe reinforcement learning with convergence guarantee. In International Conference on Machine Learning, pages 11480–11491. PMLR, 2021.

- [32] Runzhe Yang, Xingyuan Sun, and Karthik Narasimhan. A generalized algorithm for multi-objective reinforcement learning and policy adaptation. Advances in Neural Information Processing Systems, 32, 2019.
- [33] Donghao Ying, Yuhao Ding, and Javad Lavaei. A dual approach to constrained markov decision processes with entropy regularization. arXiv preprint arXiv:2110.08923, 2021.
- [34] Hao Yu and Michael J Neely. A primal-dual parallel method with $O(1/\epsilon)$ convergence for constrained composite convex programs. arXiv preprint arXiv:1708.00322, 2017.
- [35] Hao Yu and Michael J Neely. A simple parallel algorithm with an $O(1/t)$ convergence rate for general convex programs. SIAM Journal on Optimization, 27(2):759–783, 2017.
- [36] Junyu Zhang, Amrit Singh Bedi, Mengdi Wang, and Alec Koppel. Beyond cumulative returns via reinforcement learning over state-action occupancy measures. In 2021 American Control Conference (ACC), pages 894–901. IEEE, 2021.
- [37] Junyu Zhang, Alec Koppel, Amrit Singh Bedi, Csaba Szepesvari, and Mengdi Wang. Variational policy gradient method for reinforcement learning with general utilities. Advances in Neural Information Processing Systems, 33:4572–4583, 2020.
- [38] Yiming Zhang, Quan Vuong, and Keith Ross. First order constrained optimization in policy space. Advances in Neural Information Processing Systems, 33:15338–15349, 2020.
- [39] Dongruo Zhou, Jiahao Chen, and Quanquan Gu. Provable multi-objective reinforcement learning with generative models. arXiv preprint arXiv:2011.10134, 2020.

Checklist

The checklist follows the references. Please read the checklist guidelines carefully for information on how to answer these questions. For each question, change the default **[TODO]** to **[Yes]**, **[No]**, or **[N/A]**. You are strongly encouraged to include a **justification to your answer**, either by referencing the appropriate section of your paper or providing a brief inline description. For example:

- Did you include the license to the code and datasets? **[Yes]** See Section 5.
- Did you include the license to the code and datasets? **[No]** The code and the data are proprietary.
- Did you include the license to the code and datasets? **[N/A]**

Please do not modify the questions and only use the provided macros for your answers. Note that the Checklist section does not count towards the page limit. In your paper, please delete this instructions block and only keep the Checklist section heading above along with the questions/answers below.

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]**
 - (b) Did you describe the limitations of your work? **[Yes]** See the last paragraphs of Section 4.2 and 4.3.
 - (c) Did you discuss any potential negative societal impacts of your work? **[N/A]**
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]**
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[Yes]** After the definition of the problems (7),(9),(13)
 - (b) Did you include complete proofs of all theoretical results? **[Yes]** See Appendix C, D, and E.
3. If you ran experiments...

- (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See <https://github.com/tliu1997/ARNPG-MORL>.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Appendix A and B.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] We only report the mean in Figure 3 for clarity. In Appendix A, we will report error bars.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- (a) If your work uses existing assets, did you cite the creators? [Yes]
 - (b) Did you mention the license of the assets? [Yes]
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] See <https://github.com/tliu1997/ARNPG-MORL>.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [N/A]
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
5. If you used crowdsourcing or conducted research with human subjects...
- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

A Experimental settings and additional results for CMDP

A.1 Tabular CMDP

The tabular CMDP, for both exact-gradient scenario (Section 5.1) and sample-based scenario (Section 5.2), follows the same experimental setting as in [11].

The MDP with $m = 2$ objectives represented by $(\mathcal{S}, \mathcal{A}, P, \rho, \gamma, r_{1:2})$ (as the system model in Section 2) is randomly generated, where $|\mathcal{S}| = 20$, $|\mathcal{A}| = 10$, ρ is uniform distribution, and $\gamma = 0.8$. For each $(s, a) \in \mathcal{S} \times \mathcal{A}$, $P(\cdot|s, a) \in \Delta(\mathcal{S})$ is generated by normalizing a random vector $\sim \text{Unif}([0, 1]^{\mathcal{S}})$, and independent rewards $r_1(s, a), r_2(s, a) \sim \text{Unif}([0, 1])$. Choosing the constraint coefficient $b_2 = 3$, the experiments are performed on the CMDP

$$\max_{\theta \in \Theta} V_{r_1}^{\pi_\theta}(\rho) \quad \text{s.t.} \quad V_{r_2}^{\pi_\theta}(\rho) \geq b_2, \quad (15)$$

with the softmax policy class.

For both the exact-gradient scenario (Section 5.1) and the sample-based scenario (Section 5.2), we choose $\eta = 1$ and $\eta' = 1$ for ARNPG-EPD and NPG-PD (following the same hyperparameter selection as in [11]), since both rely on a primal-dual framework. Additionally, we fix $t_k = 1, \forall k = 0, 1, \dots, K - 1$ and select $\alpha = \frac{1-\gamma}{\eta} = 0.2$ for ARNPG-EPD. As for CRPO with exact gradients, we first fix the tolerance parameter as 0.01 and then choose the best learning rate 0.4 from the set $\{0.1, 0.2, \dots, 0.9, 1.0\}$, which enjoys the smallest average optimality gap after 300 iterations. For sample-based CRPO, we select the best learning rate 1.0 from the set $\{0.1, 0.5, 1, 2, 5\}$, which leads to the largest reward value after 300 iterations.

A.2 Acrobot-v1

To demonstrate the performance of the ARNPG-EPD algorithm on more complex tasks with a large state space and multiple constraints, we conduct experiments on the Acrobot-v1 from OpenAI Gym [9]. The acrobot is a planar two-link robotic arm with two joints and two links, where the joint between the links is actuated.

We follow the same experimental setup as in [31], where the task is to swing the end of the lower link to a given height, and there are two costs (objectives) corresponding to (i) the upper link swinging in a prohibited direction; (ii) the lower link swinging in a prohibited direction w.r.t. the upper link. The upper bound constraints on the cost values are (50, 50). An actor-critic framework is adopted, where the actor is defined by a neural softmax policy with two hidden layers of widths (128, 128), and there are 3 critics (one for each of the value functions) also parameterized by neural networks with two hidden layers of widths (128, 128).

Figure 3 provides the performances of the last-iterates generated by the algorithms, averaged over 10 random seeds, where the step size of the dual update (i.e., 0.0005) is tuned from the set $\{0.00001, 0.0005, 0.001, 0.005, 0.01, 0.05\}$, and the tolerance parameter 0.5 of CRPO follows from [31]. For ARNPG-EPD in Figure 3, we choose $t_k = 2, \forall k = 0, 1, \dots, K - 1$ and $\alpha = 1$. To provide more information about variance, Figure 4 includes shaded error bars (\pm standard deviation).

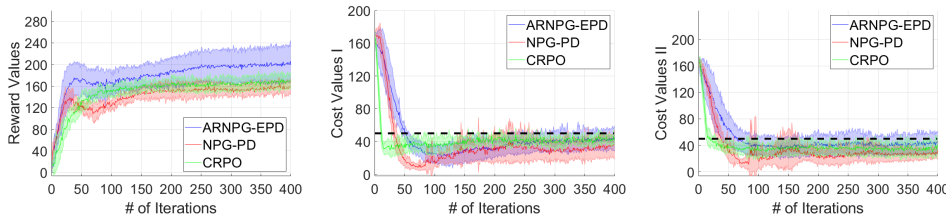


Figure 4: Last-iterate performance for sample-based ARNPG-EPD, NPG-PD, CRPO averaged over 10 random seeds with shaded error bars. The black dashed lines represent given thresholds.

To illustrate the impact of the number of InnerLoop iterations, we fix $\alpha = 1$ and take $t_k = 1, 2, 5, 10$ respectively, as shown in Figure 5. We omit the shaded error bars as in Figure 4 for better visualization and comparison of the mean performance of ARNPG-EPD algorithms under different t_k s. Figure 5

demonstrates the trade-off induced by the selection of the hyperparameter t_k . When t_k is small, e.g., $t_k = 1$, ARNPG-EPD cannot follow the underlying EPD algorithm (for convex optimization) closely. The curve is not as stable as the others, but is more adaptive (becomes feasible first) because it does not spend too much time on the KL-regularized MDP with reward \tilde{r}_k in InnerLoop. When t_k is large, e.g., $t_k = 5$ or $t_k = 10$, the curves are more stable. However, the convergence of the algorithms is relatively slow (become feasible slowly), since larger t_k may waste computation on the regularized objective $\tilde{V}_{k,\alpha}^{\pi_\theta}(\rho)$ instead of focusing on the true optimization program.

In Acrobot-v1, we find $t_k = 2$ to be the best choice empirically. Note that at iteration 300, returns of ARNPG-EPD algorithms are all feasible, and even the lowest mean reward value of 181.8 (ARNPG-EPD with $t_k = 1$), is higher than the mean reward values of NPG-PD and TRPO in Figure 3.

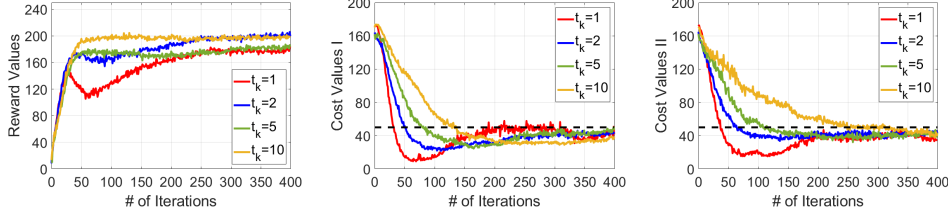


Figure 5: Last-iterate performance for sample-based ARNPG-EPD with different inner loops ($t_k = 1, 2, 5, 10$) averaged over 10 random seeds versus the total number of iterations.

A.3 Hopper-v3

We now demonstrate the performance of our approach on Hopper-v3, a more complex robotics control task, with a constraint of moving speed 82.748 [38]. Hopper-v3 is implemented via the OpenAI Gym based on the MuJoCo physical simulators [25], where both the state space and the action space are continuous. We choose the current SOTA of this task, namely the FOCOPS (First Order Constrained Optimization in Policy Space) algorithm [38], for a comparison with our approach. Since the policy update of the FOCOPS algorithm is based on the PPO (Proximal Policy Optimization) algorithm, for a fair comparison, we also revise our algorithm to a corresponding version called ‘‘ARPPO-EPD’’, where the NPG update is replaced by PPO. It is also an illustration that the anchor-changing idea is readily to be combined with other policy gradient methods. We compare performance with other baselines NPG-PD [11] and CRPO [31].

Figure 6 reports performance averaged over 5 random seeds and illustrates that our ARPPO-EPD algorithms achieve higher reward values compared with other methods, while achieving similar constraint values. (Note that 1, 2, 5 in parentheses of ARPPO-EPD represent the number of inner loops $t_k = 1, 2, 5$.)

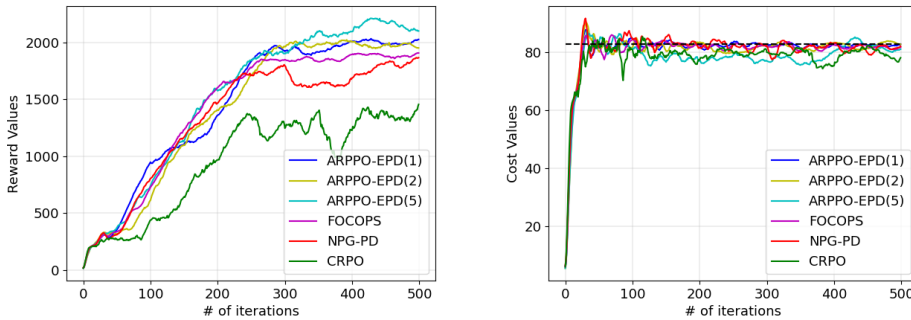


Figure 6: Last-iterate performance for sample-based ARPPO-EPD with different inner loops ($t_k = 1, 2, 5$), FOCOPS, NPG-PD, and CRPO on Hopper-v3 with a constraint on the moving speed 82.748.

B Experimental results for multi-objective MDP with scalarization

In this section, we numerically study the performance of the ARNPG-guided algorithms under the smooth concave scalarization and max-min trade-off, respectively. We first consider the tabular setting for exact-gradient and sample-based scenarios with softmax policy parameterization, and then study the Acrobot-v1 scenario with neural softmax policy parameterization.

B.1 Tabular multi-objective MDP with exact gradients

For the tabular setting, we follow the same MDP construction with $m = 2$ objectives as in Section A.1. The goal of considering exact gradients is to empirically study the theoretically established $O(1/K)$ convergences of the proposed ARNPG-IMD and ARNPG-OMDA algorithms suggested by Theorems 1 and 3. We will show that the proposed algorithms maintain $O(1/K)$ convergence in terms of the macro steps K , for different numbers of micro steps t_k .

B.1.1 Smooth concave scalarization

An example of interest is the sum-logarithmic function, which for two objectives is

$$F(V_{r_1}^\pi(\rho), V_{r_2}^\pi(\rho)) = \log(V_{r_1}^\pi(\rho) + \delta) + \log(V_{r_2}^\pi(\rho) + \delta), \quad (16)$$

where $\delta > 0$ is a small constant. This can be viewed as an approximate formulation of the proportional fairness criteria. If \mathcal{V} is convex and with v^* denoting the optimal value vector achieving the largest $F(v)$, it can be verified by the first-order optimality condition that $\frac{v_1 - v_1^*}{v_1^* + \delta} + \frac{v_2 - v_2^*}{v_2^* + \delta} \leq 0, \forall v_{1:2} \in \mathcal{V}$.

As mentioned in Section 3, when $t_k = 1$, the update in (5) reduces to an NPG update on the unregularized value function $V_{\tilde{r}_k}^{\pi_\theta}(\rho)$. In other words, when $t_k = 1$, α has no impact on the ARNPG-IMD algorithm.

Theorem 1 suggests a lower bound on the number of micro steps t_k . As with many existing algorithms in the optimization literature, the theoretically chosen hyperparameters are usually too conservative, and the convergence rate is maintained over a wider range of hyperparameters. We set up the experiments as follows. We first fix $t_k = 1$ and conduct experiments with learning rates $\eta \in \{0.5, 1.0, \dots, 9.5, 10\}$. The hyperparameter achieving the smallest average optimality gap is $\eta = 4.5$. We then fix the learning rate $\eta = 4.5$, choose the regularization parameter $\alpha = 0.01$, and show the convergence of ARNPG-IMD with $t_k \in \{1, 2, 5, 10\}$ in Figure 7. Though the learning rate $\eta = 4.5$ and the regularization parameter α are not chosen to favor $t_k > 1$, it can still be observed that larger t_k leads to faster convergence in terms of the number of macro steps (K).

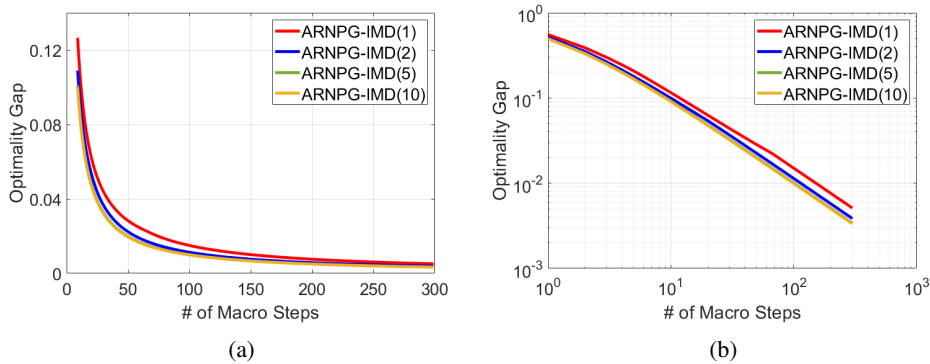


Figure 7: The average optimality gap versus the number of macro steps, for ARNPG-IMD with different number of inner loop iterations ($t_k = 1, 2, 5, 10$) on a randomly generated sum-logarithmic two-objective MDP.

Figure 7(a) illustrates that all ARNPG-IMD algorithms with different t_k s enjoy fast convergence rates and small optimality gaps. Additionally, the log-log plot in Figure 7(b) indicates a slope of approximately -1, which confirms the $O(1/K)$ convergence of ARNPG-IMD.

B.1.2 Max-min trade-off

Another multi-objective MDP of interest is max-min fairness, corresponding to a robust scalarization function, which for the case of two objectives is

$$F(V_1^\pi(\rho), V_2^\pi(\rho)) = \min_i V_i^\pi(\rho) = \min_{\lambda \in \Delta([2])} \langle V_{1:2}^\pi(\rho), \lambda \rangle. \quad (17)$$

The function F is concave but not always differentiable. We can employ a subgradient ascent-based NPG algorithm, which calculates the subgradient of F and performs one-step of NPG for the induced reward function. We call such an algorithm multi-objective NPG (MO-NPG).

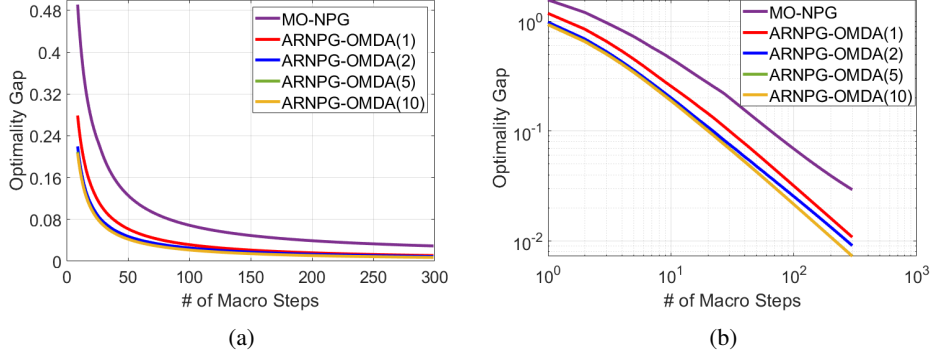


Figure 8: The average optimality gap with respect to the number of macro steps, for MO-NPG and ARNPG-OMDA with different number of inner loop iterations ($t_k = 1, 2, 5, 10$) on a randomly generated max-min two-objective MDP (17).

We conduct the experiments under learning rates $\{0.8, 0.81, \dots, 1.18, 1.19\}$ for MO-NPG and choose the best hyperparameter, 0.93. Similarly, we do a grid-search for ARNPG-OMDA ($t_k = 1$) over $\alpha \in \{1, 2, 5\}$, $\eta \in \{0.06, 0.08, 1.0\}$, $\eta' \in \{0.5, 1.0, 1.5, 2.0\}$ and select the best hyperparameters, $\alpha = 1$, $\eta = 0.08$, and $\eta' = 2$. Then, we fix α, η, η' to explore the impact of $t_k > 1$. Figure 8 shows that the ARNPG-OMDA algorithms converge faster than MO-NPG due to the better underlying optimization algorithm OMDA compared to the subgradient ascent. Moreover, larger t_k gives faster convergence in terms of macro steps K , even though the parameters are not specifically chosen to favor $t_k > 1$.

Recall that under softmax policy with exact gradients, Theorem 3 guarantees $O(1/K)$ convergence of the average optimality gap. Due to the underlining subgradient ascent of MO-NPG, the convergence rate can only be guaranteed by $O(1/\sqrt{K})$. Figure 8(a) shows that the optimality gap of the ARNPG-OMDA algorithms converges faster than that of the MO-NPG algorithm. The corresponding log-log plots in Figure 8(b) shows that the slopes are around -1 for the ARNPG-OMDA algorithms, which demonstrates that the optimality gap of ARNPG-OMDA converges at an $O(1/K)$ rate.

B.2 Sample-based tabular multi-objective MDP

We next consider the same tabular MDP with $m = 2$ objectives, but with the gradients estimated by samples from a generative model that can generate independent trajectories starting from any state and action pair.

B.2.1 Smooth concave scalarization

We conduct the experiments with learning rates $\{0.5, 1.0, \dots, 4.0, 4.5\}$ for sample-based ARNPG-IMD with $t_k = 1$ and choose its best hyperparameter 1.5. To discover the impact of the number of micro steps t_k , we fix the learning rate $\eta = 1.5$ and the regularization parameter $\alpha = 0.005$. Figure 9(a) demonstrates that $t_k > 1$ is helpful to achieve a faster convergence (i.e., during the first 50 iterations, larger t_k leads to a higher scalarized objective), though after 250 iterations they all converge to the same optimal value.

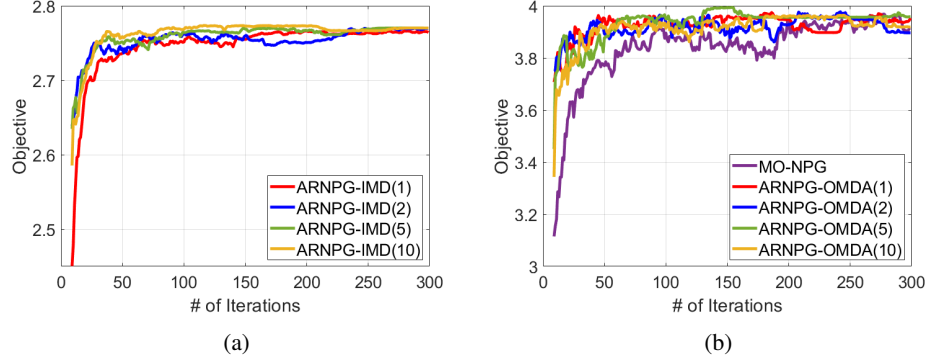


Figure 9: The last-iterate objective value versus the total number of iterations, for sample-based ARNPG-IMD, ARNPG-OMDA, MO-NPG on a randomly generated two-objective MOMDP with (a) sum-logarithmic (16) and (b) max-min trade-off (17).

B.2.2 Max-min trade-off

We conduct the experiments with learning rates $\{0.1, 0.2, \dots, 1.0\}$ for sample-based MO-NPG and choose the best hyperparameter 0.5. We also conduct a grid-search over $\alpha \in \{0.05, 0.1, 0.2\}$, $\eta \in \{0.1, 0.2, 0.5\}$, and $\eta' \in \{0.1, 0.2, 0.5\}$ for sample-based ARNPG-OMDA with $t_k = 1$ and select the best hyperparameter $\alpha = 0.1$, $\eta = 0.5$, and $\eta' = 0.2$. Fixing such α, η, η' , we explore the impact of $t_k > 1$.

Figure 9(b) shows that compared to the MO-NPG algorithm, ARNPG-OMDA algorithms converge faster (i.e., achieve larger scalarized objectives during the first 200 iterations), and that all algorithms approximately find the optimal value after 250 iterations.

B.3 Acrobot-v1

We examine two reward value functions for MOMDP scenarios where the agent is rewarded 1 when it swings the end of the lower link to the given height ranges, i.e., greater than 0.5 and 0-0.1, respectively. The goal is to maximize the objective according to the given criteria. (In this section, we focus on sum-logarithmic scalarization.) To show the impact of t_k in the ARNPG framework, we fix $\alpha = 1$ and choose InnerLoop iterations $t_k = 1, 2, 5, 10$ respectively.

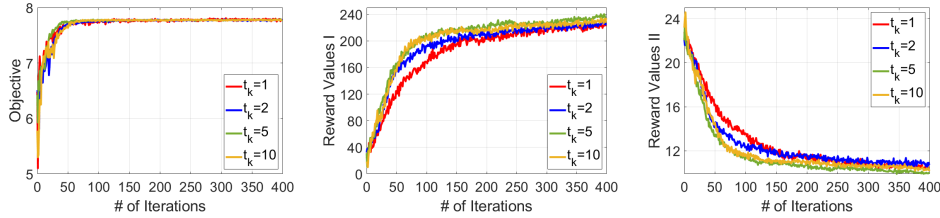


Figure 10: Last-iterate performance for sample-based ARNPG-IMD with different inner loops ($t_k = 1, 2, 5, 10$) averaged over 10 random seeds.

It can be seen from Figure 10 that the scalarized objectives are indistinguishable because the gradient of logarithmic functions is small when the input values become large. With respect to the two reward objectives (reward values I and reward values II), we can observe some trade-offs in that $t_k = 5$ converges faster than other t_k s.

C Supporting lemmas

Before delving into detailed proofs for the proposition and theorems, we introduce some supporting lemmas.

Recall that the Bregman divergence generated by a convex differentiable function $h(\cdot)$ is

$$B_h(x, y) := h(x) - h(y) - \langle \nabla h(y), x - y \rangle.$$

The fundamental inequality (2) associated with mirror ascent is formally presented in the following lemma.

Lemma 1. *Let $B_h : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be a Bregman divergence function, $\mathcal{X} \subset \mathbb{R}^n$ be a compact convex set, and $g \in \mathbb{R}^n$. Suppose $x' = \arg \max_{y \in \mathcal{X}} \{\langle g, y \rangle - \alpha B_h(y||x)\}$ for a fixed $x \in \mathcal{X}$ and $\alpha > 0$. Then for any $y \in \mathcal{X}$,*

$$\langle g, x' \rangle - \alpha B_h(x'||x) \geq \langle g, y \rangle - \alpha B_h(y||x) + \alpha B_h(y||x').$$

Inequalities of the same form have appeared in many previous works, e.g., Lemma 3.4 in [16] and a case of \mathcal{X} being a probability simplex (Lemma 2.1 in [28]). For completeness, we provide a proof of Lemma 1.

Proof of Lemma 1. Since h is proper and convex, $x' := \arg \max_{y \in \mathcal{X}} \{\langle g, y \rangle - \alpha B_h(y||x)\}$ exists and satisfies the first order condition

$$\langle g - \alpha \nabla h(x') + \alpha \nabla h(x), x' - y \rangle = \langle g - \alpha \nabla_{x'} B_h(x'||x), x' - y \rangle \geq 0, \quad \forall y \in \mathcal{X},$$

which implies $\langle g, x' - y \rangle \geq \alpha \langle \nabla h(x') - \nabla h(x), x' - y \rangle$. It can be verified that

$$\langle \nabla h(x') - \nabla h(x), x' - y \rangle = B_h(x'||x) - B_h(y||x) + B_h(y||x').$$

We can conclude the proof by substituting the equation into the previous inequality. \square

The following lemma draws a connection between the ℓ_1 difference of state-action visitation distributions and averaged KL-divergence.

Lemma 2. *Let $d_{\rho}^{\pi'}, d_{\rho}^{\pi}$ be two discounted state-action visitation distributions corresponding to policies π' and π . Then*

$$\|d_{\rho}^{\pi'} - d_{\rho}^{\pi}\|_1 \leq \frac{\gamma\sqrt{2}}{1-\gamma} \sqrt{\min \left(D_{d_{\rho}^{\pi'}}(\pi'||\pi), D_{d_{\rho}^{\pi}}(\pi||\pi'), D_{d_{\rho}^{\pi}}(\pi'||\pi), D_{d_{\rho}^{\pi}}(\pi||\pi') \right)}.$$

Proof. Let $d_{\rho,h}^{\pi}(\cdot, \cdot)$ be the state-action visitation distribution at step h , which implies $\frac{1}{1-\gamma} d_{\rho}^{\pi}(\cdot, \cdot) = \sum_{h \geq 0} \gamma^h d_{\rho,h}^{\pi}(\cdot, \cdot)$. Denote $\tilde{\pi}_h$ as the policy that implements policy π for the first h steps and then commits to policy π' thereafter. Denote its corresponding discounted state-action visitation distribution by $d_{\rho}^{\tilde{\pi}_h}$. It follows that

$$\begin{aligned} \frac{1}{1-\gamma} \|d_{\rho}^{\pi'} - d_{\rho}^{\pi}\|_1 &\stackrel{(a)}{=} \frac{1}{1-\gamma} \left\| \sum_{h=0}^{\infty} (d_{\rho}^{\tilde{\pi}_h} - d_{\rho}^{\tilde{\pi}_{h+1}}) \right\|_1 \stackrel{(b)}{\leq} \frac{1}{1-\gamma} \sum_{h=0}^{\infty} \|d_{\rho}^{\tilde{\pi}_h} - d_{\rho}^{\tilde{\pi}_{h+1}}\|_1 \\ &= \sum_{h=0}^{\infty} \left\| \sum_{t=0}^{\infty} \gamma^t (d_{\rho,t}^{\tilde{\pi}_h} - d_{\rho,t}^{\tilde{\pi}_{h+1}}) \right\|_1 = \sum_{h=0}^{\infty} \left\| \sum_{t=h+1}^{\infty} \gamma^t (d_{\rho,t}^{\tilde{\pi}_h} - d_{\rho,t}^{\tilde{\pi}_{h+1}}) \right\|_1 \\ &\stackrel{(c)}{\leq} \sum_{h=0}^{\infty} \sum_{t=h+1}^{\infty} \gamma^t \|d_{\rho,t}^{\tilde{\pi}_h} - d_{\rho,t}^{\tilde{\pi}_{h+1}}\|_1 \stackrel{(d)}{\leq} \sum_{h=0}^{\infty} \sum_{t=h+1}^{\infty} \gamma^t \|d_{\rho,h}^{\tilde{\pi}_h} - d_{\rho,h}^{\tilde{\pi}_{h+1}}\|_1 \\ &= \frac{\gamma}{1-\gamma} \sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim d_{\rho,h}^{\pi}} \|\pi(\cdot|s) - \pi'(\cdot|s)\|_1 \\ &\stackrel{(e)}{\leq} \frac{\gamma}{1-\gamma} \sqrt{\left(\sum_{h=0}^{\infty} \gamma^h \right) \left(\sum_{h=0}^{\infty} \gamma^h \mathbb{E}_{s \sim d_{\rho,h}^{\pi}} \|\pi(\cdot|s) - \pi'(\cdot|s)\|_1^2 \right)} \\ &= \frac{\gamma}{(1-\gamma)^2} \sqrt{\mathbb{E}_{s \sim d_{\rho}^{\pi}} \|\pi(\cdot|s) - \pi'(\cdot|s)\|_1^2}. \end{aligned}$$

Above, (a) holds by telescoping, (b) and (c) hold due to the triangle inequality of ℓ_1 -norm and the definition of $\tilde{\pi}_h$, (d) hold owing to the data processing inequality for f -divergence $\|\cdot\|_1$, and (e)

holds due to the Cauchy-Schwarz inequality. Due to the symmetry between π and π' , it can be similarly derived

$$\|d_{\rho}^{\pi'} - d_{\rho}^{\pi}\|_1 \leq \frac{\gamma}{1-\gamma} \sqrt{\mathbb{E}_{s \sim d_{\rho}^{\pi'}} \|\pi(\cdot|s) - \pi'(\cdot|s)\|_1^2}.$$

We can conclude the proof by further applying Pinsker's inequality. \square

An application of Lemma 2 gives an upper bound on the difference between value function vectors as follows.

Lemma 3. *For any $k = 0, 1, \dots, K-1$,*

$$\frac{1}{2} \|V_{1:m}^{\pi_k}(\rho) - V_{1:m}^{\pi_{k+1}}(\rho)\|_{\infty}^2 \leq \frac{\gamma^2}{(1-\gamma)^4} D_{d_{\rho}^{\pi_{k+1}}}(\pi_{k+1} \|\pi_k).$$

Proof. For any $i = 1, 2, \dots, m$, we have

$$\begin{aligned} |V_i^{\pi_k}(\rho) - V_i^{\pi_{k+1}}(\rho)| &= \frac{1}{1-\gamma} \left| \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} r_i(s,a) (d_{\rho}^{\pi_k}(s,a) - d_{\rho}^{\pi_{k+1}}(s,a)) \right| \\ &\leq \frac{1}{1-\gamma} \|d_{\rho}^{\pi_k} - d_{\rho}^{\pi_{k+1}}\|_1 \leq \frac{\gamma\sqrt{2}}{(1-\gamma)^2} \sqrt{D_{d_{\rho}^{\pi_{k+1}}}(\pi_{k+1} \|\pi_k)}, \end{aligned}$$

where the last inequality is due to Lemma 2. \square

D Proof in Section 3

This section presents the formal proof of Proposition 1. We begin by presenting some properties of InnerLoop. We shall omit θ in π_{θ} , since the policies are under softmax parameterization.

D.1 Linear convergence of InnerLoop

InnerLoop($\tilde{r}_k, \pi_k, \alpha, \eta, t_k$) approximately solves the following KL-regularized MDP via natural policy gradient. Note that

$$\tilde{V}_{k,\alpha}^{\pi}(s) = \mathbb{E} \left[\sum_{t \geq 0} \gamma^t (\tilde{r}_k(s_t, a_t) + \alpha \log \pi_k(a_t|s_t) - \alpha \log \pi(a_t|s_t)) \mid s_0 = s, \pi \right], \quad (18)$$

which can be viewed as an entropy regularized value with reward function $\tilde{r}_k(s, a) + \alpha \log \pi_k(a|s)$. The entropy-regularized state-action value function is then defined as [10]

$$\tilde{Q}_{k,\alpha}^{\pi}(s, a) = \tilde{r}_k(s, a) + \alpha \log \pi_k(a|s) + \gamma \mathbb{E}_{s' \sim P(\cdot|s,a)} [\tilde{V}_{k,\alpha}^{\pi}(s')]. \quad (19)$$

The convergence of NPG in entropy-regularized MDP has been well-studied in [10], with the key results summarized in the following lemma.

Lemma 4 (Linear convergence of entropy-regularized NPG, Theorem 1 in [10]). *For any learning rate $0 < \eta \leq (1-\gamma)/\alpha$ and any $k = 0, 1, \dots, K-1$, the entropy-regularized NPG updates satisfy*

$$\begin{aligned} \left\| \tilde{Q}_{k,\alpha}^{\pi_k^*} - \tilde{Q}_{k,\alpha}^{\pi_k^{(t+1)}} \right\|_{\infty} &\leq C_k \gamma (1 - \eta \alpha)^t, \\ \left\| \log \pi_k^* - \log \pi_k^{(t+1)} \right\|_{\infty} &\leq 2C_k \alpha^{-1} (1 - \eta \alpha)^t, \\ \left\| \tilde{V}_{k,\alpha}^{\pi_k^*} - \tilde{V}_{k,\alpha}^{\pi_k^{(t+1)}} \right\|_{\infty} &\leq 3C_k (1 - \eta \alpha)^t, \end{aligned}$$

for all $t \geq 0$, where C_k satisfies $C_k \geq \left\| \tilde{Q}_{k,\alpha}^{\pi_k^*} - \tilde{Q}_{k,\alpha}^{\pi_k^{(0)}} \right\|_{\infty} + 2\alpha \left(1 - \frac{\eta\alpha}{1-\gamma}\right) \left\| \log \pi_k^* - \log \pi_k^{(0)} \right\|_{\infty}$.

Remark. There is a typographical mistake in the inequality “ $\|\tilde{V}_{k,\alpha}^{\pi_k^*} - \tilde{V}_{k,\alpha}^{\pi_k^{(t+1)}}\|_\infty \leq 3\gamma C_k(1-\eta\alpha)^t$ ” in [10], and it has been corrected here. It is not hard to verify that the proofs of the inequalities in Lemma 4 [10] hold without the assumption that $0 \leq r(s, a) \leq 1$.

Denote $\tilde{V}_k^\pi(s) := V_{\tilde{r}_k}^\pi(s)$. For the regularized MDP, its optimal policy is uniformly optimal, i.e., for any state $s \in \mathcal{S}$,

$$\frac{1}{1-\gamma} \|\tilde{r}_k\|_\infty \geq \tilde{V}_k^{\pi_k^*}(s) \geq \tilde{V}_k^{\pi_k}(s) - \frac{\alpha}{1-\gamma} D_{d_s^{\pi_k^*}}(\pi_k^* || \pi_k) = \tilde{V}_{k,\alpha}^{\pi_k^*}(s) \geq \tilde{V}_{k,\alpha}^{\pi_k}(s) = \tilde{V}_k^{\pi_k}(s). \quad (20)$$

It follows that $\forall (s, a) \in \mathcal{S} \times \mathcal{A}$,

$$\left| \tilde{Q}_{k,\alpha}^{\pi_k^*}(s, a) - \tilde{Q}_{k,\alpha}^{\pi_k}(s, a) \right| = \gamma \sum_{s' \in \mathcal{S}} P(s'|s, a) \left| \tilde{V}_{k,\alpha}^{\pi_k^*}(s') - \tilde{V}_{k,\alpha}^{\pi_k}(s') \right| \stackrel{(a)}{\leq} \frac{\gamma \|\tilde{r}_k\|_\infty}{1-\gamma},$$

where (a) holds due to the relation in (20). It implies $\|\tilde{Q}_{k,\alpha}^{\pi_k^*} - \tilde{Q}_{k,\alpha}^{\pi_k}\|_\infty \leq \frac{\gamma \|\tilde{r}_k\|_\infty}{1-\gamma}$. Since $1 - \frac{\eta\alpha}{1-\gamma} = 0$ when $\eta = \frac{1-\gamma}{\alpha}$, we can apply results in Lemma 4 with $C_k = \frac{\gamma \|\tilde{r}_k\|_\infty}{1-\gamma}$, which gives

$$-\tilde{V}_k^{\pi_{k+1}}(\rho) + \frac{\alpha}{1-\gamma} D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1} || \pi_k) \leq -\tilde{V}_k^{\pi_k^*}(\rho) + \frac{\alpha}{1-\gamma} D_{d_\rho^{\pi_k^*}}(\pi_k^* || \pi_k) + 3C_k(1-\eta\alpha)^{t_k}. \quad (21)$$

D.2 Hidden convexity in state-action visitation distribution

Noting that the class of softmax policies is almost complete in the sense that its closure contains all stationary policies, we will omit the parameter θ in π_θ . The set of achievable state-action visitations is $\mathcal{D} = \{d \in \Delta(\mathcal{S} \times \mathcal{A}) : \gamma \sum_{s', a'} P(s'|s, a') d(s', a') + (1-\gamma)\rho(s) = \sum_a d(s, a), \forall s \in \mathcal{S}\}$, which is a convex compact set.

For any policies π, π' , define a pseudo KL-divergence between $d_\rho^\pi, d_\rho^{\pi'} \in \mathcal{D}_\rho$ by

$$\tilde{D}(d_\rho^\pi || d_\rho^{\pi'}) := \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} d_\rho^\pi(s, a) \log \frac{d_\rho^\pi(s, a)/d_\rho^\pi(s)}{d_\rho^{\pi'}(s, a)/d_\rho^{\pi'}(s)}. \quad (22)$$

It is not hard to verify that

$$\begin{aligned} D_{d_\rho^\pi}(\pi || \pi') &= \sum_{s \in \mathcal{S}} d_\rho^\pi(s) \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{\pi(a|s)}{\pi'(a|s)} = \sum_{s \in \mathcal{S}} d_\rho^\pi(s) \sum_{a \in \mathcal{A}} \frac{d_\rho^\pi(s, a)}{d_\rho^\pi(s)} \log \frac{d_\rho^\pi(s, a)/d_\rho^\pi(s)}{d_\rho^{\pi'}(s, a)/d_\rho^{\pi'}(s)} \\ &= \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} d_\rho^\pi(s, a) \log \frac{d_\rho^\pi(s, a)/d_\rho^\pi(s)}{d_\rho^{\pi'}(s, a)/d_\rho^{\pi'}(s)} = \tilde{D}(d_\rho^\pi || d_\rho^{\pi'}). \end{aligned} \quad (23)$$

This equation bridges the state-action visitation space and the policy space. The following lemma shows that the pseudo KL-divergence defined in (22) is actually a Bregman divergence between state-action visitation distributions.

Lemma 5. *The pseudo KL-divergence $\tilde{D}(d_\rho^\pi || d_\rho^{\pi'})$ defined in (22) is a Bregman divergence $B_h(d_\rho^\pi || d_\rho^{\pi'})$ generated by the convex function*

$$h(d_\rho^\pi) = \sum_{(s,a) \in \mathcal{S} \times \mathcal{A}} d_\rho^\pi(s, a) \log d_\rho^\pi(s, a) - \sum_{s \in \mathcal{S}} d_\rho^\pi(s) \log d_\rho^\pi(s).$$

Proof of Lemma 5. It can be verified by elementary algebra that

$$\tilde{D}(d_\rho^\pi || d_\rho^{\pi'}) = h(d_\rho^\pi) - h(d_\rho^{\pi'}) - \langle \nabla h(d_\rho^{\pi'}), d_\rho^\pi - d_\rho^{\pi'} \rangle,$$

where $\nabla_{(s,a)} h(d_\rho^{\pi'}) = \log d_\rho^{\pi'}(s, a) - \log d_\rho^{\pi'}(s)$. Hence we only need to show that $h(d_\rho^\pi)$ is convex. The Hessian matrix of function $h(d_\rho^\pi)$ can be calculated as $\text{diag}(H_1, H_2, \dots, H_{|\mathcal{S}|})$, where

$H_s = \frac{1}{d_\rho^\pi(s)} (\text{diag}(d_\rho^\pi(s)/d_\rho^\pi(s, \cdot)) - \mathbf{1}\mathbf{1}^T)$ is an $|\mathcal{A}| \times |\mathcal{A}|$ matrix corresponding to state s . For each H_s , we know for any $x_{1:|\mathcal{A}|} \in \mathbb{R}^{|\mathcal{A}|}$,

$$\begin{aligned} x^T H_s x &= \frac{1}{d_\rho^\pi(s)} \left(\sum_{a \in \mathcal{A}} \frac{d_\rho^\pi(s)}{d_\rho^\pi(s, a)} x_a^2 - \left(\sum_{a \in \mathcal{A}} x_a \right)^2 \right) \\ &= \frac{1}{d_\rho^\pi(s)} \left(\left(\sum_{a \in \mathcal{A}} \frac{d_\rho^\pi(s, a)}{d_\rho^\pi(s)} \right) \left(\sum_{a \in \mathcal{A}} \frac{d_\rho^\pi(s)}{d_\rho^\pi(s, a)} x_a^2 \right) - \left(\sum_{a \in \mathcal{A}} x_a \right)^2 \right) \\ &\stackrel{(a)}{\geq} \frac{1}{d_\rho^\pi(s)} \left(\left(\sum_{a \in \mathcal{A}} |x_a| \right)^2 - \left(\sum_{a \in \mathcal{A}} x_a \right)^2 \right) \geq 0, \end{aligned}$$

where (a) is due to the Cauchy-Schwarz inequality. Thus the Hessian matrix of $h(d_\rho^\pi)$ is positive semi-definite, which implies that $h(d_\rho^\pi)$ is convex. \square

InnerLoop of the ARNPG framework is solving a KL-regularized MDP with value as in (4),

$$\tilde{V}_{k, \alpha}^{\pi_\theta}(\rho) = V_{\tilde{r}_k}^{\pi_\theta}(\rho) - \alpha \frac{D_{d_\rho^\pi}(\pi_\theta || \pi_{\theta_k})}{1 - \gamma}.$$

This optimization can be equivalently represented by viewing state-action visitation as the decision variables:

$$\max_{\pi} V_{\tilde{r}_k}^{\pi}(\rho) - \alpha \frac{D_{d_\rho^\pi}(\pi || \pi_k)}{1 - \gamma} \Leftrightarrow \max_{d \in \mathcal{D}} \langle \tilde{r}_k, d \rangle - \alpha \tilde{D}(d || d_\rho^{\pi_k}). \quad (24)$$

Here \Leftrightarrow means that they are equivalent in the sense that the optimal policy solution π_k^* for the former optimization and the optimal visitation solution d_k^* for the latter satisfy $d_k^* = d_\rho^{\pi_k^*}$. Note that $\tilde{V}_{\tilde{r}_k}^{\pi}(\rho) = \frac{1}{1-\gamma} \langle \tilde{r}_k, d_\rho^\pi \rangle$ is a linear function of d_ρ^π , $\tilde{D}(\cdot || \cdot)$ is a Bregman divergence, and \mathcal{D} is compact. We can apply Lemma 1 on the latter optimization and have

$$\langle \tilde{r}_k, d_k^* \rangle - \alpha \tilde{D}(d_k^* || d_\rho^{\pi_k}) \geq \langle \tilde{r}_k, d \rangle - \alpha \tilde{D}(d || d_\rho^{\pi_k}) + \alpha \tilde{D}(d || d_k^*), \quad \forall d \in \mathcal{D}. \quad (25)$$

Since the policy class and the state-action visitation class are both complete, the inequality above implies that

$$V_{\tilde{r}_k}^{\pi_k^*}(\rho) - \alpha \frac{D_{d_\rho^{\pi_k^*}}(\pi_k^* || \pi_k)}{1 - \gamma} \geq V_{\tilde{r}_k}^{\pi}(\rho) - \alpha \frac{D_{d_\rho^\pi}(\pi || \pi_k) - D_{d_\rho^\pi}(\pi || \pi_k^*)}{1 - \gamma}, \quad \forall \pi. \quad (26)$$

InnerLoop does not seek to find the precise solution π_k^* but approximates it with $\pi_{k+1} = \pi_k^{(t_k)}$ via t_k micro-step iterations. Proposition 1 provides a quantitative bound regarding the approximation error of π_{k+1} .

D.3 Proof of Proposition 1

Proof of Proposition 1. Combining (21) and (26) gives

$$\begin{aligned} -\tilde{V}_k^{\pi_{k+1}}(\rho) + \alpha \frac{D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1} || \pi_k)}{1 - \gamma} &\leq -\tilde{V}_k^{\pi}(\rho) + \alpha \frac{D_{d_\rho^\pi}(\pi || \pi_k) - D_{d_\rho^\pi}(\pi || \pi_{k+1})}{1 - \gamma} \\ &\quad + 3C_k(1 - \eta\alpha)^{t_k} + \alpha \frac{D_{d_\rho^\pi}(\pi || \pi_{k+1}) - D_{d_\rho^\pi}(\pi || \pi_k^*)}{1 - \gamma}. \end{aligned}$$

Note that

$$\begin{aligned} D_{d_\rho^\pi}(\pi || \pi_{k+1}) - D_{d_\rho^\pi}(\pi || \pi_k^*) &= \left\langle d_\rho^\pi(\cdot, \cdot), \log \frac{\pi_k^*(\cdot, \cdot)}{\pi_{k+1}(\cdot, \cdot)} \right\rangle \\ &\leq \|d_\rho^\pi\|_1 \|\log \pi_k^* - \log \pi_{k+1}\|_\infty = \|\log \pi_k^* - \log \pi_{k+1}\|_\infty \leq 2C_k \alpha^{-1} (1 - \eta\alpha)^{t_k}, \end{aligned}$$

where the first inequality follows from Cauchy-Schwartz, and the last inequality is due to Lemma 4. We thus have

$$-\tilde{V}_k^{\pi_{k+1}}(\rho) + \alpha \frac{D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1}||\pi_k)}{1-\gamma} \leq -\tilde{V}_k^\pi(\rho) + \alpha \frac{D_{d_\rho^\pi}(\pi||\pi_k) - D_{d_\rho^\pi}(\pi||\pi_{k+1})}{1-\gamma} + \frac{5C_k(1-\eta\alpha)^{t_k}}{1-\gamma}.$$

We then conclude the proposition, since $\frac{5C_k(1-\eta\alpha)^t}{1-\gamma} \leq \epsilon_k$ can be guaranteed by $t_k \geq \frac{1}{1-\gamma} \log(\frac{5\gamma\|\tilde{r}_k\|_\infty}{(1-\gamma)^2\epsilon_k})$. \square

E Proof in Section 4

E.1 ARNPG-IMD for smooth scalarization

Proof of Theorem 1. By $|\tilde{r}_k(s, a)| = |\langle \tilde{G}_k, r_{1:m}(s, a) \rangle| \leq \|\tilde{G}_k\|_1 \|r_{1:m}(s, a)\|_\infty \leq L$, we know $\|\tilde{r}_k\|_\infty \leq L$. Recall $\alpha \geq \frac{\beta}{(1-\gamma)^3}$. Taking $\epsilon_k = \frac{\alpha \log(|\mathcal{A}|)}{(1-\gamma)K}$, we choose $t_k = \lceil \frac{1}{1-\gamma} \log(\frac{5LK}{\beta \log(|\mathcal{A}|)}) + 1 \rceil$. Thus by Proposition 1, for any policy π , we have the fundamental inequality

$$V_{\tilde{r}_k}^{\pi_{k+1}}(\rho) - \alpha \frac{D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1}||\pi_k)}{1-\gamma} \geq V_{\tilde{r}_k}^\pi(\rho) - \alpha \frac{D_{d_\rho^\pi}(\pi||\pi_k) - D_{d_\rho^\pi}(\pi||\pi_{k+1})}{1-\gamma} - \epsilon_k. \quad (27)$$

For the RHS of (27), by the concavity of F , we have

$$V_{\tilde{r}_k}^\pi(\rho) - V_{\tilde{r}_k}^{\pi_k}(\rho) = \langle \tilde{G}_k, V_{1:m}^\pi(\rho) - V_{1:m}^{\pi_k}(\rho) \rangle \geq F(V_{1:m}^\pi(\rho)) - F(V_{1:m}^{\pi_k}(\rho)).$$

For the LHS of (27), by the fact that F is β -smooth, we know

$$\begin{aligned} V_{\tilde{r}_k}^{\pi_{k+1}}(\rho) - V_{\tilde{r}_k}^{\pi_k}(\rho) &= \langle \tilde{G}_k, V_{1:m}^{\pi_{k+1}}(\rho) - V_{1:m}^{\pi_k}(\rho) \rangle \\ &\leq F(V_{1:m}^{\pi_{k+1}}(\rho)) - F(V_{1:m}^{\pi_k}(\rho)) + \frac{\beta}{2} \|V_{1:m}^{\pi_k}(\rho) - V_{1:m}^{\pi_{k+1}}(\rho)\|_\infty^2. \end{aligned}$$

From Lemma 3 and recalling $\alpha \geq \frac{\beta}{(1-\gamma)^3}$,

$$\frac{\beta}{2} \|V_{1:m}^{\pi_k}(\rho) - V_{1:m}^{\pi_{k+1}}(\rho)\|_\infty^2 \leq \frac{\gamma^2 \beta}{(1-\gamma)^4} D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1}||\pi_k) \leq \alpha \frac{D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1}||\pi_k)}{1-\gamma}.$$

Substituting these three inequalities into the fundamental inequality (27), telescoping from $k = 0$ to $K - 1$, and selecting $\pi = \pi^*$, we can conclude that

$$\frac{1}{K} \sum_{k=1}^K F(V_{1:m}^{\pi_k}(\rho)) \geq F(V_{1:m}^{\pi^*}(\rho)) - \frac{\alpha D_{d_\rho^{\pi^*}}(\pi^*||\pi_0)}{(1-\gamma)K} - \frac{1}{K} \sum_{k=0}^{K-1} \epsilon_k \geq F(V_{1:m}^{\pi^*}(\rho)) - \frac{2\alpha \log(|\mathcal{A}|)}{(1-\gamma)K}.$$

\square

Proof of Corollary 1. Note that $T = \sum_{k=0}^{K-1} t_k = \Theta(\frac{K}{1-\gamma} \log(K))$. It implies $\frac{K}{1-\gamma} = \Theta(T/\log(T))$. Substituting this into Theorem 1 concludes Corollary 1. \square

E.2 ARNPG-EPD for CMDP

We first introduce the properties of the Lagrange multiplier updates (10) in the following lemma.

Lemma 6 (Properties of Lagrange multiplier updates). *Based on the update of the Lagrange multipliers λ_k , for any $i \in [2 : m]$ we have:*

1. At any macro step k , $\lambda_{k,i} \geq 0$.
2. At any macro step k , $\lambda_{k,i} + \eta'(b_i - V_i^{\pi_k}(\rho)) \geq 0$.
3. At macro step 0, $|\lambda_{0,i}| \leq \eta'|V_i^{\pi_0}(\rho) - b_i|$; at any macro step $k > 0$, $|\lambda_{k,i}| \geq \eta'|V_i^{\pi_k}(\rho) - b_i|$.

Remark. The first property guarantees the feasibility of the Lagrange multipliers; the second property ensures that the Lagrangian in the inner loop can indeed maximize the constraint rewards; and the third property is a key supporting step for the analysis of the constraint violation.

Proof of Lemma 6. Taking any $i \in [2 : m]$, we prove each property respectively.

1. Note that $\lambda_{0,i} = \max\{0, \eta'(V_i^{\pi_0}(\rho) - b_i)\} \geq 0$ by initialization. Suppose $\lambda_{k,i} \geq 0$. The update is $\lambda_{k+1,i} = \max\{\eta'(V_i^{\pi_{k+1}}(\rho) - b_i), \lambda_{k,i} + \eta'(b_i - V_i^{\pi_{k+1}}(\rho))\}$.
 If $b_i - V_i^{\pi_{k+1}}(\rho) < 0$, then $\lambda_{k+1,i} \geq 0$ by the first component in the $\max\{\cdot, \cdot\}$.
 If $b_i - V_i^{\pi_{k+1}}(\rho) \geq 0$, then $\lambda_{k+1,i} \geq 0$ by the second component in the $\max\{\cdot, \cdot\}$.
 Thus, $\lambda_{k+1,i} \geq 0$, and property can be proved by induction.
2. For $k = 0$, $\lambda_{0,i} + \eta'(b_i - V_i^{\pi_0}(\rho)) = \max\{\eta'(b_i - V_i^{\pi_0}(\rho)), 0\} \geq 0$.
 The update is $\lambda_{k+1,i} = \max\{\eta'(V_i^{\pi_{k+1}}(\rho) - b_i), \lambda_{k,i} + \eta'(b_i - V_i^{\pi_{k+1}}(\rho))\}$. Thus for $k \geq 0$, $\lambda_{k+1,i} + \eta'(b_i - V_i^{\pi_{k+1}}(\rho)) = \max\{0, \lambda_{k,i} + 2\eta'(b_i - V_i^{\pi_{k+1}}(\rho))\} \geq 0$.
3. For $k = 0$, the initialization is $\lambda_{0,i} = \max\{0, \eta'(V_i^{\pi_0}(\rho) - b_i)\}$.
 If $V_i^{\pi_0}(\rho) - b_i \leq 0$, then $\lambda_{0,i} = 0$ and $|\lambda_{0,i}| \leq \eta'|V_i^{\pi_0}(\rho) - b_i|$.
 If $V_i^{\pi_0}(\rho) - b_i > 0$, then $\lambda_{0,i} = \eta'(V_i^{\pi_0}(\rho) - b_i)$ and $|\lambda_{0,i}| = \eta'|V_i^{\pi_0}(\rho) - b_i|$.
 For $k \geq 0$, the update is $\lambda_{k+1,i} = \max\{\eta'(V_i^{\pi_{k+1}}(\rho) - b_i), \lambda_{k,i} + \eta'(b_i - V_i^{\pi_{k+1}}(\rho))\}$.
 If $V_i^{\pi_{k+1}}(\rho) - b_i \leq 0$, then $\lambda_{k+1,i} = \lambda_{k,i} + \eta'(b_i - V_i^{\pi_{k+1}}(\rho))$, and $|\lambda_{k+1,i}| = \lambda_{k,i} + \eta'|V_i^{\pi_{k+1}}(\rho) - b_i| \geq \eta'|V_i^{\pi_{k+1}}(\rho) - b_i|$ by the first property that $\lambda_{k,i} \geq 0$.
 If $V_i^{\pi_{k+1}}(\rho) - b_i > 0$, then $\lambda_{k+1,i} \geq \eta'(V_i^{\pi_{k+1}}(\rho) - b_i) > 0$. Thus $|\lambda_{k+1,i}| \geq \eta'|V_i^{\pi_{k+1}}(\rho) - b_i|$.

□

We now analyze the optimality gap and constraint violation separately.

E.2.1 Optimality gap of ARNPG-EPD

Recall the definition of the reward in the ascent direction

$$\tilde{r}_k(s, a) = r_1(s, a) + \sum_{i=2}^m [\lambda_{k,i} + \eta'(b_i - V_i^{\pi_k}(\rho))] r_i(s, a). \quad (28)$$

Since $r_i(s, a) \leq 1$, we can verify that $|\tilde{r}_k(s, a)| \leq 1 + \frac{\eta'(m-1)}{1-\gamma} + \sum_{i=2}^m \lambda_{k,i} =: L_k$, which implies $\|\tilde{r}_k\|_\infty \leq L_k$. Taking $\epsilon_k = \frac{\alpha \log(|\mathcal{A}|)}{(1-\gamma)K}$, we choose $t_k = \lceil \frac{1}{1-\gamma} \log(\frac{5L_k K}{2\eta' m \log(|\mathcal{A}|)}) + 1 \rceil$.

Since $\lambda_{k,i} + \eta'(b_i - V_i^{\pi_k}(\rho)) \geq 0$ by the second property in Lemma 6, and $V_i^{\pi^*}(\rho) \geq b_i$ for any $i \in [2 : m]$, taking $\pi = \pi^*$ in Proposition 1 gives

$$\begin{aligned} & V_1^{\pi_{k+1}}(\rho) + \sum_{i=2}^m [\lambda_{k,i} + \eta'(b_i - V_i^{\pi_k}(\rho))] \cdot [V_i^{\pi_{k+1}}(\rho) - b_i] - \alpha \frac{D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1} || \pi_k)}{1-\gamma} \\ & \geq V_1^{\pi^*}(\rho) - \alpha \frac{D_{d_\rho^{\pi^*}}(\pi^* || \pi_k) - D_{d_\rho^{\pi^*}}(\pi^* || \pi_{k+1})}{1-\gamma} - \epsilon_k. \end{aligned} \quad (29)$$

Denote $\delta_{k,i} := b_i - V_i^{\pi_k}(\rho)$ as the constraint violation for the i -th constraint at macro step k . We thus have

$$[\lambda_{k,i} + \eta'(b_i - V_i^{\pi_k}(\rho))] \cdot (V_i^{\pi_{k+1}}(\rho) - b_i) = -\lambda_{k,i} \delta_{k+1,i} - \eta' \delta_{k,i} \delta_{k+1,i}.$$

We can then bound this two terms respectively.

- $\lambda_{k,i} \delta_{k+1,i}$: Note that $\lambda_{k+1,i} = \max\{-\eta' \delta_{k+1,i}, \lambda_{k,i} + \eta' \delta_{k+1,i}\}$.
 If $\lambda_{k+1,i} = -\eta' \delta_{k+1,i}$, then

$$\frac{1}{2} \lambda_{k+1,i}^2 - \frac{1}{2} \lambda_{k,i}^2 - \eta'^2 \delta_{k+1,i}^2 = -\frac{1}{2} \lambda_{k,i}^2 - \frac{\eta'^2}{2} \delta_{k+1,i}^2 \leq \eta' \lambda_{k,i} \delta_{k+1,i},$$

which implies $-\lambda_{k,i} \delta_{k+1,i} \leq \frac{\lambda_{k,i}^2 - \lambda_{k+1,i}^2}{2\eta'} + \eta' \delta_{k+1,i}^2$.

If $\lambda_{k+1,i} = \lambda_{k,i} + \eta' \delta_{k+1,i}$, then

$$\eta' \lambda_{k,i} \delta_{k+1,i} = \frac{1}{2}(\lambda_{k,i} + \eta' \delta_{k+1,i})^2 - \frac{1}{2}\lambda_{k,i}^2 - \frac{\eta'^2}{2}\delta_{k+1,i}^2 \geq \frac{1}{2}\lambda_{k+1,i}^2 - \frac{1}{2}\lambda_{k,i}^2 - \eta'^2\delta_{k+1,i}^2,$$

which also implies $-\lambda_{k,i}\delta_{k+1,i} \leq \frac{\lambda_{k,i}^2 - \lambda_{k+1,i}^2}{2\eta'} + \eta'\delta_{k+1,i}^2$.

- $\eta' \delta_{k,i} \delta_{k+1,i}$: Note that $\eta' \delta_{k,i} \delta_{k+1,i} = \frac{\eta'}{2}\delta_{k,i}^2 + \frac{\eta'}{2}\delta_{k+1,i}^2 - \frac{\eta'}{2}(\delta_{k,i} - \delta_{k+1,i})^2$, and $\frac{\eta'}{2}(\delta_{k,i} - \delta_{k+1,i})^2 \leq \frac{\gamma^2 \eta'}{(1-\gamma)^4} D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1} || \pi_k)$. We thus have $-\eta' \delta_{k,i} \delta_{k+1,i} \leq -\frac{\eta'}{2}(\delta_{k,i}^2 + \delta_{k+1,i}^2) + \frac{\gamma^2 \eta'}{(1-\gamma)^4} D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1} || \pi_k)$.

Substituting the above upper bounds into (29) leads to

$$\begin{aligned} V_1^{\pi_{k+1}}(\rho) &+ \frac{\|\lambda_k\|_2^2 - \|\lambda_{k+1}\|_2^2}{2\eta'} + \eta' \frac{\|\delta_{k+1}\|_2^2 - \|\delta_k\|_2^2}{2} + \left(\frac{\eta' \gamma^2 m}{(1-\gamma)^4} - \frac{\alpha}{1-\gamma} \right) D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1} || \pi_k) \\ &\geq V_1^{\pi^*}(\rho) - \alpha \frac{D_{d_\rho^{\pi^*}}(\pi^* || \pi_k) - D_{d_\rho^{\pi^*}}(\pi^* || \pi_{k+1})}{1-\gamma} - \epsilon_k. \end{aligned}$$

Recall $\alpha \geq \frac{2\eta' m}{(1-\gamma)^3}$, it then follows from telescoping that

$$\begin{aligned} \sum_{k=1}^K V_1^{\pi_k}(\rho) &\geq K V_1^{\pi^*}(\rho) - \alpha \frac{D_{d_\rho^{\pi^*}}(\pi^* || \pi_0) - D_{d_\rho^{\pi^*}}(\pi^* || \pi_K)}{1-\gamma} - \sum_{k=0}^{K-1} \epsilon_k \\ &\quad + \eta' \frac{\|\delta_0\|_2^2 - \|\delta_K\|_2^2}{2} + \frac{\|\lambda_K\|_2^2 - \|\lambda_0\|_2^2}{2\eta'} \end{aligned} \quad (30)$$

$$\begin{aligned} &= K V_1^{\pi^*}(\rho) - \alpha \frac{D_{d_\rho^{\pi^*}}(\pi^* || \pi_0) - D_{d_\rho^{\pi^*}}(\pi^* || \pi_K)}{1-\gamma} - \sum_{k=0}^{K-1} \epsilon_k \\ &\quad + \left(\frac{\|\lambda_K\|_2^2}{2\eta'} - \eta' \frac{\|\delta_K\|_2^2}{2} \right) + \eta' \frac{\|\delta_0\|_2^2 - \|\lambda_0\|_2^2}{2} - \frac{1/\eta' - \eta'}{2} \|\lambda_0\|_2^2 \end{aligned} \quad (31)$$

$$\begin{aligned} &\stackrel{(a)}{\geq} K V_1^{\pi^*}(\rho) - \alpha \frac{D_{d_\rho^{\pi^*}}(\pi^* || \pi_0) - D_{d_\rho^{\pi^*}}(\pi^* || \pi_K)}{1-\gamma} - \sum_{k=0}^{K-1} \epsilon_k - \frac{1/\eta' - \eta'}{2} \|\lambda_0\|_2^2 \\ &\stackrel{(b)}{\geq} K V_1^{\pi^*}(\rho) - \frac{3\alpha \log(|\mathcal{A}|)}{1-\gamma}. \end{aligned} \quad (32)$$

(a) holds due to the third property of Lemma 6, and (b) holds since π_0 is the uniformly distributed policy. Thus $D_{d_\rho^{\pi^*}}(\pi^* || \pi_0) = \sum_{s \in \mathcal{S}} d_\rho^{\pi^*}(s) \sum_{a \in \mathcal{A}} \pi^*(a|s) \log(|\mathcal{A}| \pi^*(a|s)) \leq \log(|\mathcal{A}|)$, $\sum_{k=0}^{K-1} \epsilon_k = \frac{\alpha \log(|\mathcal{A}|)}{1-\gamma}$, and $\lambda_{0,i}^2 = \eta'^2 [\delta_{0,i}]_+^2$ implying $\frac{1/\eta' - \eta'}{2} \|\lambda_0\|_2^2 \leq \frac{(\eta' - \eta'^3) \|\delta_0\|_2^2}{2} \leq \frac{\eta'}{2(1-\gamma)^2} \leq \frac{\alpha \log(|\mathcal{A}|)}{1-\gamma}$. We now obtain the bound (11), after dividing by K on both sides.

E.2.2 Violation gap of ARNPG-EPD

Recall that $\delta_{k,i} := b_i - V_i^{\pi_k}(\rho)$ is the constraint violation for the i -th constraint at macro step k . We aim to provide an upper bound on $\sum_{k=1}^K \delta_{k,i}$ to control the constraint violation.

For any $i \in [2 : m]$, since $\lambda_{k,i} = \max\{-\eta' \delta_{k,i}, \lambda_{k-1,i} + \eta' \delta_{k,i}\} \geq \lambda_{k-1,i} + \eta' \delta_{k,i}$, we have

$$\sum_{k=1}^K \delta_{k,i} \leq \frac{\lambda_{K,i} - \lambda_{0,i}}{\eta'} \leq \frac{\lambda_{K,i}}{\eta'} \leq \frac{\|\lambda_K\|_2}{\eta'} \leq \frac{\|\lambda^*\|_2 + \|\lambda_K - \lambda^*\|_2}{\eta'}. \quad (33)$$

To upper bound the constraint violation, it therefore suffices to bound the dual variables.

Consider the Lagrangian with optimal dual variable $\mathcal{L}(\pi, \lambda^*) = V_1^\pi(\rho) + \sum_{i=2}^m \lambda_i^*(V_i^\pi(\rho) - b_i)$, whose maximum value $V_1^{\pi^*}(\rho)$ is achieved by the optimal policy π^* . We know

$$\begin{aligned}
KV_1^{\pi^*}(\rho) &\stackrel{(a)}{=} K\mathcal{L}(\pi^*, \lambda^*) \geq \sum_{k=1}^K \mathcal{L}(\pi_k, \lambda^*) = \sum_{k=1}^K V_1^{\pi_k}(\rho) + \sum_{i=2}^m \lambda_i^* \sum_{k=1}^K (V_i^{\pi_k}(\rho) - b_i) \\
&= \sum_{k=1}^K V_1^{\pi_k}(\rho) - \sum_{i=2}^m \lambda_i^* \sum_{k=1}^K \delta_{k,i} \stackrel{(b)}{\geq} \sum_{k=1}^K V_1^{\pi_k}(\rho) - \frac{1}{\eta'} \sum_{i=2}^m \lambda_i^* \lambda_{K,i} \\
&\stackrel{(c)}{\geq} KV_1^{\pi^*}(\rho) - \alpha \frac{D_{d_\rho^*}(\pi^* || \pi_0) - D_{d_\rho^*}(\pi^* || \pi_K)}{1 - \gamma} + \frac{\|\lambda_K\|^2}{2\eta'} - \frac{\eta' \|\delta_K\|^2}{2} - \frac{\lambda_i^* \sum_{i=2}^m \lambda_{K,i}}{\eta'} - \Delta_K \\
&\geq KV_1^{\pi^*}(\rho) - \frac{\alpha \log(|\mathcal{A}|)}{1 - \gamma} + \frac{\alpha D_{d_\rho^*}(\pi^* || \pi_K)}{1 - \gamma} + \frac{\|\lambda_K\|_2^2}{2\eta'} - \frac{\eta' \|\delta_K\|_2^2}{2} - \frac{\lambda_i^* \sum_{i=2}^m \lambda_{K,i}}{\eta'} - \Delta_K,
\end{aligned}$$

where $\Delta_K := \frac{2\alpha \log(|\mathcal{A}|)}{1 - \gamma} \geq \sum_{k=0}^{K-1} \epsilon_k + \frac{1/\eta' - \eta'}{2} \|\lambda_0\|_2^2$. Then (a) holds due to complementary slackness $\lambda_i^*(V_i^{\pi^*}(\rho) - b_i) = 0$, (b) follows from (33), and (c) follows from (31) and the third property of Lemma 6. It then follows that

$$\frac{\|\lambda_K\|_2^2}{2\eta'} - \frac{\lambda_i^* \sum_{i=2}^m \lambda_{K,i}}{\eta'} \leq \frac{\alpha \log(|\mathcal{A}|)}{1 - \gamma} - \frac{\alpha D_{d_\rho^*}(\pi^* || \pi_K)}{1 - \gamma} + \frac{\eta' \|\delta_K\|_2^2}{2} + \Delta_K. \quad (34)$$

Denoting $\delta_i^* := b_i - V_i^{\pi^*}(\rho) \leq 0$, according to Lemma 3, we have

$$\frac{\alpha D_{d_\rho^*}(\pi^* || \pi_K)}{1 - \gamma} \geq \frac{(1 - \gamma)^3 \alpha}{2\gamma^2} \|\delta_K - \delta^*\|_\infty^2 \geq \frac{(1 - \gamma)^3 \alpha}{2\gamma^2 m} \|\delta_K - \delta^*\|_2^2. \quad (35)$$

We can also obtain

$$\begin{aligned}
-\frac{(1 - \gamma)^3 \alpha}{2\gamma^2 m} \|\delta_K - \delta^*\|_2^2 + \frac{\eta'}{2} \|\delta_K\|_2^2 &= \left(\frac{\eta'}{2} - \frac{\gamma^2 m \eta'^2}{2[\gamma^2 m \eta' - (1 - \gamma)^3 \alpha]} \right) \|\delta^*\|^2 \\
&\quad + \frac{\gamma^2 m \eta' - (1 - \gamma)^3 \alpha}{2\gamma^2 m} \left\| \delta_K - \delta^* + \frac{\gamma^2 m \eta'}{\gamma^2 m \eta' - (1 - \gamma)^3 \alpha} \delta^* \right\|^2,
\end{aligned} \quad (36)$$

by substituting $a = \frac{(1 - \gamma)^3 \alpha}{2\gamma^2 m}$, $b = \frac{\eta'}{2}$, $x = \delta_K - \delta^*$, $y = \delta^*$ into the binomial equation

$$-a\|x\|_2^2 + b\|x + y\|_2^2 = (b - \frac{b^2}{b - a})\|y\|_2^2 + (b - a)\|x\|_2^2 + \frac{b}{b - a}y\|_2^2.$$

Recalling $\alpha \geq \frac{2\eta' m}{(1 - \gamma)^3}$, we can verify that $\frac{\gamma^2 m \eta' - (1 - \gamma)^3 \alpha}{2\gamma^2 m} \leq 0$ and $\frac{\eta'}{2} - \frac{\gamma^2 m \eta'^2}{2[\gamma^2 m \eta' - (1 - \gamma)^3 \alpha]} \leq \eta'$. It follows that

$$-\frac{(1 - \gamma)^3 \alpha}{2\gamma^2 m} \|\delta_K - \delta^*\|_2^2 + \frac{\eta'}{2} \|\delta_K\|_2^2 \leq \eta' \|\delta^*\|_2^2. \quad (37)$$

Substituting (35) and (37) into (34) gives

$$\begin{aligned}
\frac{1}{2\eta'} \|\lambda_K - \lambda^*\|_2^2 &= \frac{1}{2\eta'} \|\lambda^*\|^2 + \frac{1}{2\eta'} \|\lambda_K\|^2 - \frac{1}{\eta'} \sum_{i=1}^m \lambda_i^* \lambda_{K,i} \\
&\leq \frac{1}{2\eta'} \|\lambda^*\|_2^2 + \frac{\alpha \log(|\mathcal{A}|)}{1 - \gamma} + \Delta_K + \eta' \|\delta^*\|^2 \\
&\leq \frac{1}{2\eta'} \|\lambda^*\|_2^2 + \frac{3\alpha \log(|\mathcal{A}|)}{1 - \gamma} + \frac{\eta'(m - 1)}{(1 - \gamma)^2} \\
&\leq \frac{1}{2\eta'} \|\lambda^*\|_2^2 + \frac{4\alpha \log(|\mathcal{A}|)}{1 - \gamma},
\end{aligned} \quad (38)$$

where the last inequality follows from $\frac{\eta'(m-1)}{(1-\gamma)^2} \leq \frac{\alpha \log(|\mathcal{A}|)}{1-\gamma}$. Using the above bound in (33), we get

$$\begin{aligned} \sum_{k=1}^K \delta_{k,i} &\leq \frac{\|\lambda^*\|_2}{\eta'} + \frac{\|\lambda_K - \lambda^*\|_2}{\eta'} \leq \frac{\|\lambda^*\|_2}{\eta'} + \sqrt{\frac{\|\lambda^*\|_2^2}{\eta'^2} + \frac{8\alpha \log(|\mathcal{A}|)}{(1-\gamma)\eta'}} \\ &\leq 2\frac{\|\lambda^*\|_2}{\eta'} + 3\sqrt{\frac{\alpha \log(|\mathcal{A}|)}{(1-\gamma)\eta'}}, \end{aligned} \quad (39)$$

from which we the constraint violation upper bound given in (12) follows.

Proof of Theorem 2. We can conclude Theorem 2 from the above discussion on the optimality gap and the constraint violation. \square

Proof of Corollary 2. Note that the number of iterations in the inner loop depends on the value of dual variables, i.e., $t_k = \lceil \frac{1}{1-\gamma} \log(\frac{5L_k K}{2\eta' m \log(|\mathcal{A}|)}) + 1 \rceil$ with $L_k = 1 + \frac{\eta'(m-1)}{1-\gamma} + \sum_{i=2}^m \lambda_{k,i}$. It is easy to verify that

$$\frac{1}{2\eta'} \|\lambda_k - \lambda^*\|_2^2 \leq \frac{1}{2\eta'} \|\lambda^*\|_2^2 + \frac{4\alpha \log(|\mathcal{A}|)}{1-\gamma}$$

in the same manner as the proof of inequality (38). It then follows that

$$\begin{aligned} \sum_{i=2}^m \lambda_{k,i} &= \|\lambda_k\|_1 \leq \sqrt{m} \|\lambda_k\|_2 \leq \sqrt{m} (\|\lambda_k - \lambda^*\| + \|\lambda^*\|) \\ &\leq \sqrt{2m \|\lambda^*\|_2^2 + \frac{8\eta' m \alpha \log(|\mathcal{A}|)}{1-\gamma}} = O\left(\sqrt{m} \|\lambda^*\|_2 + \frac{m \log(|\mathcal{A}|)}{(1-\gamma)^2}\right). \end{aligned}$$

We then have $t_k = \Theta\left(\frac{1}{1-\gamma} \log(K)\right)$, and $T = \sum_{k=0}^{K-1} t_k = \Theta\left(\frac{K}{1-\gamma} \log(K)\right)$. We conclude the proof by $\frac{K}{1-\gamma} = \Theta(T/\log(T))$. \square

E.3 ARNPG-OMDA for max-min trade-off

E.3.1 Smoothness property

Define $\mathcal{X} := \mathcal{V}_\rho \times \Delta([m]) \subset \mathbb{R}^{2m}$. Define a norm Ψ on \mathbb{R}^{2m} by $\Psi(v, \lambda) = \|v\|_\infty + \|\lambda\|_1$. Its dual norm is $\Psi^*(v, \lambda) = \|v\|_1 + \|\lambda\|_\infty$.

Define $G^{v, -\lambda}(X) := (\nabla_v \Phi(X), -\nabla_\lambda \Phi(X))$ for $X \in \mathcal{X}$. Assume the function Φ is β -smooth w.r.t. the Ψ -norm over its domain \mathcal{X} , i.e.,

$$\Psi^*(G^{v, -\lambda}(X) - G^{v, -\lambda}(X')) \leq \beta \Psi(X - X'), \quad \forall X, X' \in \mathcal{X}. \quad (40)$$

Define E_k , which will be an auxiliary term for the convergence analysis, as follows:

$$\begin{aligned} E_k &:= \langle \tilde{G}_k^v - \tilde{G}_{k+1}^v, V_{1:m}^{\pi_{k+1}}(\rho) - V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) \rangle + \alpha \frac{D_{d_\rho^{\tilde{\pi}_{k+1}}}(\tilde{\pi}_{k+1} \|\pi_k) + D_{d_\rho^\pi}(\pi_{k+1} \|\tilde{\pi}_{k+1})}{1-\gamma} \\ &\quad + \langle \tilde{G}_k^\lambda - \tilde{G}_{k+1}^\lambda, \tilde{\lambda}_{k+1} - \lambda_{k+1} \rangle + \frac{D(\lambda_{k+1} \|\tilde{\lambda}_{k+1}) + D(\tilde{\lambda}_{k+1} \|\lambda_k)}{\eta'}. \end{aligned} \quad (41)$$

Lemma 7 (Technical lemma for smoothness). *When $\alpha \geq \frac{6\beta}{(1-\gamma)^4}$ and $\eta' \leq \frac{1}{6\beta}$, $\sum_{k=0}^{K-1} E_k \geq 0$.*

Proof of Lemma 7. Recall the definition of E_k (41). Let $X_k := (V_{1:m}^{\pi_k}(\rho), \lambda_k) \in \mathcal{X}$ and $\tilde{X}_k := (V_{1:m}^{\tilde{\pi}_k}(\rho), \tilde{\lambda}_k) \in \mathcal{X}$; $G_k^{v, -\lambda} := G^{v, -\lambda}(X_k)$ and $\tilde{G}_k^{v, -\lambda} := G^{v, -\lambda}(\tilde{X}_k)$. We can then rewrite E_k as

$$\begin{aligned} E_k &= \langle \tilde{G}_k^{v, -\lambda} - \tilde{G}_{k+1}^{v, -\lambda}, X_{k+1} - \tilde{X}_{k+1} \rangle + \alpha \frac{D_{d_\rho^{\tilde{\pi}_{k+1}}}(\tilde{\pi}_{k+1} \|\pi_k) + D_{d_\rho^\pi}(\pi_{k+1} \|\tilde{\pi}_{k+1})}{1-\gamma} \\ &\quad + \frac{D(\lambda_{k+1} \|\tilde{\lambda}_{k+1}) + D(\tilde{\lambda}_{k+1} \|\lambda_k)}{\eta'}. \end{aligned}$$

We can obtain

$$\begin{aligned}
& \langle \tilde{G}_{k+1}^{v,-\lambda} - \tilde{G}_k^{v,-\lambda}, X_{k+1} - \tilde{X}_{k+1} \rangle \stackrel{(a)}{\leq} \Psi^*(\tilde{G}_{k+1}^{v,-\lambda} - \tilde{G}_k^{v,-\lambda}) \Psi(X_{k+1} - \tilde{X}_{k+1}) \\
& \stackrel{(b)}{\leq} \Psi^*(\tilde{G}_{k+1}^{v,-\lambda} - G_k^{v,-\lambda}) \Psi(X_{k+1} - \tilde{X}_{k+1}) + \Psi^*(G_k^{v,-\lambda} - \tilde{G}_k^{v,-\lambda}) \Psi(X_{k+1} - \tilde{X}_{k+1}) \\
& \stackrel{(c)}{\leq} \beta \Psi(\tilde{X}_{k+1} - X_k) \Psi(X_{k+1} - \tilde{X}_{k+1}) + \beta \Psi(X_k - \tilde{X}_k) \Psi(X_{k+1} - \tilde{X}_{k+1}) \\
& \stackrel{(d)}{\leq} \frac{\beta}{\sqrt{8}-2} \Psi(\tilde{X}_{k+1} - X_k)^2 + \left(\frac{\beta}{\sqrt{8}+2} + \frac{\beta}{2} \right) \Psi(X_{k+1} - \tilde{X}_{k+1})^2 + \frac{\beta}{2} \Psi(X_k - \tilde{X}_k)^2.
\end{aligned}$$

Inequality (a) follows from the Cauchy-Schwarz inequality for the Ψ -norm; (b) from the triangle inequality; (c) from the smoothness of function Φ defined in (40); and (d) from $ac + bc \leq \frac{a^2}{\sqrt{8}-2} + \frac{c^2}{\sqrt{8}+2} + \frac{b^2}{2} + \frac{c^2}{2}$.

Since $X_0 = \tilde{X}_0$, $\frac{1}{\sqrt{8}+2} + \frac{1}{2} + \frac{1}{2} = \frac{1}{\sqrt{8}-2}$, and $\Psi(v, \lambda)^2 \leq 2\|v\|_\infty^2 + 2\|\lambda\|_1^2$, we have

$$\begin{aligned}
& \sum_{k=0}^{K-1} \langle \tilde{G}_{k+1}^{v,-\lambda} - \tilde{G}_k^{v,-\lambda}, X_{k+1} - \tilde{X}_{k+1} \rangle \\
& \leq \frac{\beta}{\sqrt{8}-2} \sum_{k=0}^{K-1} \Psi(\tilde{X}_{k+1} - X_k)^2 + \frac{\beta}{\sqrt{8}-2} \sum_{k=1}^K \Psi(X_k - \tilde{X}_k)^2 \\
& \leq \frac{2\beta}{\sqrt{8}-2} \sum_{k=0}^{K-1} \left(\|V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) - V_{1:m}^{\pi_k}(\rho)\|_\infty^2 + \|\tilde{\lambda}_{k+1} - \lambda_k\|_1^2 \right. \\
& \quad \left. + \|V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) - V_{1:m}^{\pi_{k+1}}(\rho)\|_\infty^2 + \|\tilde{\lambda}_{k+1} - \lambda_{k+1}\|_1^2 \right).
\end{aligned}$$

Noting that $\frac{2\beta}{\sqrt{8}-2} \leq 3\beta$, by Lemma 3 we have

$$\frac{2\beta}{\sqrt{8}-2} \|V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) - V_{1:m}^{\pi_k}(\rho)\|_\infty^2 \leq \frac{6\gamma^2\beta}{(1-\gamma)^4} D_{d_\rho^{\tilde{\pi}_{k+1}}}(\tilde{\pi}_{k+1} \|\pi_k).$$

By Pinsker's inequality, we have

$$\frac{2\beta}{\sqrt{8}-2} \|\tilde{\lambda}_{k+1} - \lambda_{k+1}\|_1^2 \leq 6\beta D(\lambda_{k+1} \|\tilde{\lambda}_{k+1}).$$

Since $\alpha \geq \frac{6\beta}{(1-\gamma)^4}$ and $\eta' \leq \frac{1}{6\beta}$, we conclude that $\sum_{k=0}^{K-1} E_k \geq 0$. \square

E.3.2 Convergence of ARNPG-OMDA

Proof of Theorem 3. By $|\tilde{r}_k(s, a)| = |\langle \tilde{G}_k^v, r_{1:m}(s, a) \rangle| \leq \|\tilde{G}_k^v\|_1 \|r_{1:m}(s, a)\|_\infty \leq L$, we know $\|\tilde{r}_k\|_\infty \leq L$. Taking $\epsilon_k = \frac{\alpha \log(|\mathcal{A}|)}{(1-\gamma)K}$, we choose $t_k = \lceil \frac{1}{1-\gamma} \log(\frac{5LK}{6\beta \log(|\mathcal{A}|)}) + 1 \rceil$.

Then by Proposition 1, for any policy π , we have two fundamental inequalities for the updates $\tilde{\pi}_{k+1}$ and π_{k+1} respectively:

$$\begin{aligned}
V_{\tilde{r}_k}^{\tilde{\pi}_{k+1}}(\rho) - \alpha \frac{D_{d_\rho^{\tilde{\pi}_{k+1}}}(\tilde{\pi}_{k+1} \|\pi_k)}{1-\gamma} & \geq V_{\tilde{r}_k}^\pi(\rho) - \alpha \frac{D_{d_\rho^\pi}(\pi \|\pi_k) - D_{d_\rho^\pi}(\pi \|\tilde{\pi}_{k+1})}{1-\gamma} - \epsilon_k, \\
V_{\tilde{r}_{k+1}}^{\pi_{k+1}}(\rho) - \alpha \frac{D_{d_\rho^{\pi_{k+1}}}(\pi_{k+1} \|\pi_k)}{1-\gamma} & \geq V_{\tilde{r}_{k+1}}^\pi(\rho) - \alpha \frac{D_{d_\rho^\pi}(\pi \|\pi_k) - D_{d_\rho^\pi}(\pi \|\pi_{k+1})}{1-\gamma} - \epsilon_k.
\end{aligned}$$

Note that $V_{\tilde{r}_k}^\pi(\rho) = \langle \tilde{G}_k^v, V_{1:m}^\pi(\rho) \rangle$. Taking $\pi = \pi_{k+1}$ in the first inequality, and summing two inequalities gives

$$\begin{aligned}
\langle \tilde{G}_{k+1}^v, V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) - V_{1:m}^{\pi_{k+1}}(\rho) \rangle & \geq \alpha \frac{D_{d_\rho^\pi}(\pi \|\pi_{k+1}) - D_{d_\rho^\pi}(\pi \|\pi_k)}{1-\gamma} - 2\epsilon_k \\
& + \langle \tilde{G}_k^v - \tilde{G}_{k+1}^v, V_{1:m}^{\pi_{k+1}}(\rho) - V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) \rangle + \alpha \frac{D_{d_\rho^{\tilde{\pi}_{k+1}}}(\tilde{\pi}_{k+1} \|\pi_k) + D_{d_\rho^\pi}(\pi_{k+1} \|\tilde{\pi}_{k+1})}{1-\gamma}.
\end{aligned} \tag{42}$$

We can similarly get the inequality for λ that

$$\langle \tilde{G}_k^\lambda, \tilde{\lambda}_{k+1} \rangle + \frac{D(\tilde{\lambda}_{k+1} || \lambda_k)}{\eta'} \leq \langle \tilde{G}_k^\lambda, \lambda \rangle + \frac{D(\lambda || \lambda_k) - D(\lambda || \tilde{\lambda}_{k+1})}{\eta'}, \quad (43)$$

$$\langle \tilde{G}_{k+1}^\lambda, \lambda_{k+1} \rangle + \frac{D(\lambda_{k+1} || \lambda_k)}{\eta'} \leq \langle \tilde{G}_{k+1}^\lambda, \lambda \rangle + \frac{D(\lambda || \lambda_k) - D(\lambda || \lambda_{k+1})}{\eta'}. \quad (44)$$

Taking $\lambda = \lambda_{k+1}$ in the first inequality and summing two inequalities gives

$$\begin{aligned} \langle \tilde{G}_{k+1}^\lambda, \lambda - \tilde{\lambda}_{k+1} \rangle &\geq \frac{D(\lambda || \lambda_{k+1}) - D(\lambda || \lambda_k)}{\eta'} \\ &+ \langle \tilde{G}_k^\lambda - \tilde{G}_{k+1}^\lambda, \tilde{\lambda}_{k+1} - \lambda_{k+1} \rangle + \frac{D(\lambda_{k+1} || \tilde{\lambda}_{k+1}) + D(\tilde{\lambda}_{k+1} || \lambda_k)}{\eta'}. \end{aligned} \quad (45)$$

Recall the definition of E_k in (41) that

$$\begin{aligned} E_k &= \langle \tilde{G}_k^v - \tilde{G}_{k+1}^v, V_{1:m}^{\pi_{k+1}}(\rho) - V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) \rangle + \alpha \frac{D_{d_\rho^{\tilde{\pi}_{k+1}}}(\tilde{\pi}_{k+1} || \pi_k) + D_{d_\rho^\pi}(\pi_{k+1} || \tilde{\pi}_{k+1})}{1 - \gamma} \\ &+ \langle \tilde{G}_k^\lambda - \tilde{G}_{k+1}^\lambda, \tilde{\lambda}_{k+1} - \lambda_{k+1} \rangle + \frac{D(\lambda_{k+1} || \tilde{\lambda}_{k+1}) + D(\tilde{\lambda}_{k+1} || \lambda_k)}{\eta'}. \end{aligned}$$

We then have

$$\begin{aligned} &- \Phi(V_{1:m}^\pi(\rho), \tilde{\lambda}_{k+1}) + \Phi(V_{1:m}^{\tilde{\pi}_{k+1}}(\rho), \lambda) \\ &= \Phi(V_{1:m}^{\tilde{\pi}_{k+1}}(\rho), \tilde{\lambda}_{k+1}) - \Phi(V_{1:m}^\pi(\rho), \tilde{\lambda}_{k+1}) + \Phi(V_{1:m}^{\tilde{\pi}_{k+1}}(\rho), \lambda) - \Phi(V_{1:m}^{\tilde{\pi}_{k+1}}(\rho), \tilde{\lambda}_{k+1}) \\ &\stackrel{(a)}{\geq} \langle \tilde{G}_{k+1}^v, V_{1:m}^{\tilde{\pi}_{k+1}}(\rho) - V_{1:m}^\pi(\rho) \rangle + \langle \tilde{G}_{k+1}^\lambda, \lambda - \tilde{\lambda}_{k+1} \rangle \\ &\stackrel{(b)}{\geq} \alpha \frac{D_{d_\rho^\pi}(\pi || \pi_{k+1}) - D_{d_\rho^\pi}(\pi || \pi_k)}{1 - \gamma} + \frac{D(\lambda || \lambda_{k+1}) - D(\lambda || \lambda_k)}{\eta'} - 2\epsilon_k + E_k. \end{aligned}$$

Inequality (a) is by the concavity of $\Phi(\cdot, \tilde{\lambda}_{k+1})$ and convexity of $\Phi(V_{1:m}^{\tilde{\pi}_{k+1}}(\rho), \cdot)$. Inequality (b) is based on combining (42) and (45).

Taking $\pi = \pi^*$ and $\lambda = \arg \min_{\lambda' \in \Lambda} \Phi\left(\frac{1}{K} \sum_{k=1}^K V_{1:m}^{\tilde{\pi}_k}(\rho), \lambda'\right)$, we have

$$\begin{aligned} F\left(\frac{1}{K} \sum_{k=1}^K V_{1:m}^{\tilde{\pi}_k}(\rho)\right) &= \Phi\left(\frac{1}{K} \sum_{k=1}^K V_{1:m}^{\tilde{\pi}_k}(\rho), \lambda\right) \geq \frac{1}{K} \sum_{k=1}^K \Phi(V_{1:m}^{\tilde{\pi}_k}(\rho), \lambda) \\ &\geq \frac{1}{K} \sum_{k=0}^{K-1} \Phi(V_{1:m}^{\pi^*}(\rho), \tilde{\lambda}_{k+1}) + \alpha \frac{D_{d_\rho^{\pi^*}}(\pi^* || \pi_K) - D_{d_\rho^{\pi^*}}(\pi^* || \pi_0)}{(1 - \gamma)K} + \frac{D(\lambda || \lambda_K) - D(\lambda || \lambda_0)}{\eta' K} \\ &\quad - \frac{2}{K} \sum_{k=0}^{K-1} \epsilon_k + \frac{1}{K} \sum_{k=0}^{K-1} E_k \\ &\stackrel{(a)}{\geq} F(V_{1:m}^{\pi^*}(\rho)) - \frac{3\alpha \log(|\mathcal{A}|)}{(1 - \gamma)K} - \frac{\log(m)}{\eta' K}. \end{aligned}$$

Inequality (a) is due to $D_{d_\rho^{\pi^*}}(\pi^* || \pi_0) \leq \log(|\mathcal{A}|)$ and Lemma 7. \square

Proof of Corollary 3. Note that $T = \sum_{k=0}^{K-1} t_k = \Theta(\frac{K}{1-\gamma} \log(K))$. It implies $\frac{K}{1-\gamma} = \Theta(T / \log(T))$. Substituting this into Theorem 3 concludes Corollary 3. \square

F More related works

There are in general two scenarios when considering problems of multi-objective Markov decision processes (cf. survey [23]): single-policy scenario and multi-policy scenario. This work focuses on the single-policy scenario, which we will simply refer to as multi-objective MDP, and we relegate the discussion of the multi-policy scenario to the end of this section.

Multi-objective MDP In the single-policy scenario, the agent optimizes the reward objectives according to the user’s criteria. The case of linear scalarization with known weights simply collapses to a canonical MDP. Therefore, the focus of this line of research is on nonlinear scalarization. The lexicographical ordering of objectives was considered by Wray et al. [29] and Saisubramanian et al. [24]. Moffaert et al. [26] studied a Chebyshev scalarization i.e., weighted L_∞ scalarization via a Q-learning approach. Bai et al. [5] established an $O(1/\epsilon^4)$ sample complexity for a policy-gradient method under smooth concave scalarization. Besides optimization (pure exploration), the exploration-exploitation trade-off has also been studied, where with concave scalarization, an $O(\sqrt{T})$ regret was established under the Lipschitz continuous assumption on F by [3, 30].

Constrained MDP The constrained MDP (CMDP) is viewed as a special case of a multi-objective MDP in this paper. CMDPs have attracted much attention recently. Ding et al. [11] proposed an NPG-PrimalDual algorithm that uses a primal-dual approach with NPG and showed that it can achieve $O(1/\sqrt{T})$ global convergence for both the optimality gap and the constraint violation. Xu et al. [31] proposed a primal approach called constrained-rectified policy optimization (CRPO) that updates the policy alternatively between optimizing objective and decreasing constraint violation, and enjoys the same $O(1/\sqrt{T})$ global convergence. Whether policy optimization for CMDP can achieve global convergence with a rate faster than $O(1/\sqrt{T})$ was stated as an open problem in [11]. There are two concurrent works [33, 18] tackling the problem with positive answers, with some ergodic assumptions made to facilitate the proof of $\tilde{O}(1/T)$ global convergence. Ying et al. [33] propose an NPG-aided dual approach, where the dual function is smoothed by entropy regularization in the objective function. They show an $\tilde{O}(1/T)$ convergence rate to the optimal policy of the *entropy-regularized* CMDP, but not to the true optimal policy, for which there is an established slow $O(1/\sqrt{T})$ convergence rate. They also make an additional strong assumption that the initial state distribution covers the entire state space. While such an assumption was initially used in the analysis of the global convergence of PG methods for MDPs [2, 20], it is not required when analyzing the global convergence of NPG methods [2, 10]. Moreover, this assumption does not necessarily hold for safe RL or CMDP, since the algorithm may need to avoid dangerous states even at initialization, with the consequence that the optimal policy depends on the initial state distribution. Li et al. [18] propose a primal-dual approach with an $O(\log^2(T)/T)$ convergence rate to the true optimal policy by smoothing the Lagrangian with suitable regularization on both primal and dual variables. However, they assume that the Markov chain induced by any stationary policy is ergodic in order to ensure the smoothness of the dual function. This assumption, though weaker than the assumption made by [33], will generally not hold in problems where one wishes to avoid unsafe states altogether.

Multi-policy scenario In the *multi-policy* scenario, the agent is to determine not one policy but a set of policies, so that once a set of weights of a linear scalarization is subsequently provided, a policy from the determined set can be deployed [23]. This approach essentially aims at determining the convex coverage set (CCS) of the Pareto frontier. Yang et al. [32] proposed an envelope Q-learning algorithm based on the multi-objective Bellman optimality operator to utilize the convex envelope of the solution frontier. Zhou et al. [39] studied the sample complexity under a generative model setting. Kyriakis et al. [15] utilized a policy gradient solver to search for a direction that is simultaneously an ascent direction for all objectives, utilizing a new loss function called Pareto Policy Adaptation. Wu et al. [30] considered preference-free exploration, where the agent collects samples in the exploration phase and computes a near-optimal policy for any preference-weighted reward function during the planning phase.