Sample-Efficient Constrained Reinforcement Learning with General Parameterization

Washim Uddin Mondal

Department of Electrical Engineering Indian Institute of Technology Kanpur Kanpur, UP, India 208016 wmondal@iitk.ac.in

Vaneet Aggarwal

School of IE and ECE Purdue University West Lafayette, IN, USA 47906 vaneet@purdue.edu

Abstract

We consider a constrained Markov Decision Problem (CMDP) where the goal of an agent is to maximize the expected discounted sum of rewards over an infinite horizon while ensuring that the expected discounted sum of costs exceeds a certain threshold. Building on the idea of momentum-based acceleration, we develop the Primal-Dual Accelerated Natural Policy Gradient (PD-ANPG) algorithm that ensures an ϵ global optimality gap and ϵ constraint violation with $\tilde{\mathcal{O}}((1-\gamma)^{-7}\epsilon^{-2})$ sample complexity for general parameterized policies where γ denotes the discount factor. This improves the state-of-the-art sample complexity in general parameterized CMDPs by a factor of $\mathcal{O}((1-\gamma)^{-1}\epsilon^{-2})$ and achieves the theoretical lower bound in ϵ^{-1} .

1 Introduction

Reinforcement learning (RL) is a framework where an agent repeatedly interacts with an unknown Markovian environment to find a policy that maximizes the expected discounted sum of its observed rewards. Such problems, often modeled via Markov Decision Processes (MDPs), find applications in many areas, including transportation [1], communication networks [2], robotics [3], etc. In many applications, however, the agents must also obey certain constraints. For example, in a food delivery network, the orders must be delivered within a stipulated time window; the marketing decisions of a firm must satisfy its budget constraints, etc. Such constraints are incorporated into RL by introducing a cost function. In constrained MDPs (CMDPs), the agents not only maximize the expected sum of discounted rewards but also ensure that the expected sum of discounted costs does not cross a predefined boundary.

Finding an optimal policy for a CMDP is a challenging problem, especially when the environment, i.e., the state transition function, is unknown. The efficiency of a solution to a CMDP is measured by its sample complexity, which essentially states how many state-transition samples it takes to yield a policy that is ϵ close to the optimal one while also ensuring that the expected sum of discounted costs does not violate the imposed boundary by more than ϵ amount. Many articles in the literature solve the CMDP with an unknown environment. Most of these works, however, focus on the tabular case where the number of states is finite. These solutions cannot be applied to many real-life scenarios where the state space is either large or infinite. To tackle this issue, the concept of policy parameterization must be invoked. Unfortunately, as exhibited in Table 1, only a few works are available on CMDPs with parameterized policies. While, for softmax parameterization, the state-of-the-art (SOTA) sample complexity is $\mathcal{O}(\epsilon^{-2})$, the same for the general parameterization is $\tilde{\mathcal{O}}(\epsilon^{-4})$ which is far from the lower bound $\Omega(\epsilon^{-2})$. It should be noted that the number of parameters needed in softmax parameterization is $\mathcal{O}(SA)$ where S, S are the sizes of the state and action spaces of the underlying CMDP. On the other hand, general parameterization uses S a number of parameters, which makes it appropriate

38th Conference on Neural Information Processing Systems (NeurIPS 2024).

Algorithm	Sample Complexity	Parameterization
PMD-PD [9]	$\mathcal{O}(\epsilon^{-3})$	Softmax
PD-NAC [10]	$\mathcal{O}(\epsilon^{-6})$	Softmax
NPG-PD [5]	$\mathcal{O}((1-\gamma)^{-5}\epsilon^{-2})$	Softmax
CRPO [6]	$\mathcal{O}((1-\gamma)^{-7}\epsilon^{-4})$	Softmax
NPG-PD [5]	$\mathcal{O}((1-\gamma)^{-8}\epsilon^{-6})$	General
CRPO [6]	$\mathcal{O}((1-\gamma)^{-13}\epsilon^{-6})$	General
C-NPG-PDA [4]	$\tilde{\mathcal{O}}((1-\gamma)^{-8}\epsilon^{-4})^1$	General
PD-ANPG (This Work)	$\tilde{\mathcal{O}}((1-\gamma)^{-7}\epsilon^{-2})$	General
Lower Bound [11]	$\Omega((1-\gamma)^{-5}\epsilon^{-2})$	_

Table 1: Summary of sample complexity results on CMDP with parameterized policies. The parameter γ indicates the discount factor. The dependence of the sample complexities of PMD-PD and PD-NAC on γ is not depicted in [9, 10].

for dealing with large or infinite states. Given the importance of general parameterization for large state space CMDPs, the following question naturally arises: "Is it possible to solve CMDPs with general parameterization and achieve a sample complexity better than the SOTA $\tilde{\mathcal{O}}(\epsilon^{-4})$ bound?"

In this article, we provide an affirmative answer to the above question. We propose a Primal-Dual-based Accelerated Natural Policy Gradient (PD-ANPG) algorithm to solve γ -discounted CMDPs with general parameterization. We theoretically prove that PD-ANPG achieves ϵ optimality gap and ϵ constraint violation with $\tilde{\mathcal{O}}((1-\gamma)^{-7}\epsilon^{-2})$ sample complexity (Theorem 1) that improves the SOTA $\tilde{\mathcal{O}}((1-\gamma)^{-8}\epsilon^{-4})$ sample complexity result of [4]. It closes the gap between the theoretical upper and lower bounds of sample complexity in general parameterized CMDPs (in terms of ϵ^{-1}), which was an open problem for quite some time (see the results of [5], [6], [4] in Table 1).

1.1 Challenges and Key Insights

Our algorithm builds upon the idea of primal dual-based NPG [7, 4]. However, unlike the previous works, we use accelerated stochastic gradient descent (ASGD) in the inner loop to compute the estimate of the NPG. The improvement in the sample complexity results from two key observations. Firstly, we establish a global-to-local convergence lemma (Lemma 3), which dictates how the global convergence of the Lagrange function is related to the first and second-order estimation error of the NPG. Here, via careful analysis, we show that the first-order term can be written as the expected bias of the NPG estimator (i.e., the difference between the true NPG and the expectation of its estimate). Secondly, we show (Lemma 5 and its subsequent discussion) that the bias of the NPG estimate can be interpreted as the convergence error of an ASGD program with non-stochastic (i.e., deterministic) gradients. These, combined with the ASGD convergence result provided by [8], lead to a convergence result of the Lagrange function (Corollary 1).

Finally, Theorem 1 segregates the objective and constraint violation rates from Lagrange convergence. Corollary 1 shows that Lagrange convergence error is bounded by an independent function of ζ (the dual learning rate). One might, hence, be tempted to make ζ arbitrarily small. However, our analysis shows that although small ζ leads to better objective convergence, it worsens the constraint violation rate. We demonstrate how to optimally choose ζ to reach the middle ground which eventually leads us to $\tilde{\mathcal{O}}(\epsilon^{-2})$ sample complexity.

1.2 Related Works

Unconstrained RL: Many algorithms solve MDPs with exact gradients [12, 13, 14, 15, 16]. Moreover, many works use generative models to show either first-order [17, 18, 19, 20, 21] or global convergence [22, 23, 24, 25, 7, 26, 27].

 $^{^1}$ We would like to point out that the sample complexity of C-NPG-PDA reported in [4] is $\tilde{\mathcal{O}}((1-\gamma)^{-6}\epsilon^{-4})$. However, the result is erroneous. The authors have subsequently corrected their result, and the sample complexity has been modified to $\tilde{\mathcal{O}}((1-\gamma)^{-8}\epsilon^{-4})$ in the arXiv version (updated May 2024).

Constrained RL: The tabular setting is well investigated, and many model-based [28, 29, 30, 31] and model-free [6, 30, 32, 33] algorithms are available in the literature. In comparison, there are relatively fewer works on parameterized policies. Policy mirror descent-primal dual (PMD-PD) algorithm was proposed by [9] that achieves $\mathcal{O}(\epsilon^{-3})$ sample complexity for softmax policies. For the same parameterization, [10] achieved $\mathcal{O}(\epsilon^{-6})$ sample complexity via their proposed Online Primal-Dual Natural Actor-Critic Algorithm. The primal-dual Natural Policy Gradient algorithm suggested by [5] yields $\mathcal{O}(\epsilon^{-2})$ and $\mathcal{O}(\epsilon^{-6})$ sample complexities for softmax and general parametrization respectively. [6] also proposed a primal policy-based algorithm that works for both the softmax and general function approximation cases. The state of the art for the general parameterization is given by [4], where the $\tilde{\mathcal{O}}(\epsilon^{-4})$ sample complexity is obtained. This work improves upon this direction to obtain $\tilde{\mathcal{O}}(\epsilon^{-2})$ sample complexity. The comparisons are summarized in Table 1.

2 Formulation

Let us consider a constrained Markov Decision Process (CMDP) characterized by the tuple $\mathcal{M}=(\mathcal{S},\mathcal{A},r,c,P,\gamma,\rho)$ where \mathcal{S},\mathcal{A} denote the state space and the action space respectively, $r:\mathcal{S}\times\mathcal{A}\to[0,1]$ defines the reward function, $c:\mathcal{S}\times\mathcal{A}\to[-1,1]$ is the cost function, $P:\mathcal{S}\times\mathcal{A}\to\Delta(\mathcal{S})$ indicates the state transition kernel (where $\Delta(\mathcal{S})$ is the collection of all probability distributions over \mathcal{S}), γ is the discount factor, and $\rho\in\Delta(\mathcal{S})$ is the initial state distribution. Note that the state space \mathcal{S} , in our setting, can potentially be a compact set of infinite size. However, for simplicity, we assume it to be countable. The action space, \mathcal{A} , is assumed to be of finite size. The range of the cost function is chosen to be [-1,1], rather than [0,1], to ensure that the constraint in our central optimization problem (defined later in (2)) is non-trivial. A policy, $\pi:\mathcal{S}\to\Delta(\mathcal{A})$ is defined as a distribution over the action space for a given state of the environment. For a given policy, π , and a state-action pair (s,a), we define the Q-value associated with $g\in\{r,c\}$ as follows.

$$Q_g^{\pi}(s, a) \triangleq \mathbf{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t g(s_t, a_t) \middle| s_0 = s, a_0 = a \right]$$

where \mathbf{E}_{π} is the expectation computed over all π -induced trajectories $\{(s_t,a_t)\}_{t=0}^{\infty}$ where $s_{t+1} \sim P(s_t,a_t)$ and $a_t \sim \pi(s_t)$, $\forall t \in \{0,1,\cdots\}$. Similarly, the V-value associated with policy π , state s, and $g \in \{r,c\}$ is defined below.

$$V_g^{\pi}(s) \triangleq \mathbf{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t g(s_t, a_t) \middle| s_0 = s \right] = \sum_{a} \pi(a|s) Q_g^{\pi}(s, a)$$

Below we define the advantage value for a policy π , a state-action pair (s, a), and $g \in \{r, c\}$.

$$A_q^{\pi}(s, a) \triangleq Q_q^{\pi}(s, a) - V_q^{\pi}(s)$$

Define a function $J_{q,\rho}^{\pi}$, $\forall g \in \{r,c\}$ as follows.

$$J_{g,\rho}^{\pi} \triangleq \mathbf{E}_{s \sim \rho}[V_g^{\pi}(s)] = \frac{1}{1 - \gamma} \sum_{s,a} d_{\rho}^{\pi}(s) \pi(a|s) g(s,a)$$

where $d_{\rho}^{\pi} \in \Delta(\mathcal{S})$ is the state occupancy measure given by,

$$d_{\rho}^{\pi}(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} \Pr(s_{t} = s | s_{0} \sim \rho, \pi), \ \forall s \in \mathcal{S}$$

Similarly, the state-action occupancy measure is defined as,

$$\nu_o^{\pi}(s, a) = d_o^{\pi}(s)\pi(a|s), \ \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$
 (1)

Our goal is to maximize the function $J^{\pi}_{r,\rho}$ over all policies π while ensuring that $J^{\pi}_{c,\rho}$ does not lie below a predefined threshold. Without loss of generality, we can formally express this problem as,

$$\max_{\pi} J^{\pi}_{r,\rho} \quad \text{subject to: } J^{\pi}_{c,\rho} \ge 0 \tag{2}$$

If the state space, S, is large or infinite (which is the case in many application scenarios), the policies can no longer be represented in the tabular format; rather, they are indexed by a parameter, $\theta \in \Theta$. In

this paper, we assume $\Theta = \mathbb{R}^d$. Such indexing can be done via, for example, neural networks (NNs). Let $J_{g,\rho}(\theta) \triangleq J_{g,\rho}^{\pi_{\theta}}$. This allows us to redefine the constrained optimization problem as follows.

$$\max_{\theta \in \Theta} J_{r,\rho}(\theta) \quad \text{subject to: } J_{c,\rho}(\theta) \ge 0 \tag{3}$$

We assume the existence of at least one interior point solution of the above optimization. This is also known as Slater condition which can be formally expressed as follows.

Assumption 1. There exists $\bar{\theta}$ such that $J_{c,\rho}(\bar{\theta}) \geq c_{\text{slater}}$ for some $c_{\text{slater}} \in (0, \frac{1}{1-\gamma}]$.

3 Algorithm

The dual problem associated with the constraint optimization (3) can be written as follows.

$$\min_{\lambda>0} \max_{\theta \in \Theta} J_{L,\rho}(\theta,\lambda) \text{ where } J_{L,\rho}(\theta,\lambda) \triangleq J_{r,\rho}(\theta) + \lambda J_{c,\rho}(\theta)$$
 (4)

The function, $J_{\mathrm{L},\rho}(\cdot,\cdot)$ is called the Lagrange function while λ is said to be the Lagrange multiplier. The above problem can be solved by iteratively applying the following update rule $\forall k \in \{0,\cdots,K-1\}$, starting with (θ_0,λ_0) where $\theta_0 \in \Theta$ is arbitrary and $\lambda_0=0$.

$$\theta_{k+1} = \theta_k + \eta F_{\rho}(\theta_k)^{\dagger} \nabla_{\theta} J_{L,\rho}(\theta_k, \lambda_k)$$
(5)

$$\lambda_{k+1} = \mathcal{P}_{\Lambda} \left[\lambda_k - \zeta J_{c,\rho}(\theta_k) \right] \tag{6}$$

where η, ζ are learning rates, \mathcal{P}_{Λ} denotes the projection function onto the set, $\Lambda \triangleq [0, \lambda_{\max}]$, and \dagger is the Moore-Penrose pseudoinverse operator. The choice of λ_{\max} will be specified later. Note that the update rule of θ is similar to that of the standard policy gradient method except here the learning rate, η is modulated by the inverse of the Fisher matrix, $F_{\rho}(\theta)$ which is defined below.

$$F_{\rho}(\theta) \triangleq \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi_{\theta}}} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \otimes \nabla_{\theta} \log \pi_{\theta}(a|s) \right] \tag{7}$$

where \otimes indicates the outer product. Using a variation of the classical policy gradient theorem [34], one can obtain the gradient of the Lagrange function as follows.

$$\nabla_{\theta} J_{\mathcal{L},\rho}(\theta,\lambda) = \frac{1}{1-\gamma} H_{\rho}(\theta,\lambda), \text{ where } H_{\rho}(\theta,\lambda) \triangleq \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi_{\theta}}} \left[A_{\mathcal{L},\lambda}^{\pi_{\theta}}(s,a) \nabla_{\theta} \log \pi_{\theta}(a|s) \right]$$

$$\text{and } A_{\mathcal{L},\lambda}^{\pi_{\theta}}(s,a) \triangleq A_{r}^{\pi_{\theta}}(s,a) + \lambda A_{c}^{\pi_{\theta}}(s,a)$$

$$(8)$$

In most application scenarios, the learner is unaware of the state transition function, P, and thereby, of the advantage function, $A_{\mathrm{L},\lambda}^{\pi_{\theta}}$ and the occupancy measure, $\nu_{\rho}^{\pi_{\theta}}$. This makes the exact computation of $F_{\rho}(\theta)$ and $H_{\rho}(\theta,\lambda)$ an impossible task. Fortunately, there is a way to obtain an approximate value of the *natural policy gradient* $\omega_{\theta,\lambda}^* \triangleq F_{\rho}(\theta)^{\dagger} \nabla_{\theta} J_{\mathrm{L},\rho}(\theta,\lambda)$ that does not require the knowledge of P. Invoking (8), one can prove that $\omega_{\theta,\lambda}^*$ is a solution of a quadratic optimization. Formally, we have,

$$\omega_{\theta,\lambda}^* \in \arg\min_{\omega \in \mathbb{R}^d} L_{\nu_{\rho}^{\pi_{\theta}}}(\omega, \theta, \lambda),$$
where $L_{\nu_{\rho}^{\pi_{\theta}}}(\omega, \theta, \lambda) \triangleq \frac{1}{2} \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi_{\theta}}} \left[\left(\frac{1}{1 - \gamma} A_{\mathrm{L},\lambda}^{\pi_{\theta}}(s, a) - \omega^{\mathrm{T}} \nabla_{\theta} \log \pi_{\theta}(a|s) \right)^2 \right]$ (9)

The above reformulation opens up the possibility to compute $\omega_{\theta,\lambda}^*$ via a gradient descent-type iterative procedure. Observe that the gradient of $L_{\nu_{\alpha}^{\pi_{\theta}}}(\cdot,\theta,\lambda)$ can be calculated as follows.

$$\nabla_{\omega} L_{\nu_{\rho}^{\pi_{\theta}}}(\omega, \theta, \lambda) = F_{\rho}(\theta)\omega - \frac{1}{1 - \gamma} H_{\rho}(\theta, \lambda)$$
(10)

Algorithm 1 describes a procedure to obtain unbiased estimates of this gradient. This is inspired by Algorithm 3 of [12]. Additionally, observe from (6) that the update of the Lagrange variable, λ requires the computation of $J_{c,\rho}(\theta)$ which is also difficult to accomplish without having an explicit knowledge about P. Algorithm 1 also provides an unbiased estimation of the above quantity.

Algorithm 1 first samples a horizon length, T, from the geometric distribution with success probability $(1-\gamma)$ and executes the CMDP for T instances following the policy, π_{θ} , starting from a state $s_0 \sim \rho$.

Algorithm 1 Unbiased Sampling

```
1: Input: Parameters (\theta, \omega, \lambda, \gamma), Initial Distribution \rho
```

2:
$$T \sim \text{Geo}(1-\gamma)$$
, $s_0 \sim \rho$, $a_0 \sim \pi_{\theta}(s_0)$

3: **for**
$$j \in \{0, \cdots, T-1\}$$
 do

3: **for** $j \in \{0, \dots, T-1\}$ **do**4: Execute a_j , observe $s_{j+1} \sim P(s_j, a_j)$ and sample $a_{j+1} \sim \pi_{\theta}(s_{j+1})$ 5: $\hat{J}_{c,\rho}(\theta) \leftarrow \sum_{j=0}^{T} c(s_j, a_j), (\hat{s}, \hat{a}) \leftarrow (s_T, a_T)$

5:
$$\hat{J}_{c,\rho}(\theta) \leftarrow \sum_{i=0}^{T} c(s_i, a_i), (\hat{s}, \hat{a}) \leftarrow (s_T, a_T)$$

6:
$$T \sim \text{Geo}(1-\gamma), (s_0, a_0) \leftarrow (\hat{s}, \hat{a})$$

7: **for** $j \in \{0, \dots, T-1\}$ **do**

7: **for**
$$j \in \{0, \cdots, T-1\}$$
 do

Execute a_j , observe $s_{j+1} \sim P(s_j, a_j)$, and sample $a_{j+1} \sim \pi_{\theta}(s_{j+1})$

9: for
$$g \in \{r, c\}$$
 do

10:
$$\hat{Q}_g^{\pi_\theta}(\hat{s}, \hat{a}) \leftarrow \sum_{j=0}^T g(s_j, a_j)$$

11:
$$T \sim \text{Geo}(1-\gamma), s_0 \leftarrow \hat{s}, a_0 \sim \pi_{\theta}(s_0)$$
 \triangleright V-function Estimation 12: **for** $j \in \{0, \dots, T-1\}$ **do**

12: **for**
$$j \in \{0, \cdots, T-1\}$$
 do

Execute a_j , observe $s_{j+1} \sim P(s_j, a_j)$, and sample $a_{j+1} \sim \pi_{\theta}(s_{j+1})$

14: **for**
$$g \in \{r, c\}$$
 do

15:
$$\hat{V}_q^{\pi_\theta}(\hat{s}) \leftarrow \sum_{j=0}^T g(s_j, a_j)$$

15:
$$\hat{V}_{g}^{\pi_{\theta}}(\hat{s}) \leftarrow \sum_{j=0}^{T} g(s_{j}, a_{j})$$

16: $\hat{A}_{g}^{\pi_{\theta}}(\hat{s}, \hat{a}) \leftarrow \hat{Q}_{g}^{\pi_{\theta}}(\hat{s}, \hat{a}) - \hat{V}_{g}^{\pi_{\theta}}(\hat{s})$

Estimation of Relevant functions

$$\hat{A}_{L,\lambda}^{\pi_{\theta}}(\hat{s}, \hat{a}) \leftarrow \hat{A}_{r}^{\pi_{\theta}}(\hat{s}, \hat{a}) + \lambda \hat{A}_{c}^{\pi_{\theta}}(\hat{s}, \hat{a}) \tag{11}$$

$$\hat{F}_{\rho}(\theta) \leftarrow \nabla_{\theta} \log \pi_{\theta}(\hat{a}|\hat{s}) \otimes \nabla_{\theta} \log \pi_{\theta}(\hat{a}|\hat{s})$$
(12)

$$\hat{H}_{\rho}(\theta, \lambda) \leftarrow \hat{A}_{L,\lambda}^{\pi_{\theta}}(\hat{s}, \hat{a}) \nabla_{\theta} \log \pi_{\theta}(\hat{a}|\hat{s})$$
(13)

Gradient Estimate 18:

$$\hat{\nabla}_{\omega} L_{\nu_{\rho}^{\pi_{\theta}}}(\omega, \theta, \lambda) \leftarrow \hat{F}_{\rho}(\theta)\omega - \frac{1}{1 - \gamma} \hat{H}_{\rho}(\theta, \lambda) \tag{14}$$

19: **Output:**
$$\hat{J}_{c,\rho}(\theta), \hat{\nabla}_{\omega} L_{\nu_{\rho}^{\pi_{\theta}}}(\omega, \theta, \lambda)$$

The total cost observed in the resulting trajectory is assigned as the estimate $\hat{J}_{c,\rho}(\theta)$. The state-action pair (s_T, a_T) can be assumed to be an arbitrary sample (\hat{s}, \hat{a}) chosen from the occupancy measure $\nu_{\rho}^{\pi_{\theta}}$. The algorithm then generates a π_{θ} -induced trajectory of length $T \sim \text{Geo}(1-\gamma)$, taking (\hat{s}, \hat{a}) as the starting point. The total reward and cost observed in this trajectory are assigned as $\hat{Q}_{\pi}^{\pi_{\theta}}(\hat{s},\hat{a})$ and $\hat{Q}_{e}^{\pi_{\theta}}(\hat{s},\hat{a})$ respectively. Next, another π_{θ} -induced trajectory of length $T\sim \text{Geo}(1-\gamma)$ is generated assuming the state, \hat{s} as the initiation point. The total reward and cost of this trajectory are assigned as $\hat{V}_r^{\pi_{\theta}}(\hat{s})$ and $\hat{V}_c^{\pi_{\theta}}(\hat{s})$ respectively. For $g \in \{r, c\}$, an estimate of the advantage value is computed as $\hat{A}_g^{\pi_{\theta}}(\hat{s},\hat{a}) = \hat{Q}_g^{\pi_{\theta}}(\hat{s},\hat{a}) - \hat{V}_g^{\pi_{\theta}}(\hat{s})$. Finally, the estimates of $F_{\rho}(\theta)$ and $H_{\rho}(\theta,\lambda)$ are obtained via (12) and (13) respectively which produces an estimation of the desired gradient in (14). The following Lemma demonstrates that the estimates produced by Algorithm 1 are unbiased.

Lemma 1. Let $\hat{J}_{c,\rho}(\theta)$, $\hat{\nabla}_{\omega}L_{\nu_{\alpha}^{\pi_{\theta}}}(\omega,\theta,\lambda)$ be the estimates produced by Algorithm 1 for a predefined set of parameters $(\omega, \theta, \lambda)$. The following equations hold.

$$\mathbf{E}\left[\hat{J}_{c,\rho}(\theta)\big|\theta\right] = J_{c,\rho}(\theta) \ \ \text{and} \ \ \mathbf{E}\left[\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega,\theta,\lambda)\big|\omega,\theta,\lambda\right] = \nabla_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega,\theta,\lambda)$$

In the absence of knowledge about the transition model, P, one can utilize the estimates generated by Algorithm 1 as good proxies for their true values. In particular, one can obtain an approximate value of the natural policy gradient $\omega_{\theta,\lambda}^*$ by iteratively minimizing the function $L_{\nu_{\alpha}^{\pi\theta}}(\cdot,\theta,\lambda)$ using the gradient estimate $\nabla_{\omega} L_{\nu_{\rho}^{\pi_{\theta}}}(\omega, \theta, \lambda)$. On the other hand, using the estimate $\hat{J}_{c,\rho}(\theta)$, an approximate update equation of the Lagrange parameter can be formed. Algorithm 2 uses these two ideas to obtain a policy that is close to the optimal one.

Algorithm 2 Primal-Dual Accelerated Natural Policy Gradient (PD-ANPG)

1: **Input:** Parameters (θ_0, λ_0) , Distribution ρ , Run-time Parameters K, H, Learning Parameters $\eta, \zeta, \alpha, \beta, \xi, \delta$

2: for
$$k \in \{0, \cdots, K-1\}$$
 do

3: $\mathbf{x}_0, \mathbf{v}_0 \leftarrow \mathbf{0}$

4: **for**
$$h \in \{0, \dots, H-1\}$$
 do \triangleright Inner Loop

$$\mathbf{y}_h \leftarrow \alpha \mathbf{x}_h + (1 - \alpha) \mathbf{v}_h \tag{15}$$

6: $\hat{G} \leftarrow \hat{\nabla}_{\omega} L_{\nu_{\rho}^{\pi_{\theta_{k}}}}(\omega, \theta_{k}, \lambda_{k}) \big|_{\omega = \mathbf{y}_{h}} \text{ (Algorithm 1)}$

$$\mathbf{x}_{h+1} \leftarrow \mathbf{y}_h - \delta \hat{G} \tag{16}$$

$$\mathbf{z}_h \leftarrow \beta \mathbf{y}_h + (1 - \beta) \mathbf{v}_h \tag{17}$$

$$\mathbf{v}_{h+1} \leftarrow \mathbf{z}_h - \xi \hat{G} \tag{18}$$

7: Tail Averaging:

$$\omega_k \leftarrow \frac{2}{H} \sum_{\frac{H}{2} < h \le H} \mathbf{x}_h \tag{19}$$

- 8: Obtain $\hat{J}_{c,\rho}(\theta_k)$ via Algorithm 1.
- 9: Parameter Updates:

$$\theta_{k+1} \leftarrow \theta_k + \eta \omega_k \tag{20}$$

$$\lambda_{k+1} \leftarrow \mathcal{P}_{\Lambda}[\lambda_k - \zeta \hat{J}_{c,\rho}(\theta_k)] \tag{21}$$

10: **Output:** $\{\theta_k\}_{k=0}^{K-1}$

Algorithm 2 has a nested loop structure. The *outer loop* runs K number of times. At a given instance, k, of the outer loop, the policy parameter θ_k and the Lagrange parameter, λ_k are updated via (20) and (21). The estimate, $\hat{J}_{c,\rho}(\theta_k)$ is computed via Algorithm 1. On the other hand, ω_k , the approximate value of the natural policy gradient $\omega_{\theta_k,\lambda_k}^*$ is obtained by iteratively minimizing $L_{\nu_\rho^{\pi_{\theta_k}}}(\cdot,\theta_k,\lambda_k)$ in H number of *inner loop* steps via the Accelerated Stochastic Gradient Descent (ASGD) procedure as stated in [8]. ASGD comprises the iterative updates (15)–(18) with tunable learning parameters $(\alpha,\beta,\xi,\delta)$ followed by a tail-averaging step (19). The gradient estimate utilized in (16) and (18) is obtained via Algorithm 1. It is worth mentioning that existing NPG algorithms such as that given in [7] typically apply the SGD, rather than the ASGD procedure, to obtain ω_k . The difference between these subroutines is that while SGD uses only the current gradient estimate to update ω_k , ASGD considers the contribution of all previous gradient estimates (momentum) using its convoluted iteration and tail-averaging steps.

4 Analysis

Our goal in this section is to characterize the rate of convergence of the objective function and the constraint violation if policy parameters are generated via Algorithm 2. We start by stating a few assumptions needed for the analysis.

Assumption 2. The log-likelihood function is G-Lipschitz and B-smooth where B, G > 0. Formally, the following relations hold $\forall \theta, \theta_1, \theta_2 \in \Theta$, and $\forall (s, a) \in S \times A$.

$$\|\nabla_{\theta} \log \pi_{\theta}(a|s)\| \leq G$$
, and $\|\nabla_{\theta} \log \pi_{\theta_1}(a|s) - \nabla_{\theta} \log \pi_{\theta_2}(a|s)\| \leq B\|\theta_1 - \theta_2\|$

Remark 1. Assumption 2 is commonly applied in proving convergence guarantees of policy gradient-type algorithms [35, 12, 7]. This assumption is obeyed by many widely used policy classes such as the class of neural networks with bounded weights.

Assumption 2 implies the boundedness of the gradient of the Lagrange function. This can be formally expressed as follows.

Lemma 2. If Assumption 2 holds, then the following inequality is true $\forall \theta \in \Theta$ and $\forall \lambda \in \Lambda$.

$$\|\nabla_{\theta} J_{L,\rho}(\theta,\lambda)\| \le \frac{G(1+\lambda_{\max})}{(1-\gamma)^2}$$

Proof. Statement (a) can be proven using (8) along with Assumption 2 and the facts that $|A_{L,\lambda}^{\pi_{\theta}}(s,a)|$ is bounded by $(1+\lambda)/(1-\gamma)$ and $\lambda \leq \lambda_{\max}$, $\forall \lambda \in \Lambda$.

The result established by Lemma 2 will be pivotal in our further analysis.

Assumption 3. The compatible function approximation error defined in (9) satisfies the inequality $L_{\nu_{\rho}^{\pi^*}}(\omega_{\theta,\lambda}^*,\theta,\lambda) \leq \epsilon_{\text{bias}}/2$, $\forall \theta \in \Theta$ and $\forall \lambda \in \Lambda$ where π^* is a solution to the original constrained optimization problem (2) and $\omega_{\theta,\lambda}^*$ is defined in (9). The term ϵ_{bias} is a non-negative constant. The factor 2 is used for notational convenience.

Remark 2. The term $\epsilon_{\rm bias}$ quantifies the expressivity of the parameterized policy class. For example, if the parameterization is complete i.e., includes all possible policies (such as in direct or softmax parameterization), then $\epsilon_{\rm bias} = 0$ [12]. A similar result can be proven for linear MDPs [36]. For incomplete policy classes, we have $\epsilon_{\rm bias} > 0$. However, if the class is sufficiently rich (such as neural networks with a large number of parameters), $\epsilon_{\rm bias}$ can be assumed to be negligibly small [37].

Assumption 4. There exists a positive constant μ_F such that $F_{\rho}(\theta) - \mu_F I_d$ is positive semidefinite i.e., $F_{\rho}(\theta) \succeq \mu_F I_d$, $\forall \theta \in \Theta$ where I_d is a $d \times d$ identity matrix and $F_{\rho}(\cdot)$ is defined in (7).

The property of the policy classes laid out in Assumption 4 is called Fisher Non-Degeneracy (FND) which essentially ensures that $\forall \theta \in \Theta$, the Fisher matrix $F_{\rho}(\theta)$ is away from the zero matrix by a certain amount. Observe that the Hessian of the function, $l_{\theta,\lambda}(\cdot) \triangleq L_{\nu_{\rho}^{\pi_{\theta}}}(\cdot,\theta,\lambda)$ is $F_{\rho}(\theta)$. Therefore, Assumption 4 also indicates that $l_{\theta,\lambda}$ is μ_F -strongly convex. This assumption is commonly applied in analyzing policy-gradient algorithms [4, 38, 7]. Assumption 4 also ensures that the matrix $F_{\rho}(\theta)$ is invertible, which, in turn, implies the uniqueness of the maximizer $\omega_{\theta,\lambda}^* = \arg\min_{\omega \in \mathbb{R}^d} l_{\theta,\lambda}(\omega)$. [39] describes a concrete set of policies that obeys Assumption 2-4.

4.1 Local-to-Global Convergence Lemma

Recall that our goal is to establish the global convergence rates. Lemma 3 (stated below) is the first step in that direction. Specifically, it demonstrates how the average optimality gap of the Lagrange function can be bounded by the first and second-order error of the gradient estimates.

Lemma 3. If the parameters $\{\theta_k, \lambda_k\}_{k=0}^{K-1}$ are updated via (20) and (21) and assumptions 2-4 hold, then the following inequality holds for any K.

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left(J_{\mathrm{L},\rho}(\pi^*, \lambda_k) - J_{\mathrm{L},\rho}(\theta_k, \lambda_k) \right) \leq \sqrt{\epsilon_{\mathrm{bias}}} + \frac{G}{K} \sum_{k=0}^{K-1} \mathbf{E} \| (\mathbf{E} \left[\omega_k \middle| \theta_k, \lambda_k \right] - \omega_k^*) \| + \frac{B\eta}{2K} \sum_{k=0}^{K-1} \mathbf{E} \| \omega_k \|^2 + \frac{1}{\eta K} \mathbf{E}_{s \sim d_\rho^{\pi^*}} \left[KL(\pi^*(\cdot | s) \middle| \pi_{\theta_0}(\cdot | s)) \right]$$
(22)

where $\omega_k^* \triangleq \omega_{\theta_k,\lambda_k}^*$, $\omega_{\theta_k,\lambda_k}^*$ is the natural policy gradient defined in (9), and π^* is the solution to the constrained optimization (2). Finally, ω_k is the approximation of ω_k^* given by (19) and $KL(\cdot||\cdot|)$ is the KL-divergence.

Note the presence of the term, $\epsilon_{\rm bias}$ in (22). It shows that due to the incompleteness of the parameterized policy class, the average optimality error cannot be made arbitrarily small. It is worth mentioning that many existing CMDP analyses (such as [4]) follow a path similar to that of Lemma 3. However, while the first order term in those works turns out to be $\mathbf{E} \|\omega_k - \omega_k^*\|$, we improved it to $\mathbf{E} \|\mathbf{E}[\omega_k|\theta_k,\lambda_k] - \omega_k^*\|$. Such seemingly insignificant improvement has important ramifications for

our analysis as explained later in the paper. The second order term in (22) can be expanded as,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \|\omega_{k}\|^{2} \leq \frac{2}{K} \sum_{k=0}^{K-1} \mathbf{E} \|\omega_{k} - \omega_{k}^{*}\|^{2} + \frac{2}{K} \sum_{k=0}^{K-1} \mathbf{E} \|\omega_{k}^{*}\|^{2}
\leq \frac{2}{K} \sum_{k=0}^{K-1} \mathbf{E} \|\omega_{k} - \omega_{k}^{*}\|^{2} + \frac{2}{\mu_{F}^{2} K} \sum_{k=0}^{K-1} \mathbf{E} \|\nabla_{\theta} J_{L,\rho}(\theta_{k}, \lambda_{k})\|^{2}
\leq \frac{2}{K} \sum_{k=0}^{K-1} \mathbf{E} \|\omega_{k} - \omega_{k}^{*}\|^{2} + \frac{2G^{2}(1 + \lambda_{\max})^{2}}{\mu_{F}^{2}(1 - \gamma)^{4}}$$
(23)

where (a) utilises $\omega_k^* = F_\rho(\theta_k)^\dagger \nabla_\theta J_{\mathrm{L},\rho}(\theta_k,\lambda_k)$ and Assumption 4. The second inequality applies Lemma 2 together with the fact that $\lambda_k \in \Lambda$. Our next subsection provides a bound on $\mathbf{E} \|\omega_k - \omega_k^*\|^2$ and the first order term $\mathbf{E} \|\mathbf{E}[\omega_k|\theta_k,\lambda_k] - \omega_k^*\|$.

4.2 Local Convergence of the Natural Policy Gradient

To deliver the promised bounds, we first provide some characterization of the gradient estimate.

Lemma 4. Let $\nabla_{\omega} L_{\nu_{\rho}^{\pi_{\theta}}}(\omega, \theta, \lambda)$ be the estimate produced by Algorithm 1. Under assumptions 2 and 4, the following semidefinite inequality holds for any $\theta \in \Theta$ and $\lambda \in \Lambda$.

$$\mathbf{E}\left[\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega_{\theta,\lambda}^{*},\theta,\lambda)\otimes\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega_{\theta,\lambda}^{*},\theta,\lambda)\right] \preceq \sigma^{2}F_{\rho}(\theta)$$

where ω_{θ}^* , $F_{\rho}(\theta)$ are given by (9) and (7) respectively and σ^2 is defined below.

$$\sigma^2 \triangleq \frac{1}{(1-\gamma)^4} \left[\frac{2G^4}{\mu_F^2} + 32 \right] (1+\lambda_{\text{max}})^2$$
 (24)

The term σ^2 defined in Lemma 4 can be described as the scaled variance of the gradient estimate, $\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega,\theta,\lambda)$. Note that if the estimates were non-stochastic (i.e., deterministic), we would have $\sigma^2=0$ since $\nabla_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega_{\theta,\lambda}^*,\theta,\lambda)=0$. The last equation can be proved using the definition of $\omega_{\theta,\lambda}^*$ given in (9), the gradient expression provided in (10), and observing that the Fisher matrix, $F_{\rho}(\theta)$ is invertible due to Assumption 4. The above information is crucial in bounding the first-order error, as stated in the following lemma.

Lemma 5. If assumptions 2 and 4 hold, then the following relations are satisfied $\forall k \in \{0, \cdots, K-1\}$ with learning rates $\alpha = \frac{3\sqrt{5}G^2}{\mu_F + 3\sqrt{5}G^2}$, $\beta = \frac{\mu_F}{9G^2}$, $\xi = \frac{1}{3\sqrt{5}G^2}$, and $\delta = \frac{1}{5G^2}$ provided that the inner loop length of Algorithm 2 obeys $H > \bar{C}\frac{G^2}{\mu_F}\log\left(\sqrt{\mathrm{d}\frac{G^2}{\mu_F}}\right)$ for some universal constant, \bar{C} .

$$\mathbf{E}\|\omega_k - \omega_k^*\|^2 \le 22 \frac{\sigma^2 d}{\mu_F H} + C \exp\left(-\frac{\mu_F}{20G^2} H\right) \left[\frac{(1 + \lambda_{\text{max}})^2}{\mu_F (1 - \gamma)^4} \right], \tag{25}$$

$$\mathbf{E}\|\mathbf{E}\left[\omega_{k}|\theta_{k},\lambda_{k}\right] - \omega_{k}^{*}\| \leq \sqrt{C}\exp\left(-\frac{\mu_{F}}{40G^{2}}H\right)\left[\frac{1+\lambda_{\max}}{\sqrt{\mu_{F}}(1-\gamma)^{2}}\right]$$
(26)

where C denotes a universal constant, ω_k is given by (19), and σ^2 is defined in (24).

The first bound (25) is a consequence of Lemma 4 and the ASGD convergence result provided in [8] (Corollary 2). To gain intuition about the second result (26), note that by taking the conditional expectation $\mathbf{E}[\cdot|\theta_k,\lambda_k]$ on both sides of the ASGD iterations (15)–(18), and applying the unbiasedness of the gradient estimate (Lemma 1) we obtain the following $\forall h \in \{0,\dots,H-1\}$.

$$\bar{\mathbf{y}}_{h} = \alpha \bar{\mathbf{x}}_{h} + (1 - \alpha) \bar{\mathbf{v}}_{h},
\bar{\mathbf{x}}_{h+1} = \bar{\mathbf{y}}_{h} - \delta \nabla_{\omega} L_{\nu_{\rho}^{\pi_{\theta_{k}}}}(\omega, \theta_{k}, \lambda_{k}) \big|_{\omega = \bar{\mathbf{y}}_{h}}
\bar{\mathbf{z}}_{h} = \beta \bar{\mathbf{y}}_{h} + (1 - \beta) \bar{\mathbf{v}}_{h},
\bar{\mathbf{v}}_{h+1} = \bar{\mathbf{z}}_{h} - \xi \nabla_{\omega} L_{\nu_{\rho}^{\pi_{\theta_{k}}}}(\omega, \theta_{k}, \lambda_{k}) \big|_{\omega = \bar{\mathbf{y}}_{h}}$$
(27)

where $\bar{l}_h = \mathbf{E}[l_h|\theta_k, \lambda_k], l \in \{\mathbf{v}, \mathbf{x}, \mathbf{y}, \mathbf{z}\}$. Moreover, taking conditional expectation on both sides of the tail averaging process (19), we arrive at the following.

$$\bar{\omega}_k \triangleq \mathbf{E} \left[\omega_k \middle| \theta_k, \lambda_k \right] = \frac{2}{H} \sum_{\frac{H}{2} < h \le H} \bar{\mathbf{x}}_h \tag{28}$$

Note that the steps (27)–(28) resemble the iterative updates of a deterministic ASGD. This allows us to obtain $\mathbf{E}\|\bar{\omega}_k - \omega_k^*\|$ by substituting $\sigma^2 = 0$ in (25) and applying the Cauchy-Schwarz inequality.

4.3 Global Convergence of the Lagrange

Combining Lemma 3, (23) and using the expected gradient errors provided by Lemma 5, we bound the average Lagrange optimality gap as a function of tunable parameters H and K as stated in the following corollary.

Corollary 1. Consider the same setup and the choice of parameters described in Lemma 3-5. The following inequality holds if assumptions 2-4 are met.

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left(J_{L,\rho}(\pi^*, \lambda_k) - J_{L,\rho}(\theta_k, \lambda_k) \right) \leq \sqrt{\epsilon_{\text{bias}}} + \left[\frac{G\sqrt{C}(1 + \lambda_{\text{max}})}{\sqrt{\mu_F}(1 - \gamma)^2} \right] \exp\left(-\frac{\mu_F}{40G^2} H \right) \\
+ B\eta \left[\frac{22\sigma^2 d}{\mu_F H} + \exp\left(-\frac{\mu_F}{20G^2} H \right) \left[\frac{C(1 + \lambda_{\text{max}})^2}{\mu_F(1 - \gamma)^4} \right] + \frac{G^2(1 + \lambda_{\text{max}})^2}{\mu_F^2(1 - \gamma)^4} \right] \\
+ \frac{1}{\eta K} \mathbf{E}_{s \sim d_\rho^{\pi^*}} \left[KL(\pi^*(\cdot|s) \| \pi_{\theta_0}(\cdot|s)) \right] \tag{29}$$

Corollary 1 bounds the optimality error of the Lagrange function as $\mathcal{O}(\sqrt{\epsilon_{\text{bias}}} + \exp{(-C_0 H)} + \eta + \frac{1}{\eta K})$ where C_0 is some problem specific constant. Interestingly, the dual learning parameter, ζ does not appear in this bound. However, the next section shows that ζ plays a pivotal role in deciding the objective and constraint violation rates.

4.4 Decoupling the Objective and the Constraint Violation Rates

The goal of the following theorem is to choose optimal values of the tunable parameters and decouple the objective and constraint violation rates from the Lagrange convergence result given in (29).

Theorem 1. Consider the same setup and the choice of parameters given in Lemma 3–5. Assume $\eta = (1-\gamma)^2(1+\lambda_{\max})^{-1}/\sqrt{K}$, $\zeta = \lambda_{\max}(1-\gamma)/\sqrt{K}$, and $\lambda_{\max} = 2/[(1-\gamma)c_{\text{slater}}]$. For sufficiently small $\epsilon > 0$, the following inequalities hold

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k) \right] \le \sqrt{\epsilon_{\text{bias}}} + \epsilon,$$

$$\mathbf{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} -J_{c,\rho}(\theta_k) \right] \le (1-\gamma) c_{\text{slater}} \sqrt{\epsilon_{\text{bias}}} + \epsilon$$
(30)

whenever assumptions 1-4 are met, $H=\mathcal{O}(\log(\epsilon^{-1}))$ and $K=\mathcal{O}((1-\gamma)^{-6}\epsilon^{-2})$. Therefore, the sample complexity to ensure (30) is $\mathcal{O}((1-\gamma)^{-1}HK)=\tilde{\mathcal{O}}((1-\gamma)^{-7}\epsilon^{-2})$. It is to be clarified that the $(1-\gamma)^{-1}$ factor in the sample complexity calculation appears due to the fact that it requires $\mathcal{O}((1-\gamma)^{-1})$ samples on an average to obtain a gradient estimate via Algorithm 1.

Theorem 1 dictates that, with appropriate choice of the parameters, the rate of convergence of the objective and that of constraint violation can be bounded as $\mathcal{O}(\sqrt{\epsilon_{\mathrm{bias}}} + \epsilon)$ with $\tilde{\mathcal{O}}((1-\gamma)^{-7}\epsilon^{-2})$ samples. This beats the SOTA $\tilde{\mathcal{O}}(\epsilon^{-4})$ sample complexity of [4] and achieves the theoretical lower bound. Our derived sample complexity is also dependent on c_{slater} . However, this is not explicitly mentioned in Theorem 1. Interested readers can find such details in the appendix.

Remark 3. Note the importance of the nested expectation in the first-order term $\mathbf{E}\|\mathbf{E}[\omega_k|\theta_k,\lambda_k]-\omega_k^*\|$ in Lemma 3. Lemma 5 bounds this term as $\mathcal{O}(\exp(-C_0H))$ where C_0 is a problem dependent constant. Other terms in the Lagrange optimality bound (Lemma 3) can be expressed as $\mathcal{O}(\eta+\frac{1}{\eta K})$. Moreover, following the proof of Theorem 1, one sees that decoupling the objective optimality error and constraint violation bounds incurs additional $\mathcal{O}(\zeta+\frac{1}{\zeta K})$ terms. Choosing η,ζ as prescribed in Theorem 1, makes both the objective and constraint violation errors as $\mathcal{O}(\sqrt{\epsilon_{\text{bias}}}+\exp(-C_0H)+K^{-0.5})$. This allows us to take $H=\tilde{\mathcal{O}}(1)$ and $K=\mathcal{O}(\epsilon^{-2})$, leading to $\tilde{\mathcal{O}}(\epsilon^{-2})$ sample complexity. Had the first-order term been $\mathbf{E}\|\omega_k-\omega_k^*\|$, it would have resulted in $\mathcal{O}(\sqrt{\epsilon_{\text{bias}}}+H^{-0.5}+K^{-0.5})$ objective and constraint violation errors which would have lead to $\mathcal{O}(\epsilon^{-4})$ sample complexity.

5 Conclusions and Limitations

This paper considers the problem of learning a CMDP where the goal is to maximize the objective value function while guaranteeing that the cost value exceeds a predefined threshold. We propose an acceleration-based primal-dual natural policy gradient algorithm that ensures ϵ optimality gap and ϵ constraint violation with $\tilde{\mathcal{O}}(\epsilon^{-2})$ sample complexity. This improves upon the previous state-of-the-art sample complexity of $\mathcal{O}(\epsilon^{-4})$ and achieves the theoretical lower bound. Future works include applying the idea of acceleration-based NPG to improve sample complexities in other related domains of constrained reinforcement learning, e.g., non-linear CMDP, average reward CMDP, etc.

References

- [1] Al-Abbasi, A. O., A. Ghosh, V. Aggarwal. Deeppool: Distributed model-free algorithm for ridesharing using deep reinforcement learning. *IEEE Transactions on Intelligent Transportation Systems*, 20(12):4714–4727, 2019.
- [2] Geng, N., T. Lan, et al. A multi-agent reinforcement learning perspective on distributed traffic engineering. In 2020 IEEE 28th International Conference on Network Protocols (ICNP), pages 1–11. IEEE, 2020.
- [3] Gonzalez, G., M. Balakuntala, M. Agarwal, et al. Asap: A semi-autonomous precise system for telesurgery during communication delays. *IEEE Transactions on Medical Robotics and Bionics*, 5(1):66–78, 2023.
- [4] Bai, Q., A. S. Bedi, V. Aggarwal. Achieving zero constraint violation for constrained reinforcement learning via conservative natural policy gradient primal-dual algorithm. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, pages 6737–6744. 2023.
- [5] Ding, D., K. Zhang, T. Basar, M. Jovanovic. Natural policy gradient primal-dual method for constrained markov decision processes. *Advances in Neural Information Processing Systems*, 33:8378–8390, 2020.
- [6] Xu, T., Y. Liang, G. Lan. Crpo: A new approach for safe reinforcement learning with convergence guarantee. In *International Conference on Machine Learning*, pages 11480–11491. 2021.
- [7] Liu, Y., K. Zhang, T. Basar, W. Yin. An improved analysis of (variance-reduced) policy gradient and natural policy gradient methods. *Advances in Neural Information Processing Systems*, 33:7624–7636, 2020.
- [8] Jain, P., S. M. Kakade, R. Kidambi, P. Netrapalli, A. Sidford. Accelerating stochastic gradient descent for least squares regression. In *Conference On Learning Theory*, pages 545–604. 2018.
- [9] Liu, T., R. Zhou, D. Kalathil, P. Kumar, C. Tian. Policy optimization for constrained mdps with provable fast global convergence. *arXiv preprint arXiv:2111.00552*, 2021.
- [10] Zeng, S., T. T. Doan, J. Romberg. Finite-time complexity of online primal-dual natural actorcritic algorithm for constrained markov decision processes. In 2022 IEEE 61st Conference on Decision and Control (CDC), pages 4028–4033. IEEE, 2022.
- [11] Vaswani, S., L. Yang, C. Szepesvári. Near-optimal sample complexity bounds for constrained mdps. Advances in Neural Information Processing Systems, 35:3110–3122, 2022.

- [12] Agarwal, A., S. M. Kakade, J. D. Lee, G. Mahajan. On the theory of policy gradient methods: Optimality, approximation, and distribution shift. *The Journal of Machine Learning Research*, 22(1):4431–4506, 2021.
- [13] Bhandari, J., D. Russo. On the linear convergence of policy gradient methods for finite mdps. In *International Conference on Artificial Intelligence and Statistics*, pages 2386–2394. 2021.
- [14] Cen, S., C. Cheng, Y. Chen, Y. Wei, Y. Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 70(4):2563–2578, 2022.
- [15] Lan, G. Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *Mathematical programming*, 198(1):1059–1106, 2023.
- [16] Zhan, W., S. Cen, B. Huang, Y. Chen, J. D. Lee, Y. Chi. Policy mirror descent for regularized reinforcement learning: A generalized framework with linear convergence. *SIAM Journal on Optimization*, 33(2):1061–1091, 2023.
- [17] Xu, P., F. Gao, Q. Gu. An improved convergence analysis of stochastic variance-reduced policy gradient. In *Uncertainty in Artificial Intelligence*, pages 541–551. 2020.
- [18] Gargiani, M., A. Zanelli, A. Martinelli, T. Summers, J. Lygeros. Page-pg: A simple and loopless variance-reduced policy gradient method with probabilistic gradient estimation. In *International Conference on Machine Learning*, pages 7223–7240. 2022.
- [19] Huang, F., S. Gao, J. Pei, H. Huang. Momentum-based policy gradient methods. In *International conference on machine learning*, pages 4422–4433. 2020.
- [20] Salehkaleybar, S., M. Khorasani, N. Kiyavash, N. He, P. Thiran. Momentum-based policy gradient with second-order information. *Transactions on Machine Learning Research*, 2024.
- [21] Shen, Z., A. Ribeiro, H. Hassani, H. Qian, C. Mi. Hessian aided policy gradient. In *International conference on machine learning*, pages 5729–5738. 2019.
- [22] Chen, Z., S. Khodadadian, S. T. Maguluri. Finite-sample analysis of off-policy natural actorcritic with linear function approximation. *IEEE Control Systems Letters*, 6:2611–2616, 2022.
- [23] Chen, Z., S. T. Maguluri. Sample complexity of policy-based methods under off-policy sampling and linear function approximation. In *International Conference on Artificial Intelligence and Statistics*, pages 11195–11214. 2022.
- [24] Khodadadian, S., T. T. Doan, J. Romberg, S. T. Maguluri. Finite sample analysis of two-time-scale natural actor-critic algorithm. *IEEE Transactions on Automatic Control*, 2022.
- [25] Fatkhullin, I., A. Barakat, A. Kireeva, N. He. Stochastic policy gradient methods: Improved sample complexity for fisher-non-degenerate policies. In *International Conference on Machine Learning*, pages 9827–9869. 2023.
- [26] Masiha, S., S. Salehkaleybar, N. He, N. Kiyavash, P. Thiran. Stochastic second-order methods improve best-known sample complexity of SGD for gradient-dominated functions. *Advances in Neural Information Processing Systems*, 35:10862–10875, 2022.
- [27] Mondal, W. U., V. Aggarwal. Improved sample complexity analysis of natural policy gradient algorithm with general parameterization for infinite horizon discounted reward Markov decision processes. In *International Conference on Artificial Intelligence and Statistics*, pages 3097–3105. 2024.
- [28] Efroni, Y., S. Mannor, M. Pirotta. Exploration-exploitation in constrained mdps. *arXiv preprint arXiv:2003.02189*, 2020.
- [29] Liu, T., R. Zhou, D. Kalathil, P. Kumar, C. Tian. Learning policies with zero or bounded constraint violation for constrained mdps. Advances in Neural Information Processing Systems, 34:17183–17193, 2021.
- [30] Ding, D., X. Wei, Z. Yang, Z. Wang, M. Jovanovic. Provably efficient safe exploration via primal-dual policy optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 3304–3312. 2021.
- [31] He, J., D. Zhou, Q. Gu. Nearly minimax optimal reinforcement learning for discounted mdps. Advances in Neural Information Processing Systems, 34:22288–22300, 2021.

- [32] Wei, H., X. Liu, L. Ying. Triple-q: A model-free algorithm for constrained reinforcement learning with sublinear regret and zero constraint violation. In *International Conference on Artificial Intelligence and Statistics*, pages 3274–3307. 2022.
- [33] Bai, Q., A. S. Bedi, M. Agarwal, A. Koppel, V. Aggarwal. Achieving zero constraint violation for constrained reinforcement learning via primal-dual approach. In *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, pages 3682–3689. 2022.
- [34] Sutton, R. S., D. McAllester, S. Singh, Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. *Advances in neural information processing systems*, 12, 1999.
- [35] Zhang, J., C. Ni, C. Szepesvari, M. Wang. On the convergence and sample efficiency of variance-reduced policy gradient method. Advances in Neural Information Processing Systems, 34:2228–2240, 2021.
- [36] Jin, C., Z. Yang, Z. Wang, M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. In J. Abernethy, S. Agarwal, eds., *Proceedings of Thirty Third Conference on Learning Theory*, vol. 125 of *Proceedings of Machine Learning Research*, pages 2137–2143. PMLR, 2020.
- [37] Wang, L., Q. Cai, Z. Yang, Z. Wang. Neural policy gradient methods: Global optimality and rates of convergence. In *International Conference on Learning Representations*. 2019.
- [38] Zhang, K., A. Koppel, H. Zhu, T. Basar. Global convergence of policy gradient methods to (almost) locally optimal policies. *SIAM Journal on Control and Optimization*, 58(6):3586–3612, 2020.
- [39] Mondal, W. U., V. Aggarwal, S. Ukkusuri. Mean-field control based approximation of multiagent reinforcement learning in presence of a non-decomposable shared global state. *Transac*tions on Machine Learning Research, 2023.
- [40] Ding, D., K. Zhang, J. Duan, T. Başar, M. R. Jovanović. Convergence and sample complexity of natural policy gradient primal-dual methods for constrained mdps. *arXiv* preprint *arXiv*:2206.02346, 2022.

A Helper Lemma

Lemma 6. With slight abuse of notation, define $J_{L,\rho}(\pi,\lambda) \triangleq J_{r,\rho}^{\pi} + \lambda J_{c,\rho}^{\pi}$. The following relation holds for any two policies π_1, π_2 and $\lambda \in \Lambda$.

$$J_{L,\rho}(\pi_1, \lambda) - J_{L,\rho}(\pi_2, \lambda) = \frac{1}{1 - \gamma} \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi_1}} \left[A_{L,\lambda}^{\pi_2}(s, a) \right]$$
(31)

Proof. This can be proved using Lemma 2 of [12] and the definition of the Lagrange function. \Box

B Proof of Lemma 1

Proof. Fix arbitrary θ and λ . Note that the following equation holds due to the definition of $\hat{J}_{c,\rho}(\theta)$.

$$\mathbf{E}\left[\hat{J}_{c,\rho}(\theta)\middle|\theta\right] = (1-\gamma)\mathbf{E}\left[\sum_{t=0}^{\infty} \gamma^{t} \sum_{j=0}^{t} c(s_{j}, a_{j})\middle|s_{0} \sim \rho, \pi_{\theta}\right]$$

$$= (1-\gamma)\mathbf{E}\left[\sum_{j=0}^{\infty} c(s_{j}, a_{j}) \sum_{t=j}^{\infty} \gamma^{t}\middle|s_{0} \sim \rho, \pi_{\theta}\right]$$

$$= \mathbf{E}\left[\sum_{j=0}^{\infty} \gamma^{j} c(s_{j}, a_{j})\middle|s_{0} \sim \rho, \pi_{\theta}\right] = J_{c,\rho}(\theta)$$
(32)

To prove the unbiasedness of the gradient, we first prove that the distribution of the sample pair (\hat{s}, \hat{a}) produced by Algorithm 1 is indeed $\nu_{\rho}^{\pi_{\theta}}$. Observe the following.

$$\Pr(\hat{s} = s, \hat{a} = a | \rho, \pi_{\theta}) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^{t} \Pr(s_{t} = s, a_{t} = a | s_{0} \sim \rho, \pi_{\theta}) = \nu_{\rho}^{\pi_{\theta}}(s, a)$$
(33)

Next, we show that for a given pair (s, a), the estimate $\hat{Q}_g^{\pi_\theta}(s, a)$, $g \in \{r, c\}$ is unbiased.

$$\mathbf{E}\left[\hat{Q}_{g}^{\pi_{\theta}}(s,a)\middle|\theta,s,a\right] = (1-\gamma)\mathbf{E}\left[\sum_{t=0}^{\infty} \gamma^{t} \sum_{j=0}^{t} g(s_{i},a_{i})\middle|s_{0} = s, a_{0} = a, \pi_{\theta}\right]$$

$$= (1-\gamma)\mathbf{E}\left[\sum_{j=0}^{\infty} g(s_{i},a_{i}) \sum_{t=j}^{t} \gamma^{t}\middle|s_{0} = s, a_{0} = a, \pi_{\theta}\right]$$

$$= \mathbf{E}\left[\sum_{j=0}^{\infty} \gamma^{i} g(s_{i},a_{i})\middle|s_{0} = s, a_{0} = a, \pi_{\theta}\right] = Q_{g}^{\pi_{\theta}}(s,a)$$

$$(34)$$

In a similar fashion, one can establish that $\mathbf{E}[\hat{V}_g^{\pi_{\theta}}(s) | \theta, s] = V_g^{\pi_{\theta}}(s), \forall g \in \{r, c\}$. Combining this with (34) leads to: $\mathbf{E}[\hat{A}_{\mathrm{L},\lambda}^{\pi_{\theta}}(s,a) | \theta, \lambda, s, a] = A_{\mathrm{L},\lambda}^{\pi_{\theta}}(s,a)$. We arrive at,

$$\mathbf{E}\left[\hat{F}_{\rho}(\theta)\big|\theta\right] = \mathbf{E}_{(\hat{s},\hat{a})\sim\nu_{\rho}^{\pi_{\theta}}}\left[\nabla_{\theta}\log\pi_{\theta}(\hat{a}|\hat{s})\otimes\nabla_{\theta}\log\pi_{\theta}(\hat{a}|\hat{s})\big|\theta\right] = F_{\rho}(\theta),\tag{35}$$

$$\mathbf{E}\left[\hat{H}_{\rho}(\theta,\lambda)\big|\theta,\lambda\right] = \mathbf{E}_{(\hat{s},\hat{a})\sim\nu_{\rho}^{\pi_{\theta}}} \left[\mathbf{E}\left[\hat{A}_{\mathrm{L},\lambda}^{\pi_{\theta}}(\hat{s},\hat{a})\bigg|\theta,\lambda,\hat{s},\hat{a}\right]\nabla_{\theta}\log\pi_{\theta}(\hat{a}|\hat{s})\bigg|\theta,\lambda\right] \\ = \mathbf{E}_{(\hat{s},\hat{a})\sim\nu_{\rho}^{\pi_{\theta}}} \left[A_{\mathrm{L},\lambda}^{\pi_{\theta}}(\hat{s},\hat{a})\nabla_{\theta}\log\pi_{\theta}(\hat{a}|\hat{s})\big|\theta,\lambda\right] = H_{\rho}(\theta,\lambda)$$
(36)

Finally, we arrive at the following by utilizing the definitions (10) and (14).

$$\mathbf{E}\left[\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega,\theta,\lambda)\big|\omega,\theta,\lambda\right] = \mathbf{E}\left[\hat{F}_{\rho}(\theta)\big|\theta\right]\omega - \frac{1}{1-\gamma}\mathbf{E}\left[\hat{H}_{\rho}(\theta,\lambda)\big|\theta,\lambda\right]$$

$$= F_{\rho}(\theta)\omega - \frac{1}{1-\gamma}H_{\rho}(\theta,\lambda) = \nabla_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega,\theta,\lambda)$$
(37)

This concludes the proof.

C Proof of Lemma 3

Proof. Invoking the definition of KL divergence, we arrive at the following series of inequalities.

$$\begin{split} \mathbf{E}_{s \sim d_{\rho}^{\pi^*}} [KL(\pi^*(\cdot|s) || \pi_{\theta_k}(\cdot|s)) - KL(\pi^*(\cdot|s) || \pi_{\theta_{k+1}}(\cdot|s))] \\ &= \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} \left[\log \frac{\pi_{\theta_{k+1}}(a|s)}{\pi_{\theta_k}(a|s)} \right] \\ \overset{(a)}{\geq} \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot (\theta_{k+1} - \theta_k)] - \frac{B}{2} || \theta_{k+1} - \theta_k ||^2 \\ &= \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot \omega_k] - \frac{B\eta^2}{2} || \omega_k ||^2 \\ &= \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot \omega_k^*] + \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot (\omega_k - \omega_k^*)] - \frac{B\eta^2}{2} || \omega_k ||^2 \\ &= \eta [J_{\mathbf{L},\rho}(\pi^*, \lambda_k) - J_{\mathbf{L},\rho}(\theta_k, \lambda_k)] + \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot (\omega_k - \omega_k^*)] - \frac{B\eta^2}{2} || \omega_k ||^2 \\ &= \eta [J_{\mathbf{L},\rho}(\pi^*, \lambda_k) - J_{\mathbf{L},\rho}(\theta_k, \lambda_k)] + \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot (\omega_k - \omega_k^*)] - \frac{B\eta^2}{2} || \omega_k ||^2 \\ &\stackrel{(b)}{=} \eta [J_{\mathbf{L},\rho}(\pi^*, \lambda_k) - J_{\mathbf{L},\rho}(\theta_k, \lambda_k)] + \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot \omega_k^* - \frac{1}{1 - \gamma} A_{\mathbf{L},\lambda_k}^{\pi_{\theta_k}}(s,a)] \\ &+ \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot (\omega_k - \omega_k^*)] - \frac{B\eta^2}{2} || \omega_k ||^2 \\ &\stackrel{(c)}{=} \eta [J_{\mathbf{L},\rho}(\pi^*, \lambda_k) - J_{\mathbf{L},\rho}(\theta_k, \lambda_k)] - \eta \sqrt{\mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot \omega_k^* - \frac{1}{1 - \gamma} A_{\mathbf{L},\lambda_k}^{\pi_{\theta_k}}(s,a)]^2 \\ &+ \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot (\omega_k - \omega_k^*)] - \frac{B\eta^2}{2} || \omega_k ||^2 \\ &\stackrel{(d)}{=} \eta [J_{\mathbf{L},\rho}(\pi^*, \lambda_k) - J_{\mathbf{L},\rho}(\theta_k, \lambda_k)] - \eta \sqrt{\epsilon_{\text{bias}}} \\ &+ \eta \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi^*}} [\nabla_{\theta} \log \pi_{\theta_k}(a|s) \cdot (\omega_k - \omega_k^*)] - \frac{B\eta^2}{2} || \omega_k ||^2 \end{aligned}$$

where the step (a) holds by Assumption 2 and step (b) holds by Lemma 6. Step (c) uses the convexity of the function $f(x) = x^2$. Finally, step (d) comes from the Assumption 3. Rearranging the terms and taking expectations on both sides, we have,

$$\mathbf{E}\left[J_{\mathcal{L},\rho}(\pi^{*},\lambda_{k}) - J_{\mathcal{L},\rho}(\theta_{k},\lambda_{k})\right] \leq -\mathbf{E}_{(s,a)\sim\nu_{\rho}^{\pi^{*}}}\mathbf{E}\left[\nabla_{\theta}\log\pi_{\theta_{k}}(a|s)\cdot(\mathbf{E}\left[\omega_{k}|\theta_{k},\lambda_{k}\right] - \omega_{k}^{*})\right]
+ \frac{B\eta}{2}\mathbf{E}\|\omega_{k}\|^{2} + \frac{1}{\eta}\mathbf{E}_{s\sim d_{\rho}^{\pi^{*}}}\left[\mathbf{E}\left[KL(\pi^{*}(\cdot|s)\|\pi_{\theta_{k}}(\cdot|s))\right] - \mathbf{E}\left[KL(\pi^{*}(\cdot|s)\|\pi_{\theta_{k+1}}(\cdot|s))\right]\right] + \sqrt{\epsilon_{\text{bias}}}
\stackrel{(a)}{\leq} \sqrt{\epsilon_{\text{bias}}} + G\mathbf{E}\left\|(\mathbf{E}\left[\omega_{k}|\theta_{k},\lambda_{k}\right] - \omega_{k}^{*})\right\| + \frac{B\eta}{2}\mathbf{E}\|\omega_{k}\|^{2}
+ \frac{1}{\eta}\mathbf{E}_{s\sim d_{\rho}^{\pi^{*}}}\left[\mathbf{E}\left[KL(\pi^{*}(\cdot|s)\|\pi_{\theta_{k}}(\cdot|s))\right] - \mathbf{E}\left[KL(\pi^{*}(\cdot|s)\|\pi_{\theta_{k+1}}(\cdot|s))\right]\right]$$
(39)

where (a) follows from Assumption 2. Summing from k = 0 to K - 1, using the non-negativity of KL divergence and dividing the resulting expression by K, we obtain,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left(J_{\mathrm{L},\rho}(\pi^*, \lambda_k) - J_{\mathrm{L},\rho}(\theta_k, \lambda_k) \right) \leq \sqrt{\epsilon_{\mathrm{bias}}} + \frac{G}{K} \sum_{k=0}^{K-1} \mathbf{E} \| (\mathbf{E} \left[\omega_k \middle| \theta_k, \lambda_k \right] - \omega_k^*) \| + \frac{B\eta}{2K} \sum_{k=0}^{K-1} \mathbf{E} \| \omega_k \|^2 + \frac{1}{\eta K} \mathbf{E}_{s \sim d_\rho^{\pi^*}} \left[KL(\pi^*(\cdot | s) \middle| \pi_{\theta_0}(\cdot | s)) \right]$$
(40)

This concludes the proof.

D Proof of Lemma 4

Proof. Fix a $\theta \in \Theta$ and a $\lambda \in \Lambda$. Observe the following equation.

$$\mathbf{E}\left[\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega_{\theta,\lambda}^{*},\theta,\lambda)\otimes\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega_{\theta,\lambda}^{*},\theta,\lambda)\right]$$

$$=\mathbf{E}_{(s,a)\sim\nu_{\rho}^{\pi_{\theta}}}\left[\mathbf{E}\left[\underbrace{\nabla_{\theta}\log\pi_{\theta}(a|s)\cdot\omega_{\theta,\lambda}^{*}-\frac{1}{1-\gamma}\hat{A}_{\mathrm{L},\lambda}^{\pi_{\theta}}(s,a)}^{2}\right]^{2}\nabla_{\theta}\log\pi_{\theta}(a|s)\otimes\nabla_{\theta}\log\pi_{\theta}(a|s)\right]$$

$$\triangleq \zeta_{\theta,\lambda}(s,a)$$

To prove the lemma, it is sufficient to demonstrate that $\mathbf{E}[\zeta_{\theta,\lambda}(s,a)] \leq \sigma^2$, $\forall (s,a)$. Notice the chain of inequalities stated below.

$$\mathbf{E} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot \omega_{\theta,\lambda}^{*} - \frac{1}{1-\gamma} \hat{A}_{L,\lambda}^{\pi_{\theta}}(s,a) \right]^{2} \\
\leq 2 \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \cdot \omega_{\theta,\lambda}^{*} \right]^{2} + \frac{2}{(1-\gamma)^{2}} \mathbf{E} \left[\hat{A}_{L,\lambda}^{\pi_{\theta}}(s,a) \right]^{2} \\
\stackrel{(a)}{\leq} 2 \|\nabla_{\theta} \log \pi_{\theta}(a|s)\|^{2} \|\omega_{\theta,\lambda}^{*}\|^{2} + \frac{4(1+\lambda)^{2}}{(1-\gamma)^{2}} \max_{g \in \{r,c\}} \left\{ \mathbf{E} \left[\hat{A}_{g}^{\pi_{\theta}}(s,a) \right]^{2} \right\} \\
\stackrel{(b)}{\leq} 2G^{2} \|F_{\rho}(\theta)^{\dagger} \nabla_{\theta} J_{L,\rho}(\theta,\lambda)\|^{2} + \frac{8(1+\lambda)^{2}}{(1-\gamma)^{2}} \max_{g \in \{r,c\}} \left\{ \mathbf{E} \left[\hat{Q}_{g}^{\pi_{\theta}}(s,a) \right]^{2} + \mathbf{E} \left[\hat{V}_{g}^{\pi_{\theta}}(s,a) \right]^{2} \right\} \\
\stackrel{(c)}{\leq} \frac{2G^{4}(1+\lambda)^{2}}{\mu_{F}^{2}(1-\gamma)^{4}} + \frac{32(1+\lambda)^{2}}{(1-\gamma)^{4}} \leq \frac{1}{(1-\gamma)^{4}} \left[\frac{2G^{4}}{\mu_{F}^{2}} + 32 \right] (1+\lambda_{\max})^{2} \tag{41}$$

Inequality (a) follows from the Cauchy-Schwarz inequality and the fact that $(a+b)^2 \leq 2(a^2+b^2)$ for any two numbers a,b. The same argument is also applied in inequality (b). Additionally, it uses Assumption 2 and the definition of $\omega_{\theta,\lambda}^*$. Finally, (c) is a consequence of Assumption 4, Lemma 2, and the following two bounds.

$$\mathbf{E}\left[\hat{Q}_g^{\pi_\theta}(s,a)\right]^2 \leq \frac{2}{(1-\gamma)^2}, \text{ and } \mathbf{E}\left[\hat{V}_g^{\pi_\theta}(s)\right]^2 \leq \frac{2}{(1-\gamma)^2}, \ \forall (s,a) \in \mathcal{S} \times \mathcal{A}, \forall g \in \{r,c\} \ \ (42)$$

To establish the first bound, note that $|\hat{Q}_g^{\pi_\theta}(s,a)|$ is assigned a value of at most (j+1) with probability $(1-\gamma)\gamma^j, \forall g \in \{r,c\}$. Therefore,

$$\mathbf{E}\left[\hat{Q}_g^{\pi_{\theta}}(s,a)\right]^2 \le \sum_{j=0}^{\infty} (1-\gamma)(j+1)^2 \gamma^j = \frac{(1+\gamma)}{(1-\gamma)^2} < \frac{2}{(1-\gamma)^2}$$
(43)

The second bound in (42) can be proven similarly. This concludes the lemma.

E Proof of Lemma 5

We establish Lemma 5 applying Corollary 2 of [8]. Note the following statements.

S1: The following quantities exist and are finite $\forall \theta \in \Theta$.

$$F_{\rho}(\theta) \triangleq \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi_{\theta}}} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \otimes \nabla_{\theta} \log \pi_{\theta}(a|s) \right], \tag{44}$$

$$G_{\rho}(\theta) \triangleq \mathbf{E}_{(s,a) \sim \nu_{\rho}^{\pi_{\theta}}} \left[\nabla_{\theta} \log \pi_{\theta}(a|s) \otimes \nabla_{\theta} \log \pi_{\theta}(a|s) \otimes \nabla_{\theta} \log \pi_{\theta}(a|s) \otimes \nabla_{\theta} \log \pi_{\theta}(a|s) \right]$$
(45)

 $\mathbf{S2}$: There exists σ^2 such that the following is obeyed $\forall \theta \in \Theta$ where $\omega_{\theta,\lambda}^*$ minimizes $L_{\nu_o^{\pi_\theta}}(\cdot,\theta,\lambda)$.

$$\mathbf{E}\left[\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega_{\theta,\lambda}^{*},\theta,\lambda)\otimes\hat{\nabla}_{\omega}L_{\nu_{\rho}^{\pi_{\theta}}}(\omega_{\theta,\lambda}^{*},\theta,\lambda)\right] \preccurlyeq \sigma^{2}F_{\rho}(\theta) \tag{46}$$

S3: There exists μ_F , G > 0 such that the following statements hold $\forall \theta \in \Theta$.

$$(a) F_{\rho}(\theta) \succcurlyeq \mu_F I_{\rm d},$$
 (47)

(b)
$$\mathbf{E}_{(s,a)\sim\nu_{\rho}^{\pi_{\theta}}} \left[\|\nabla_{\theta}\log \pi_{\theta}(a|s)\|^{2} \nabla_{\theta}\log \pi_{\theta}(a|s) \otimes \nabla_{\theta}\log \pi_{\theta}(a|s) \right] \preccurlyeq G^{2} F_{\rho}(\theta),$$
 (48)

$$(c) \mathbf{E}_{(s,a)\sim\nu_{\rho}^{\pi_{\theta}}} \left[\|\nabla_{\theta} \log \pi_{\theta}(a|s)\|_{F_{\rho}(\theta)^{\dagger}}^{2} \nabla_{\theta} \log \pi_{\theta}(a|s) \otimes \nabla_{\theta} \log \pi_{\theta}(a|s) \right] \preccurlyeq \frac{G^{2}}{\mu_{F}} F_{\rho}(\theta)$$
(49)

Statement S1 follows from Assumption 2 whereas S2 is a consequence of Lemma 4. Statement S3(a) is identical to Assumption 4, S3(b) results from Assumption 2, and finally, S3(c) follows from Assumption 2 and 4. We can, therefore, apply Corollary 2 of [8] with $\kappa = \tilde{\kappa} = G^2/\mu_F$ and deduce the following convergence result whenever $H > \bar{C}\sqrt{\kappa\tilde{\kappa}}\log(\sqrt{\mathrm{d}}\sqrt{\kappa\tilde{\kappa}})$ and the learning rates are set as $\alpha = \frac{3\sqrt{5}\sqrt{\kappa\tilde{\kappa}}}{1+3\sqrt{5}\kappa\tilde{\kappa}}, \beta = \frac{1}{9\sqrt{\kappa\tilde{\kappa}}}, \xi = \frac{1}{3\sqrt{5}\mu_F\sqrt{\kappa\tilde{\kappa}}}$, and $\delta = \frac{1}{5G^2}$.

$$\mathbf{E}\left[l_{k}(\omega_{k})\right] - l_{k}(\omega_{k}^{*}) \leq \frac{C}{2} \exp\left(-\frac{H}{20\sqrt{\kappa\tilde{\kappa}}}\right) \left[l_{k}(\mathbf{0}) - l_{k}(\omega_{k}^{*})\right] + 11\frac{\sigma^{2}d}{H},$$
where $l_{k}(\omega) \triangleq L_{\nu_{k}^{\pi_{\theta_{k}}}}(\omega, \theta_{k}, \lambda_{k}), \ \forall \omega \in \mathbb{R}^{d}$

$$(50)$$

The term, C is a universal constant. Note that $l_k(\omega_k^*) \geq 0$ and $l_k(\mathbf{0})$ is bounded above as follows.

$$l_k(\mathbf{0}) = \frac{1}{2} \mathbf{E}_{(s,a) \sim \nu_\rho^{\pi_{\theta_k}}} \left[\frac{1}{1 - \gamma} A_{\mathrm{L},\lambda_k}^{\pi_{\theta_k}}(s,a) \right]^2 \stackrel{(a)}{\leq} \frac{(1 + \lambda_{\max})^2}{2(1 - \gamma)^4}$$
 (51)

where (a) is a result of the fact that $|A_{\mathrm{L},\lambda}^{\pi_{\theta}}(s,a)| \leq (1+\lambda_{\mathrm{max}})/(1-\gamma), \forall (s,a) \in \mathcal{S} \times \mathcal{A}, \forall \theta \in \Theta,$ and $\forall \lambda \in \Lambda$. Combining (50), (51), and the fact that $l_k(\cdot)$ is μ_F -strongly convex, we establish,

$$\mathbf{E}\|\omega_{k} - \omega_{k}^{*}\|^{2} \leq \frac{2}{\mu_{F}} \left[\mathbf{E} \left[l_{k}(\omega_{k}) \right] - l_{k}(\omega_{k}^{*}) \right] \leq 22 \frac{\sigma^{2} d}{\mu_{F} H} + C \exp \left(-\frac{\mu_{F}}{20G^{2}} H \right) \left[\frac{(1 + \lambda_{\max})^{2}}{\mu_{F} (1 - \gamma)^{4}} \right]$$
(52)

This proves the first statement. We get the following for noiseless ($\sigma^2=0$) gradient updates.

$$\mathbf{E}\|(\mathbf{E}[\omega_k|\theta_k] - \omega_k^*)\|^2 \le C \exp\left(-\frac{\mu_F}{20G^2}H\right) \left[\frac{(1+\lambda_{\max})^2}{\mu_F(1-\gamma)^4}\right]$$
(53)

The second statement can be established from (53) by applying Jensen's inequality on the function $f(x) = x^2$.

F Proof of Theorem 1

Applying the inequalities $H \ge 1$, $\exp(-(\mu_F/20G^2)H) \le 1$ and substituting the values of σ^2 and η (as stated in Lemma 4 and Theorem 1 respectively), we can rewrite (29) as follows.

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left(J_{\mathrm{L},\rho}(\pi^*, \lambda_k) - J_{\mathrm{L},\rho}(\theta_k, \lambda_k) \right) \\
\leq \sqrt{\epsilon_{\mathrm{bias}}} + f_0 \frac{(1 + \lambda_{\mathrm{max}})}{(1 - \gamma)^2} \exp\left(-\frac{\mu_F}{40G^2} H \right) + f_1 \frac{(1 + \lambda_{\mathrm{max}})}{(1 - \gamma)^2} \frac{1}{\sqrt{K}}$$
(54)

The terms f_0 and f_1 are defined below.

$$f_{0} \triangleq \frac{G\sqrt{C}}{\sqrt{\mu_{F}}},$$

$$f_{1} \triangleq B \left[\frac{44d}{\mu_{F}} \left(\frac{G^{4}}{\mu_{F}^{2}} + 16 \right) + \frac{C}{\mu_{F}} + \frac{G^{2}}{\mu_{F}^{2}} \right] + \mathbf{E}_{s \sim d_{\rho}^{\pi^{*}}} \left[KL(\pi^{*}(\cdot|s) \| \pi_{\theta_{0}}(\cdot|s)) \right]$$

$$(55)$$

Using the definition of $J_{L,\rho}(\cdot,\cdot)$, one can write,

$$J_{L,\rho}(\pi^*, \lambda_k) - J_{L,\rho}(\theta_k, \lambda_k) = (J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k)) + \lambda_k (J_{c,\rho}^{\pi^*} - J_{c,\rho}(\theta_k))$$

$$\stackrel{(a)}{\geq} (J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k)) + \lambda_k (-J_{c,\rho}(\theta_k))$$
(56)

where (a) follows from the fact that $\lambda_k \geq 0$ and $J_{c,\rho}^* \geq 0$ due to feasibility. Combining (54) and (56), we obtain the following.

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k) \right] + \frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[-\lambda_k J_{c,\rho}(\theta_k) \right] \\
\leq \sqrt{\epsilon_{\text{bias}}} + f_0 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \exp\left(-\frac{\mu_F}{40G^2} H \right) + f_1 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \frac{1}{\sqrt{K}}$$
(57)

F.1 Convergence Rate of the Objective Function

Note that (57) can be alternatively written as,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k) \right] \le \sqrt{\epsilon_{\text{bias}}} + f_0 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \exp\left(-\frac{\mu_F}{40G^2} H \right) + f_1 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \frac{1}{\sqrt{K}} + \frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[\lambda_k J_{c,\rho}(\theta_k) \right]$$
(58)

To obtain the convergence rate of the objective function, we need to bound the last term in the above expression. Observe the following chain of inequalities.

$$0 \le (\lambda_K)^2 \stackrel{(a)}{=} \sum_{k=0}^{K-1} \left((\lambda_{k+1})^2 - (\lambda_k)^2 \right)$$

$$\stackrel{(b)}{\le} \sum_{k=0}^{K-1} \left(\left[\lambda_k - \zeta \hat{J}_{c,\rho}(\theta_k) \right]^2 - (\lambda_k)^2 \right) = -2\zeta \sum_{k=0}^{K-1} \lambda_k \hat{J}_{c,\rho}(\theta_k) + \zeta^2 \sum_{k=0}^{K-1} \hat{J}_{c,\rho}^2(\theta_k)$$
(59)

where (a) uses $\lambda_0 = 0$ and (b) follows from the contraction property of the projection operator, \mathcal{P}_{Λ} . Using the above inequality, one can write,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[\lambda_k \hat{J}_{c,\rho}(\theta_k) \right] \le \frac{\zeta}{2K} \sum_{k=0}^{K-1} \mathbf{E} \left[\hat{J}_{c,\rho}^2(\theta_k) \right]$$
 (60)

Using the unbiasedness of $\hat{J}_{c,\rho}(\theta_k)$ (Lemma 1), we deduce the following.

$$\mathbf{E} \left[\lambda_k \hat{J}_{c,\rho}(\theta_k) \right] \stackrel{(a)}{=} \mathbf{E} \left[\lambda_k \mathbf{E} \left[\hat{J}_{c,\rho}(\theta_k) \middle| \theta_k \right] \right] = \mathbf{E} \left[\lambda_k J_{c,\rho}(\theta_k) \right]$$
(61)

where (a) is a consequence of the fact that $\hat{J}_{c,\rho}(\theta_k)$ and λ_k are conditionally independent given θ_k . Note that, $\hat{J}_{c,\rho}^2(\theta_k)$ is assigned a value of at most $(j+1)^2$ with probability $(1-\gamma)\gamma^j$, $j\in\{0,1,\cdots\}$. Therefore, the RHS of (60) can be bounded as follows.

$$\mathbf{E}\Big[\hat{J}_{c,\rho}^{2}(\theta_{k})\Big] \leq \sum_{j=0}^{\infty} (1-\gamma)(j+1)^{2} \gamma^{j} = \frac{1+\gamma}{(1-\gamma)^{2}} < \frac{2}{(1-\gamma)^{2}}$$
 (62)

Combining (58), (60), (61), and (62), we finally obtain,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k) \right]
\leq \sqrt{\epsilon_{\text{bias}}} + f_0 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \exp \left(-\frac{\mu_F}{40G^2} H \right) + f_1 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \frac{1}{\sqrt{K}} + \frac{\zeta}{(1 - \gamma)^2}
\stackrel{(a)}{=} \sqrt{\epsilon_{\text{bias}}} + f_0 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \exp \left(-\frac{\mu_F}{40G^2} H \right) + f_1 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \frac{1}{\sqrt{K}} + \frac{\lambda_{\text{max}}}{(1 - \gamma)} \frac{1}{\sqrt{K}}
\leq \sqrt{\epsilon_{\text{bias}}} + f_0 \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \exp \left(-\frac{\mu_F}{40G^2} H \right) + (f_1 + 1) \frac{(1 + \lambda_{\text{max}})}{(1 - \gamma)^2} \frac{1}{\sqrt{K}}$$
(63)

where (a) uses the substitution $\zeta = \lambda_{\max}(1 - \gamma)/\sqrt{K}$.

F.2 Rate of Constraint Violation

The following inequality is satisfied for any $k \in \{0, 1, \dots, K-1\}$.

$$|\lambda_{k+1} - \lambda_{\max}|^2 = |\mathcal{P}_{\Lambda}(\lambda_k - \zeta \hat{J}_{c,\rho}(\theta_k)) - \lambda_{\max}|^2$$

$$\stackrel{(a)}{\leq} |\lambda_k - \zeta \hat{J}_{c,\rho}(\theta_k) - \lambda_{\max}|^2 = |\lambda_k - \lambda_{\max}|^2 - 2\zeta \hat{J}_{c,\rho}(\theta_k) (\lambda_k - \lambda_{\max}) + \zeta^2 \hat{J}_{c,\rho}^2(\theta_k)$$
(64)

where (a) is due to the contractive property of \mathcal{P}_{Λ} . Performing an average over $k \in \{0, \dots, K-1\}$, and applying expectations on both sides, we get,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[(\lambda_k - \lambda_{\max}) \hat{J}_{c,\rho}(\theta_k) \right] \stackrel{(a)}{\leq} \frac{|\lambda_{\max}|^2 - |\lambda_K - \lambda_{\max}|^2}{2\zeta K} + \frac{\zeta}{2K} \sum_{k=0}^{K-1} \mathbf{E} \left[\hat{J}_{c,\rho}^2(\theta_k) \right] \\
\stackrel{(b)}{\leq} \frac{\lambda_{\max}^2}{2\zeta K} + \frac{\zeta}{(1-\gamma)^2} \stackrel{(c)}{=} \frac{3\lambda_{\max}}{2(1-\gamma)} \frac{1}{\sqrt{K}}$$
(65)

where (a) utilises $\lambda_0=0$, (b) applies (62), and (c) is derived using $\zeta=\lambda_{\max}(1-\gamma)/\sqrt{K}$. Note that one can write the following using Lemma 1 and the observation that $\hat{J}_{c,\rho}(\theta_k)$ and λ_k are conditionally independent given θ_k .

$$\mathbf{E}\left[(\lambda_k - \lambda_{\max})\hat{J}_{c,\rho}(\theta_k)\right] = \mathbf{E}\left[(\lambda_k - \lambda_{\max})\mathbf{E}\left[\hat{J}_{c,\rho}(\theta_k)\big|\theta_k\right]\right] = \mathbf{E}\left[(\lambda_k - \lambda_{\max})J_{c,\rho}(\theta_k)\right]$$
(66)

Combining (65) and (66), we get,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[(\lambda_k - \lambda_{\max}) J_{c,\rho}(\theta_k) \right] \le \frac{3\lambda_{\max}}{2(1-\gamma)} \frac{1}{\sqrt{K}} \le \frac{2(1+\lambda_{\max})}{(1-\gamma)^2 \sqrt{K}}$$
(67)

Finally, combining (58) and (67), we arrive at,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k) \right] + \lambda_{\max} \frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[-J_{c,\rho}(\theta_k) \right] \\
\leq \sqrt{\epsilon_{\text{bias}}} + f_0 \frac{(1 + \lambda_{\max})}{(1 - \gamma)^2} \exp\left(-\frac{\mu_F}{40G^2} H \right) + (f_1 + 2) \frac{(1 + \lambda_{\max})}{(1 - \gamma)^2} \frac{1}{\sqrt{K}}$$
(68)

Since the functions $\{J_{g,\rho}(\theta_k)\}$, $g \in \{r,c\}$, $k \in \{0,\cdots,K-1\}$ are linear in occupancy measure, there exists a policy $\bar{\pi}$ such that the following holds $\forall g \in \{r,c\}$.

$$\frac{1}{K} \sum_{k=0}^{K-1} J_{g,\rho}(\theta_k) = J_{g,\rho}^{\bar{\pi}}$$

This allows us to rewrite (68) as,

$$J_{r,\rho}^{\pi^*} - \mathbf{E}\left[J_{r,\rho}^{\bar{\pi}}\right] + \lambda_{\max} \mathbf{E}\left[-J_{c,\rho}^{\bar{\pi}}\right] \le \sqrt{\epsilon_{\text{bias}}} + f_0 \frac{(1+\lambda_{\max})}{(1-\gamma)^2} \exp\left(-\frac{\mu_F}{40G^2}H\right) + (f_1+2)\frac{(1+\lambda_{\max})}{(1-\gamma)^2} \frac{1}{\sqrt{K}}$$

$$(69)$$

Applying Lemma 9 and verifying (via Lemma 8) that $\lambda_{\text{max}} \geq 2\lambda^*$ where λ^* is the non-negative minimizer of the dual function corresponding to the unparameterized constrained optimization (2), we can write the constraint violation rate as follows.

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[-J_{c,\rho}(\theta_k) \right] \le \frac{2\sqrt{\epsilon_{\text{bias}}}}{\lambda_{\text{max}}} + 2f_0 \frac{(1+\lambda_{\text{max}})}{\lambda_{\text{max}}(1-\gamma)^2} \exp\left(-\frac{\mu_F}{40G^2}H\right) + 2(f_1+2) \frac{(1+\lambda_{\text{max}})}{\lambda_{\text{max}}(1-\gamma)^2} \frac{1}{\sqrt{K}}$$
(70)

F.3 Final Result

Substituting $\lambda_{\rm max}=2/[(1-\gamma)c_{\rm slater}]\geq 2$ in (63), we get the rate of convergence of the objective as follows.

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k) \right] \\
\leq \sqrt{\epsilon_{\text{bias}}} + \frac{3f_0 c_{\text{slater}}^{-1}}{(1-\gamma)^3} \exp\left(-\frac{\mu_F}{40G^2} H \right) + \frac{3(f_1+1)c_{\text{slater}}^{-1}}{(1-\gamma)^3} \frac{1}{\sqrt{K}}$$
(71)

Similarly, we obtain the constraint violation rate as,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[-J_{c,\rho}(\theta_k) \right] \\
\leq (1 - \gamma) c_{\text{slater}} \sqrt{\epsilon_{\text{bias}}} + \frac{3f_0}{(1 - \gamma)^2} \exp\left(-\frac{\mu_F}{40G^2} H \right) + \frac{3(f_1 + 2)}{(1 - \gamma)^2} \frac{1}{\sqrt{K}}$$
(72)

Let, H and K are chosen as follows for an arbitrary $\epsilon > 0$.

$$H = \frac{40G^2}{\mu_F} \log \left(\max \left\{ \frac{6f_0 c_{\text{slater}}^{-1}}{(1 - \gamma)^3} \epsilon^{-1}, \frac{6f_0}{(1 - \gamma)^2} \epsilon^{-1} \right\} \right) = \mathcal{O}(\log(\epsilon^{-1})),$$

$$K = \max \left\{ \frac{36(f_1 + 1)^2 c_{\text{slater}}^{-2}}{(1 - \gamma)^6} \epsilon^{-2}, \frac{36(f_1 + 2)^2}{(1 - \gamma)^4} \epsilon^{-2} \right\} = \mathcal{O}((1 - \gamma)^{-6} \epsilon^{-2})$$
(73)

For the above choice of H and K, we have,

$$\frac{1}{K} \sum_{k=0}^{K-1} \mathbf{E} \left[J_{r,\rho}^{\pi^*} - J_{r,\rho}(\theta_k) \right] \le \sqrt{\epsilon_{\text{bias}}} + \epsilon,$$

$$\mathbf{E} \left[\frac{1}{K} \sum_{k=0}^{K-1} -J_{c,\rho}(\theta_k) \right] \le (1-\gamma) c_{\text{slater}} \sqrt{\epsilon_{\text{bias}}} + \epsilon$$
(74)

Note that the expected number of steps required in Algorithm 1 is $\mathcal{O}((1-\gamma)^{-1})$. Therefore, the sample complexity required to ensure (74) is $\mathcal{O}((1-\gamma)^{-1}HK) = \tilde{\mathcal{O}}((1-\gamma)^{-7}\epsilon^{-2})$. Finally, note that Lemma 5 requires $H > \bar{C}\frac{G^2}{\mu_F}\log\left(\sqrt{\mathrm{d}}\frac{G^2}{\mu_F}\right)$. This can be ensured if ϵ is sufficiently small.

G Strong Duality and Related Lemmas

Define the dual function associated with the unparameterized constrained optimization (2) as follows.

$$J_{\mathrm{D},\rho}^{\lambda} = \max_{\pi} \left\{ J_{r,\rho}^{\pi} + \lambda J_{c,\rho}^{\pi} \right\} \tag{75}$$

The following lemma formally describes the strong duality result.

Lemma 7. [40] [Lemma 3] If $\lambda^* \triangleq \arg\min_{\lambda \geq 0} J_{D,\rho}^{\lambda}$ and π^* is a solution of (2), then the following holds whenever Assumption I is true.

$$J_{r,\rho}^{\pi^*} = J_{\mathrm{D},\rho}^{\lambda^*} \tag{76}$$

It is to be mentioned that strong duality, in general, does not hold for the parameterized optimization (3). The following lemma established a bound on λ^* which becomes the foundation for choosing the value of λ_{\max} in Algorithm 2.

Lemma 8. [40][Lemma 3] Let λ^* be the optimal dual variable as defined in Lemma 7. The following inequalities hold where c_{slater} is defined in Assumption 1.

$$0 \le \lambda^* \le \frac{1}{(1 - \gamma)c_{\text{slater}}}$$

The following lemma is the main tool in decoupling the objective and constraint violation rates.

Lemma 9. Let Slater's condition (Assumption 1) hold. If $C \ge 2\lambda^*$ where λ^* is defined in Lemma 7 and $\bar{\pi}$ is a policy such that $J^{\pi^*}_{r,\rho} - J^{\bar{\pi}}_{r,\rho} + C[-J^{\bar{\pi}}_{c,\rho}] \le \zeta$ for some $\zeta > 0$ then

$$-J_{c,\rho}^{\bar{\pi}} \le \frac{2\zeta}{C} \tag{77}$$

Proof. Define the function $v(\cdot)$ as follows.

$$v(\tau) = \max_{\pi} \left\{ J_{r,\rho}^{\pi} \middle| J_{c,\rho}^{\pi} \ge \tau \right\}, \ \tau \in \mathbb{R}^{d}$$
 (78)

Let $au=J_{c,\rho}^{ar{\pi}}.$ Therefore, one can write the following chain of inequalities.

$$J_{r,\rho}^{\bar{\pi}} \stackrel{(a)}{\leq} v(\tau) = \max_{\pi} \left\{ J_{r,\rho}^{\pi} \middle| J_{c,\rho}^{\pi} \geq \tau \right\} \stackrel{(b)}{\leq} \max_{\pi} \left\{ J_{r,\rho}^{\pi} + \lambda^{*} (J_{c,\rho}^{\pi} - \tau) \middle| J_{c,\rho}^{\pi} \geq \tau \right\}$$

$$= \max_{\pi} \left\{ J_{r,\rho}^{\pi} + \lambda^{*} J_{c,\rho}^{\pi} \middle| J_{c,\rho}^{\pi} \geq \tau \right\} - \tau \lambda^{*}$$

$$\stackrel{(c)}{\leq} J_{D,\rho}^{\lambda^{*}} - \tau \lambda^{*} \stackrel{(d)}{=} J_{r,\rho}^{\pi^{*}} - \tau \lambda^{*}$$

$$(79)$$

where (a) follows from the definition of $v(\tau)$, (b) uses $\lambda^* \geq 0$, (c) follows from the definition of the dual function, and (d) is a result of the strong duality (Lemma 7). Utilizing (79), we get,

$$(C - \lambda^*)(-\tau) \le J_{r,\rho}^{\pi^*} - J_{r,\rho}^{\bar{\pi}} + C(-\tau) \le \zeta$$
(80)

This leads to $-\tau \le \zeta/(C-\lambda^*) \le 2\zeta/C$.

NeurIPS Paper Checklist

1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The abstract and introduction clearly state the paper's contribution and scope. Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals
 are not attained by the paper.

2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: Some limitations of the paper have been discussed in the conclusion section. Moreover, some assumptions are also accompanied by elaborate discussion on their limitations.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: We have clearly stated the required assumptions and an accompanying complete proof in the appendix for each theory result.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and crossreferenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: Our paper is primarily of theoretical nature and does not include experiments. Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
- (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
- (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
- (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: Our paper is primarily of theoretical nature and does not include experiments. Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so "No" is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (https://nips.cc/public/guides/CodeSubmissionPolicy) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: Our paper is primarily of theoretical nature and does not include experiments. Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental
 material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Our paper is primarily of theoretical nature and does not include experiments. Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error
 of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: Our paper is primarily of theoretical nature and does not include experiments. Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics https://neurips.cc/public/EthicsGuidelines?

Answer: [Yes]

Justification: Our research conforms, in every respect, to the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. **Broader Impacts**

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: Our work is primarily of theoretical nature and has no immediate societal impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: No high risk data or model have been used.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: No existing asset has been used in the paper.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

 If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: No new asset is introduced in the paper.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: No experiments with human subjects were conducted.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: We conducted no experiments with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.