Provably Efficient Safe Exploration via Primal-Dual Policy Optimization

Dongsheng Ding University of Southern California dongshed@usc.edu Xiaohan Wei Facebook, Inc. ubimeteor@fb.com Zhuoran Yang Princeton University zy6@princeton.edu

Zhaoran Wang Northwestern University zhaoranwang@gmail.com Mihailo R. Jovanović
University of Southern California
mihailo@usc.edu

${f Abstract}$

We study the safe reinforcement learning problem using the constrained Markov decision processes in which an agent aims to maximize the expected total reward subject to a safety constraint on the expected total value of a utility function. We focus on an episodic setting with the function approximation where the Markov transition kernels have a linear structure but do not impose any additional assumptions on the sampling model. Designing safe reinforcement learning algorithms with provable computational and statistical efficiency is particularly challenging under this setting because of the need to incorporate both the safety constraint and the function approximation into the fundamental exploitation/exploration tradeoff. To this end, we present an Optimistic Primal-Dual Proximal Policy OPtimization (OPDOP) algorithm where the value function is estimated by combining the least-squares policy evaluation and an additional bonus term for safe exploration. We prove that the proposed algorithm achieves an $\tilde{O}(dH^{2.5}\sqrt{T})$ regret and an $\widetilde{O}(dH^{2.5}\sqrt{T})$ constraint violation, where d is the dimension of the feature mapping, His the horizon of each episode, and T is the total number of steps. These bounds hold when the reward/utility functions are fixed but the feedback after each episode is ban-

Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS) 2021, San Diego, California, USA. PMLR: Volume 130. Copyright 2021 by the author(s).

dit. Our bounds depend on the capacity of the state-action space only through the dimension of the feature mapping and thus our results hold even when the number of states goes to infinity. To the best of our knowledge, we provide the first provably efficient online policy optimization algorithm for constrained Markov decision processes in the function approximation setting, with safe exploration.

1 Introduction

Reinforcement Learning (RL) studies how an agent learns to maximize its expected total reward by interacting with an unknown environment over time [60]. Safe RL augments RL with a practical consideration of safety to deal with restrictions/constraints arising from real-world problems [33, 6, 28]. Examples include collision-avoidance in autonomous robots [31, 32], cost limitations in medical applications [34, 11], and legal and business restrictions in financial management [2]. A standard environment model for safe RL is the Constrained Markov Decision Processes (CMDPs) [5] that generalize the classical MDPs to maximizing the expected total reward under a safety-related constraint on the expected total utility [3, 65]. presence of constraints makes the fundamental exploration/exploitation trade-off more challenging.

There is considerable growth in safe RL, especially those studies on CMDPs, e.g., constrained policy gradient [63, 59], Lagrangian-based actor-critic [15, 14, 61, 46, 74], constrained policy optimization [3, 72, 78], primal-dual policy optimization [53, 52]. A key highlight of their developments is the successful integration of the constrained optimization and the policy-based RL for addressing constraints. Notwithstanding many successes, these safe RL algorithms either do not have

a convergence theory or are limited to asymptotic convergence. In practice, only a finite amount of data is available. Hence, it is imperative to design safe RL algorithms with computational and statistical efficiency guarantees. For this purpose, we must address the exploration/exploitation trade-off under constraints.

In this work, we look at the challenging problem of finding a sequence of policies in response to online streaming samples of transition, reward functions, and utility functions. We attempt to provide theoretical guarantees on the regret of an algorithm approaching the best policy in hindsight, and feasibility region determined by constraints. The task of safe exploration is to explore the unknown environment and learn to adapt the policy to the constraint set. Our problem setting deviates from existing scenarios, where good priors on constraints or transition models are more focused, e.g., references [62, 13, 25, 66, 23, 24, 65]. Recent policy-based safe RL algorithms for CMDPs, e.g., constrained policy optimization [3, 72, 78] and primaldual policy optimization [53, 52], seek a single safe policy via the constrained policy optimization whose sample efficiency guarantees do not have a theory.

In this paper, we aim to answer a theoretical question:

Can we design a provably sample efficient online policy optimization algorithm for CMDPs in the function approximation setting?

Contribution. We propose a provably efficient safe RL algorithm for CMDPs with an unknown transition model in the linear episodic setting – an Optimistic Primal-Dual Proximal Policy OPtimization (OPDOP) algorithm – where the value function is estimated by combining the least-squares policy evaluation and an additional bonus term for safe exploration. Theoretically, we prove that the proposed algorithm achieves an $O(dH^{2.5}\sqrt{T})$ regret and the same $O(dH^{2.5}\sqrt{T})$ constraint violation, where d is the dimension of the feature mapping, H is the horizon of each episode, and Tis the total number of steps. We establish these bounds in the setting where the reward/utility functions are fixed but the feedback after each episode is bandit. Our bounds depend on the capacity of the state space only through the dimension of the feature mapping and thus hold even when the number of states goes to infinity. To the best of our knowledge, our result is the first provably efficient online policy optimization for CMDPs in the function approximation setting, with safe exploration.

Related Work. Our work is related to a line of provably efficient RL algorithms based on the linear function approximation, e.g., references [70, 71, 37, 20, 76]. Using the optimism in the face of uncertainty [7, 19], these references address the exploration/exploitation

trade-off by adding the Upper Confidence Bound (UCB) bonus, and proposed algorithms are provably sample-efficient. A closely-related reference [20] connects policy-based RL with optimism, and proposes an optimistic proximal policy optimization with UCB exploration. However, all these references only study some particular MDPs in unconstrained RL problems. Additional efforts need to pay for making them work for CMDPs. Our work seeks to design an optimistic variant of proximal policy optimization for CMDPs. For the large CMDPs with unknown transition models, there is a line of literature that is related to the policy optimization under constraints, e.g., references [63, 3, 72, 61, 48, 78]. However, the exploration under constraints is less studied and their theoretical guarantees are unknown. Our work fills in this gap.

The study of RL algorithms for CMDPs has received growing attention, especially those on learning CMDPs with unknown transitions and rewards. As we know, most of them are model-based and only apply to finite state-action spaces. References [58, 29] leverage upper confidence bound (UCB) on fixed reward, utility, and transition probability to propose sampleefficient algorithms for tabular CMDPs; reference [58] establishes an $\widetilde{O}(\sqrt{|\mathcal{A}|T^{1.5}\log T})$ regret and constraint violation via linear program in the average-cost case in time T; reference [29] achieves an $O(|\mathcal{S}|\sqrt{H^3T})$ regret and constraint violation in the episodic setting via linear program and primal-dual policy optimization, where S is a state-space, A is an action space, and H is a fixed horizon of episode. In reference [55], the authors study an adversarial stochastic shortest path problem under constraints and unknown transitions with $O(|\mathcal{S}|\sqrt{|\mathcal{A}|}H^2T)$ regret and constraint violation. Reference [10] extends Q-learning with optimism for finite state-action CMDPs with peak constraints. Reference [18] proposes UCB-based convex planning for episodic tabular CMDPs in dealing with convex or hard constraints. References [40, 35] establish probably approximately correct (PAC) guarantees that enjoy better problem-dependent sample-complexity. In contrast, our proposed algorithm can potentially apply to scenarios with infinite state-space, and our provided sublinear regret and constraint violation bounds only depend on the implicit dimension instead of the true dimension of the state space. Compared to more recent references [26, 69, 21, 77], our development attacks the exploration directly and does not rely on any policy 'simulators' (or generative models).

2 Problem Setup

We consider an episodic Markov decision process (MDP) – MDP($\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r$) – where \mathcal{S} is a state

space, \mathcal{A} is an action space, H is a fixed length of each episode, $\mathbb{P} = \{\mathbb{P}_h\}_{h=1}^H$ is a collection of transition probability measures, and $r = \{r_h\}_{h=1}^H$ is a collection of reward functions. We assume that \mathcal{S} is a measurable space with possibly infinite number of elements. Moreover, for each step $h \in [H]$, $\mathbb{P}_h(\cdot | x, a)$ is a transition kernel over next state if action a is taken for state x and r_h : $\mathcal{S} \times \mathcal{A} \to [0,1]$ is a reward function. The constrained MDP – CMDP($\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r, g$) – additionally contains utility functions $g = \{g_h\}_{h=1}^H$ where g_h : $\mathcal{S} \times \mathcal{A} \to [0,1]$. We assume that reward/utility functions are deterministic. Our analysis readily generalizes to the setting where reward/utility are random.

Let the policy space $\Delta(\mathcal{A} \mid \mathcal{S}, H)$ be $\{\{\pi_h(\cdot \mid \cdot)\}_{h=1}^H$: $\pi_h(\cdot | x) \in \Delta(\mathcal{A}), \ \forall x \in \mathcal{S} \text{ and } h \in [H] \}, \text{ where } \Delta(\mathcal{A})$ denotes a probability simplex over the action space. Let $\pi^k \in \Delta(\mathcal{A} | \mathcal{S}, H)$ be a policy taken by the agent at episode k, where $\pi_h^k(\cdot | x_h^k)$: $\mathcal{S} \to \mathcal{A}$ is the action that the agent takes at state x_h^k . For simplicity, we assume the initial state x_1^k to be fixed as x_1 in different episodes for brevity. The agent interacts with the environment in the kth episode as follows. At the beginning, the agent determines a policy π^k . Then, at each step $h \in [H]$, the agent observes the state $x_h^k \in \mathcal{S}$, determines an action a_h^k following the policy $\pi_h^k(\cdot | x_h^k)$, and receives a reward $r_h(x_h^k, a_h^k)$ together with an utility $g_h(x_h^k, a_h^k)$. Meanwhile, the MDP evolves into next state x_{h+1}^k drawing from the probability $\mathbb{P}_h(\cdot | x_h^k, a_h^k)$. The episode terminates at state x_H^k in which no control action is taken and both reward and utility functions are equal to zero. In this paper, we focus a bandit setting where the agent only observes the values of reward/utility functions, $r_h(x_h^k, a_h^k)$, $g_h(x_h^k, a_h^k)$, at visited state-action pair (x_h^k, a_h^k) . We assume that reward/utility functions are fixed over episodes.

Given a policy $\pi \in \Delta(\mathcal{A} \mid \mathcal{S}, H)$, the value function $V_{r,h}^{\pi}$ associated with the reward function r at each step h are the expected values of total rewards,

$$V_{r,h}^{\pi}(x) = \mathbb{E}_{\pi} \left[\sum_{i=h}^{H} r_i(x_i, a_i) \, \middle| \, x_h = x \right]$$

for all $x \in \mathcal{S}$, $h \in [H]$, where the expectation \mathbb{E}_{π} is taken over the random state-action sequence $\{(x_h, a_h)\}_{h=i}^H$; the action a_h follows the policy $\pi_h(\cdot|x_h)$ at the state x_h and the next state x_{h+1} follows the transition dynamics $\mathbb{P}_h(\cdot|x_h, a_h)$. Thus, the action-value function $Q_{r,h}^{\pi}(x, a) \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ associated with the reward function r is the expected value of total rewards when the agent starts from state-action pair (x, a) at step h and follows policy π ,

$$Q_{r,h}^{\pi}(x,a) = \mathbb{E}_{\pi} \left[\sum_{i=h}^{H} r_i(x_i, a_i) \, \big| \, x_h = x, a_h = a \right]$$

for all $(x,a) \in \mathcal{S} \times \mathcal{A}$ and $h \in [H]$. Similarly, we define the value function $V_{g,h}^{\pi} \colon \mathcal{S} \to \mathbb{R}$ and the action-value function $Q_{g,h}^{\pi}(x,a) \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ associated with the utility function g. Denote symbol $\diamond = r$ or g. For brevity, we take the shorthand $\mathbb{P}_h V_{\diamond,h+1}^{\pi}(x,a) := \mathbb{E}_{x' \sim \mathbb{P}_h(\cdot \mid x,a)} V_{\diamond,h+1}^{\pi}(x')$. The Bellman equations associated with a policy π are given by

$$Q_{\diamond,h}^{\pi}(x,a) = \left(\diamond_h + \mathbb{P}_h V_{\diamond,h+1}^{\pi}\right)(x,a) \tag{1}$$

where $V_{\diamond,h}^{\pi}(x) = \langle Q_{\diamond,h}^{\pi}(x,\cdot), \pi_h(\cdot|x) \rangle_{\mathcal{A}}$, for all $(x,a) \in \mathcal{S} \times \mathcal{A}$. Here, the inner product of a function $f: \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ with $\pi(\cdot|x) \in \Delta(\mathcal{A})$ at fixed $x \in \mathcal{S}$ represents $\langle f(x,\cdot), \pi(\cdot|x) \rangle_{\mathcal{A}} := \sum_{a \in \mathcal{A}} \langle f(x,a), \pi(a|x) \rangle$.

2.1 Learning Performance

The learning agent aims to find a solution of a constrained problem in which the objective function is the expected total rewards and the constraint is on the expected total utilities,

$$\underset{\pi \in \Delta(\mathcal{A} \mid \mathcal{S}, H)}{\text{maximize}} V_{r,1}^{\pi}(x_1) \text{ subject to } V_{g,1}^{\pi}(x_1) \geq b \quad (2)$$

where we take $b \in (0, H]$ to avoid triviality. It is readily generalized to the problem with multiple constraints. Let $\pi^* \in \Delta(A \mid S, H)$ be a solution to problem (2). Since the policy π^* is computed from knowing the transition model and all reward and utility functions, we refer it as an optimal policy in-hindsight.

The associated Lagrangian of problem (2) is given by

$$\mathcal{L}(\pi, Y) := V_{r,1}^{\pi}(x_1) + Y(V_{q,1}^{\pi}(x_1) - b)$$

where π is the primal variable and $Y \geq 0$ is the dual variable. We can cast (2) into a saddle-point problem,

$$\max_{\pi \in \Delta(\mathcal{A} \mid \mathcal{S}, H)} \text{ minimize } \mathcal{L}(\pi, Y)$$

where $\mathcal{L}(\pi, Y)$ is convex in Y and is non-concave in π in general. To address the non-concavity, we will exploit the structure of value functions to propose a variant of Lagrange multipliers method for constrained RL problems in Section 3, which warrants a new line of primal-dual mirror descent type analysis in sequel. This distinguishes from unconstrained RL, e.g., [4, 20].

Another key feature of constrained RL is the safe exploration under constraints [33]. Without any constraint information $a\ priori$, it is infeasible for each policy to satisfy the constraint since utility information on constraints is only revealed after a policy is decided. Instead, we allow each policy to violate the constraint in each episode and minimize regret while minimizing total constraint violations for safe exploration over K episodes. We define the regret as the

difference between the total reward value of policy π^* in hindsight and that of the agent's policy π^k over K episodes, and the constraint violation as a difference between the offset Kb and the total utility value of the agent's policy π^k over K episodes,

$$\operatorname{Regret}(K) = \sum_{k=1}^{K} \left(V_{r,1}^{\pi^{\star}}(x_{1}) - V_{r,1}^{\pi^{k}}(x_{1}) \right)$$

$$\operatorname{Violation}(K) = \sum_{k=1}^{K} \left(b - V_{g,1}^{\pi^{k}}(x_{1}) \right).$$
(3)

In this paper, we design algorithms, taking bandit feedback of the reward/utility functions, with both regret and constraint violation being sublinear in the total number of steps T := HK. Put differently, the algorithm should ensure that given $\epsilon > 0$, if $T = O(1/\epsilon^2)$, then both $\operatorname{Regret}(K) = O(\epsilon)$ and $\operatorname{Violation}(K) = O(\epsilon)$ hold with high probability.

Let $\mathcal{D}(Y) := \operatorname{maximize}_{\pi} \mathcal{L}(\pi, Y)$ be the dual function and $Y^* := \operatorname{argmin}_{Y \geq 0} \mathcal{D}(Y)$ be the optimal dual variable. We assume feasibility for problem (2) in Assumption 1 that is known as the Slater condition [53, 29, 55]. It is convenient to establish the strong duality [53] and the boundedness of the optimal dual variable Y^* that can be found in Appendix E.

Assumption 1 (Feasibility). There exists $\gamma > 0$ and $\bar{\pi} \in \Delta(\mathcal{A} | \mathcal{S}, H)$ such that $V_{a,1}^{\bar{\pi}}(x_1) \geq b + \gamma$.

Lemma 1 (Strong Duality, Boundedness of Y^*). Let Assumption 1 hold. Then $V_{r,1}^{\pi^*}(x_1) = \mathcal{D}(Y^*)$. Moreover, $0 \leq Y^* \leq (V_{r,1}^{\pi^*}(x_1) - V_{r,1}^{\bar{\pi}}(x_1))/\gamma$.

Lemma 1 provides useful optimization properties of (2) for our algorithm design and analysis.

2.2 Linear Function Approximation

We focus on a class of CMDPs, where transition kernels are linear in feature maps.

Assumption 2. The CMDP($\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r, g$) is a linear MDP with a kernel feature map $\psi \colon \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}^{d_1}$, if for any $h \in [H]$, there exists a vector $\theta_h \in \mathbb{R}^{d_1}$ with $\|\theta_h\|_2 \leq \sqrt{d_1}$ such that for any $(x, a, x') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$,

$$\mathbb{P}_h\left(x'\,|\,x,a\right) = \langle \psi\left(x,a,x'\right),\theta_h\rangle;$$

there exists a feature map $\varphi \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}^{d_2}$ and vectors $\theta_{r,h}, \theta_{g,h} \in \mathbb{R}^{d_2}$ such that for any $(x,a) \in \mathcal{S} \times \mathcal{A}$,

$$r_h(x,a) = \langle \varphi(x,a), \theta_{r,h} \rangle$$
 and $g_h(x,a) = \langle \varphi(x,a), \theta_{a,h} \rangle$

where $\max(\|\theta_{r,h}\|_2, \|\theta_{g,h}\|_2) \leq \sqrt{d_2}$. Moreover, we assume that for any function $V \colon \mathcal{S} \to [0, H]$, $\|\int_{\mathcal{S}} \psi(x, a, x') V(x') dx'\|_2 \leq \sqrt{d_1} H$ for all $(x, a) \in \mathcal{S} \times \mathcal{A}$, and $\max(d_1, d_2) \leq d$.

Assumption 2 adapts the definition of linear kernel MDP [8, 79, 20] for CMDPs. Linear kernel MDP examples include tabular MDPs [79], feature embedded transition models [71], and linear combinations of base models [50]. We can construct related examples of CMDPs with linear structure by adding adding proper constraints. For usefulness of linear structure, see discussions in the literature [27, 64, 43]. For more general transition dynamics, see factored MDPs [54].

Although our definition in Assumption 2 and linear MDPs [70, 37] all contain tabular MDPs as special cases, they define transition dynamics using different feature maps. They are not comparable since one cannot be implied by the other [79]. We provide more details on the tabular case of Assumption 2 in Section 5.

3 Proposed Algorithm

In Algorithm 1, we present a new variant of proximal policy optimization [57] – an Optimistic Primal—Dual Proximal Policy OPtimization (OPDOP) algorithm. We effectuate the optimism through the Upper-Confidence Bounds (UCB) [70, 71, 37], and address the constraints via the union of the Lagrange multipliers method with the value function structure that is captured by the performance difference lemma [38, 20].

Lemma 2 (Performance Difference Lemma). For any two policies π , $\pi' \in \Delta(\mathcal{A} | \mathcal{S}, H)$, $\diamond = r$ or q,

$$V_{\diamond,1}^{\pi'}(x_1) - V_{\diamond,1}^{\pi}(x_1)$$

$$= \mathbb{E}_{\pi'} \left[\sum_{h=1}^{H} \left\langle Q_{\diamond,h}^{\pi}(x_h,\cdot), \pi'_h(\cdot \mid x_h) - \pi_h(\cdot \mid x_h) \right\rangle \mid x_1 \right].$$

In each episode, our algorithm consists of three main stages. The first stage (lines 4–8) is $Policy\ Improvement$: we receive a new policy π^k by improving previous π^{k-1} via a mirror descent type optimization; The second stage (line 9) is $Dual\ Update$: we update dual variable Y^k based on the constraint violation induced by previous policy π^k ; The third stage (line 10) is $Policy\ Evaluation$: we optimistically evaluate newly obtained policy via the least-squares policy evaluation with an additional UCB bonus term for exploration.

3.1 Policy Improvement

In the k-th episode, a natural attempt of obtaining a policy π^k is to solve a Lagrangian-based policy optimization problem,

$$\max_{\pi \in \Delta(\mathcal{A}|\mathcal{S}, H)} \mathcal{L}(\pi, Y^{k-1}) := V_{r, 1}^{\pi}(x_1) - Y^{k-1}(b - V_{g, 1}^{\pi}(x_1))$$

where $\mathcal{L}(\pi, Y)$ is the Lagrangian and the dual variable $Y^{k-1} \geq 0$ is from the last episode; we show that Y^{k-1}

can be updated efficiently in Section 3.2. This type update also finds in references [45, 53, 52, 61]. They rely on an oracle solver, e.g., Q-learning [30], proximal policy optimization [57], or trust region policy optimization [56], to deliver a near-optimal policy, making overall algorithmic complexity expensive. Hence, they are not suitable for online use. In contrast, this work utilizes RL problem structure and shows that only an easily-computable proximal step is sufficient for efficiently achieving near-optimal performance.

Recall symbol $\diamond = r$ or g. Via the performance difference lemma, we can expand value functions $V_{\diamond,1}^{\pi}(x_1)$ at the previously known policy π^{k-1} ,

$$V_{\diamond,1}^{\pi}(x_1) = V_{\diamond,1}^{\pi^{k-1}}(x_1^k) + \\ \mathbb{E}_{\pi^{k-1}} \left[\sum_{h=1}^{H} \left\langle Q_{\diamond,h}^{\pi}(x_h,\cdot), (\pi_h - \pi_h^{k-1})(\cdot \mid x_h) \right\rangle \right]$$

where $\mathbb{E}_{\pi^{k-1}}$ is taken over the random state-action sequence $\{(x_h, a_h)\}_{h=1}^H$. Thus, we introduce an approximation of $V_{\diamond,1}^{\pi}(x_1)$ for any state-action sequence $\{(x_h, a_h)\}_{h=1}^H$ induced by π ,

$$L_{\diamond}^{k-1}(\pi) = V_{\diamond,1}^{k-1}(x_1) + \sum_{h=1}^{H} \left\langle Q_{\diamond,h}^{k-1}(x_h,\cdot), (\pi_h - \pi_h^{k-1})(\cdot \mid x_h) \right\rangle$$

where $V_{\diamond,h}^{k-1}$ and $Q_{\diamond,h}^{k-1}$ can be estimated from an optimistic policy evaluation that will be discussed in Section 3.3. With this notion, in each episode, instead of solving a Lagrangian-based policy optimization, we perform a simple policy update in online mirror descent fashion,

$$\max_{\pi \in \Delta(\mathcal{A}|\mathcal{S}, H)} L_r^{k-1}(\pi) - Y^{k-1}(b - L_g^{k-1}(\pi))$$
$$-\frac{1}{\alpha} \sum_{h=1}^{H} D(\pi_h(\cdot \mid x_h) \mid \widetilde{\pi}_h^{k-1}(\cdot \mid x_h))$$

where $\widetilde{\pi}_h^{k-1}(\cdot | x_h) = (1-\theta) \pi_h^{k-1}(\cdot | x_h) + \theta \operatorname{Unif}(\mathcal{A})$ is a mixed policy of the previous one and the uniform distribution $\operatorname{Unif}(\mathcal{A})$ with $\theta \in (0,1]$. The constant $\alpha > 0$ is a trade-off parameter, $D(\pi | \widetilde{\pi}^{k-1})$ is the KL divergence between π and $\widetilde{\pi}^{k-1}$ in which π is absolutely continuous in $\widetilde{\pi}^{k-1}$. The policy mixing step ensures such absolute continuity and implies uniformly bounded KL divergence; see Lemma 15 in Appendix F. Ignoring other π -irrelevant terms, we update π^k in terms of previous policy π^{k-1} by

$$\underset{\pi \in \Delta(\mathcal{A}|\mathcal{S}, H)}{\operatorname{argmax}} \sum_{h=1}^{H} \left\langle (Q_{r,h}^{k-1} + Y^{k-1} Q_{g,h}^{k-1})(x_h, \cdot), \pi_h(\cdot \mid x_h) \right\rangle$$
$$-\frac{1}{\alpha} \sum_{h=1}^{H} D(\pi_h(\cdot \mid x_h) \mid \widetilde{\pi}_h^{k-1}(\cdot \mid x_h)).$$

Since the above update is separable over H steps, we can update the policy π^k as line 6 in Algorithm 1, a closed-form solution for any step $h \in [H]$. If we set $Y^{k-1} = 0$ and $\theta = 0$, the above update reduces to one step in an optimistic proximal policy optimization [20]. The idea of KL-divergence regularization in policy optimization has been widely used in many unconstrained scenarios [39, 57, 56, 67, 47]. Our method is distinct in that it is based on the performance difference lemma and the optimistically estimated value functions.

Algorithm 1 Optimistic Primal-Dual Proximal Policy OPtimization (OPDOP)

- 1: **Initialization**: Let $\{Q_{r,h}^0, Q_{g,h}^0\}_{h=1}^H$ be zero functions, $\{\pi_h^0\}_{h\in[H]}$ be uniform distributions on \mathcal{A} , $V_{g,1}^0$ be b, Y^0 be 0, χ be $2H/\gamma$, $\alpha, \eta > 0$, $\theta \in (0,1]$.
- 2: **for** episode k = 1, ..., K + 1 **do**
- 3: Set the initial state $x_1^k = x_1$.
- 4: **for** step h = 1, 2, ..., H **do**
- 5: Mix the policy

$$\widetilde{\pi}_h^{k-1}(\cdot|\cdot) = (1-\theta)\pi_h^{k-1}(\cdot|\cdot) + \theta \operatorname{Unif}(\mathcal{A}).$$

6: Update the policy

$$\pi_h^k(\cdot|\cdot) \propto \widetilde{\pi}_h^{k-1}(\cdot|\cdot) \operatorname{e}^{\left(\alpha\left(Q_{r,h}^{k-1} + Y^{k-1}Q_{g,h}^{k-1}\right)(\cdot,\cdot)\right)}.$$

- 7: Take an action $a_h^k \sim \pi_h^k(\cdot | x_h^k)$ and recieve reward/utility, $r_h(x_h^k, a_h^k)$, $g_h(x_h^k, a_h^k)$.
- 8: Observe the next state x_{h+1}^k .
- 9: Update the dual variable Y^k by

$$Y^{k} = \operatorname{Proj}_{[0,\chi]}(Y^{k-1} + \eta(b - V_{g,1}^{k-1}(x_{1}))).$$

10: Estimate the action-value or value functions $\{Q_{r,h}^k(\cdot,\cdot),Q_{q,h}^k(\cdot,\cdot),V_{q,h}^k(\cdot)\}_{h=1}^H$ via

$$\text{LSTD}\Big(\{x_h^{\tau}, a_h^{\tau}, r_h(x_h^{\tau}, a_h^{\tau}), g_h(x_h^{\tau}, a_h^{\tau})\}_{h, \tau = 1}^{H, k}\Big).$$

3.2 Dual Update

To infer the constraint violation for the dual update, we estimate $V_{g,1}^{\pi^k}(x_1)$ via an optimistic policy evaluation by $V_{g,1}^{k-1}(x_1)$ that is discussed in Section 3.3. We update the Lagrange multiplier Y by moving Y^k to the direction of minimizing the Lagrangian $\mathcal{L}(\pi,Y)$ over $Y \geq 0$ in line 9 of Algorithm 1, where $\eta > 0$ is a stepsize and $\operatorname{Proj}_{[0,\chi]}$ is a projection onto $[0,\chi]$ with an upper bound χ on Y^k . By Lemma 1, we choose $\chi = 2H/\gamma \geq 2Y^*$ so that projection interval $[0,\chi]$ includes the optimal dual variable Y^* . This type design also finds in references [29,51].

The dual update works as a trade-off between the reward maximization and the constraint violation reduction. If the current policy π^k satisfies the approximated constraint, i.e., $b-L_g^{k-1}(\pi^k) \leq 0$, we put less weight on the action-value function associated with the utility and maximize the reward; otherwise, we sacrifice the reward a bit to satisfy the constraint. The dual update has a similar use in dealing with constraints in CMDPs, e.g., Lagrangian-based actor-critic [22, 46], and online constrained optimization [73, 68, 75]. In contrast, we handle the dual update via the optimistic policy evaluation, yielding a simple, but efficient estimation on the constraint violation.

Algorithm 2 Least-Squares Temporal Difference (LSTD) with UCB exploration

- 1: **Input**: $\{x_h^{\tau}, a_h^{\tau}, r_h(x_h^{\tau}, a_h^{\tau}), g_h(x_h^{\tau}, a_h^{\tau})\}_{h, \tau = 1}^{H, k}$.
- 2: Initialization: Set $\{V_{r,H+1}^k, V_{g,H+1}^k\}$ be zero functions and $\lambda = 1, \beta = O(\sqrt{dH^2 \log{(dT/p)}})$.
- 3: for step $h=H,H-1,\cdots,1$ do $\Rightarrow \Rightarrow =r,g$

4:
$$\Lambda_{\diamond,h}^{k} = \sum_{\tau=1}^{k-1} \phi_{\diamond,h}^{\tau}(x_{h}^{\tau}, a_{h}^{\tau}) \phi_{\diamond,h}^{\tau}(x_{h}^{\tau}, a_{h}^{\tau})^{\top} + \lambda I.$$

5:
$$w_{\diamond,h}^k = (\Lambda_{\diamond,h}^k)^{-1} \sum_{\tau=1}^{k-1} \phi_{\diamond,h}^{\tau}(x_h^{\tau}, a_h^{\tau}) V_{\diamond,h+1}^{\tau}(x_{h+1}^{\tau}).$$

6:
$$\phi_{\diamond,h}^k(\cdot,\cdot) = \int_{\mathcal{S}} \psi(\cdot,\cdot,x') V_{\diamond,h+1}^k(x') dx'.$$

7:
$$\Gamma^{k}_{\diamond,h}(\cdot,\cdot) = \beta(\phi^{k}_{\diamond,h}(\cdot,\cdot)^{\top}(\Lambda^{k}_{\diamond,h})^{-1}\phi^{k}_{\diamond,h}(\cdot,\cdot))^{1/2}.$$

8:
$$\Lambda_h^k = \sum_{\tau=1}^{k-1} \varphi(x_h^{\tau}, a_h^{\tau}) \varphi(x_h^{\tau}, a_h^{\tau})^{\top} + \lambda I.$$

9:
$$u_{\diamond,h}^k = (\Lambda_h^k)^{-1} \sum_{\tau=1}^{k-1} \varphi(x_h^{\tau}, a_h^{\tau}) \diamond_h (x_h^{\tau}, a_h^{\tau}).$$

10:
$$\Gamma_h^k(\cdot,\cdot) = \beta(\varphi(\cdot,\cdot)^{\top}(\Lambda_h^k)^{-1}\varphi(\cdot,\cdot))^{1/2}$$

11:
$$Q_{\diamond,h}^{k}(\cdot,\cdot) = \min\left(H - h + 1, \, \varphi(\cdot,\cdot)^{\top} u_{\diamond,h}^{k} + \phi_{\diamond,h}^{k}(\cdot,\cdot)^{\top} w_{\diamond,h}^{k} + (\Gamma_{h}^{k} + \Gamma_{\diamond,h}^{k})(\cdot,\cdot)\right)^{+}.$$

12:
$$V_{\diamond,h}^k(\cdot) = \langle Q_{\diamond,h}^k(\cdot,\cdot), \pi_h^k(\cdot|\cdot) \rangle_{\mathcal{A}}.$$

13: Return: $\{Q_{\diamond,h}^k(\cdot,\cdot),V_{\diamond,h}^k(\cdot,\cdot)\}_{h=1}^H$.

3.3 Policy Evaluation

The last stage of the kth episode takes the Least-Squares Temporal Difference (LSTD) [17, 16, 44, 42] to evaluate the policy π^k based on previous k-1 historical trajectories. For each step $h \in [H]$, instead of $\mathbb{P}_h V_{r,h+1}^{\pi^k}$ in the Bellman equations (1), we estimate $\mathbb{P}_h V_{r,h+1}^k$ by $(\phi_{r,h}^k)^\top w_{r,h}^k$ where $w_{r,h}^k$ is updated by the minimizer of the regularized least-squares problem over w,

$$\sum_{\tau=1}^{k-1} \left(V_{r,h+1}^{\tau}(x_{h+1}^{\tau}) - \phi_{r,h}^{\tau}(x_{h}^{\tau}, a_{h}^{\tau})^{\top} w \right)^{2} + \lambda \|w\|_{2}^{2}$$
 (4)

where $\phi_{r,h}^{\tau}(\cdot,\cdot) := \int_{\mathcal{S}} \psi(\cdot,\cdot,x') V_{r,h+1}^{\tau}(x') dx',$ $V_{r,h+1}^{\tau}(\cdot) = \langle Q_{r,h+1}^{\tau}(\cdot,\cdot), \pi_{h+1}^{\tau}(\cdot|\cdot) \rangle_{\mathcal{A}}$ for $h \in [H-1]$ and $V_{H+1}^{\tau} = 0$, and $\lambda > 0$ is the regularization parameter. Similarly, we estimate $\mathbb{P}_h V_{g,h+1}^k$ by $(\phi_{g,h}^k)^{\top} w_{g,h}^k$. We display the least-squares solution in line 4–6 of Algorithm 2, where symbol $\diamond = r$ or g. We also estimate $r_h(\cdot,\cdot)$ by $\varphi^{\top} u_{r,h}^k$, where $u_{r,h}^k$ is updated by the minimizer of another regularized least-squares problem,

$$\sum_{\tau=1}^{k-1} \left(r_h(x_h^{\tau}, a_h^{\tau}) - \varphi(x_h^{\tau}, a_h^{\tau})^{\top} u \right)^2 + \lambda \|u\|_2^2$$
 (5)

where $\lambda > 0$ is the regularization parameter. Similarly, we estimate $g_h(\cdot,\cdot)$ by $\varphi^{\top}u_{g,h}^k$. The least-squares solutions lead to line 8–9 of Algorithm 2.

After obtaining estimates of $\mathbb{P}_h V_{\diamond,h+1}^k$ and $\diamond_h(\cdot,\cdot)$ for $\diamond = r$ or g, we update the estimated action-value function $\{Q_{\diamond,h}^k\}_{h=1}^H$ iteratively in line 11 of Algorithm 2, where $\varphi^\top u_{\diamond,h}^k$ is an estimate of \diamond_h and $(\phi_{\diamond,h}^k)^\top w_{\diamond,h}^k$ is an estimate of $\mathbb{P}_h V_{\diamond,h+1}^k$; we add UCB bonus terms $\Gamma_h^k(\cdot,\cdot), \Gamma_{\diamond,h}^k(\cdot,\cdot)$: $\mathcal{S} \times \mathcal{A} \to \mathbb{R}^+$ so that

$$\varphi^{\top} u_{\diamond,h}^k + \Gamma_h^k$$
 and $(\phi_{\diamond,h}^k)^{\top} w_{\diamond,h}^k + \Gamma_{\diamond,h}^k$

all become their upper confidence bounds. Here, the bonus terms take $\Gamma_h^k = \beta(\varphi^\top (\Lambda_h^k)^{-1}\varphi)^{1/2}$ and $\Gamma_{\diamond,h}^k = \beta((\phi_{\diamond,h}^k)^\top (\Lambda_{\diamond,h}^k)^{-1}\phi_{\diamond,h}^k)^{1/2}$ and we leave the parameter $\beta>0$ to be tuned later. Moreover, the bounded reward/utility $\diamond_h \in [0,1]$ implies $Q_{\diamond,h}^k \in [0,H-h+1]$.

We remark the computational efficiency of Algorithm 1. For the time complexity, since line 6 is a scalar update, they need $O(d|\mathcal{A}|T)$ time. A dominating calculation is from lines 5/9 in Algorithm 2. If we use the Sherman–Morrison formula for computing $(\Lambda_h^k)^{-1}$, it takes $O(d^2T)$ time. Another important calculation is the integration from line 6 in Algorithm 2. We can either compute it analytically if it is tractable or approximate it via the Monte Carlo integration [79] that assumes polynomial time. Therefore, the time complexity is $O(\text{poly}(d)|\mathcal{A}|T)$ in total. For the space complexity, we don't need to store policy since it is recursively calculated via line 6 of Algorithm 1. By updating Y^k , Λ_h^k , $\Lambda_{\diamond,h}^k$, $w_{\diamond,h}^k$, $u_{\diamond,h}^k$, and $\diamond_h(x_h^k, a_h^k)$ recursively, it takes $O((d^2 + |\mathcal{A}|)H)$ space.

4 Regret and Constraint Violation Analysis

We now prove that the regret and the constraint violation for Algorithm 1 are sublinear in T := KH, the total number of steps taken by the algorithm, where K is the total number of episodes and H is the episode horizon. We recall that $|\mathcal{A}|$ is the cardinality of action space \mathcal{A} and d is the dimension of the feature map.

Theorem 1 (Linear Kernal MDP: Regret and Constraint Violation). Let Assumptions 1 and 2 hold. Fix $p \in (0,1)$. We set $\alpha = \sqrt{\log |\mathcal{A}|}/(H^2K)$, $\beta = C_1\sqrt{dH^2\log(dT/p)}$, $\eta = 1/\sqrt{K}$, $\theta = 1/K$, and $\lambda = 1$ in Algorithm 1, where C_1 is an absolute constant. Suppose $\log |\mathcal{A}| = O\left(d^2\log^2(dT/p)\right)$. Then, with probability 1-p, the regret and the constraint violation in (3) satisfy

$$\begin{aligned} \operatorname{Regret}(K) & \leq C \, dH^{2.5} \sqrt{T} \, \log \left(\frac{dT}{p} \right) \\ \left[\operatorname{Violation}(K) \right]_{+} & \leq C' \, dH^{2.5} \sqrt{T} \log \left(\frac{dT}{p} \right) \end{aligned}$$

where C and C' are absolute constants.

The above result establishes that Algorithm 1 enjoys an $\widetilde{O}(dH^{2.5}\sqrt{T})$ regret and an $\widetilde{O}(dH^{2.5}\sqrt{T})$ constraint violation if we set algorithm parameters $\{\alpha, \beta, \eta, \theta, \lambda\}$ properly. Our results have the optimal dependence on the total number of steps T up to some logarithmic factors. The d dependence occurs due to the uniform concentration for controlling the fluctuations in the least-squares policy evaluation. This matches the existing bounds in the linear MDP setting without any constraints [20, 8, 79]. Our bounds differ from them only by H dependence, which is a price introduced by the uniform bound on the constraint violation. It is noticed that our algorithm works for bandit feedback of reward/utility functions after each episode.

Regarding safe exploration, our violation bound provides finite-time convergence to the feasibility region defined by constraints. In the interaction with an unknown environment, the UCB exploration in the utility value function adds optimism towards constraint satisfaction. The dual update regularizes the policy improvement for governing actual constraint violation. Our regret and violation bounds readily lead to PAC guarantees [36]. Compared to most recent references [26, 69, 21, 77], our algorithm is sample-efficient in exploration and does not need simulations of policy.

We remark the tabular setting for Algorithm 1; see Appendix C for details. The tabular CMDP is a special case of Assumption 2 by taking canonical bases as feature mappings; see them in Section 5. The feature map has dimension $d = |\mathcal{S}|^2 |\mathcal{A}|$ and thus Theorem 1 automatically provides $O(|\mathcal{S}|^2 |\mathcal{A}| H^{2.5} \sqrt{T})$ regret and constraint violation for the tabular CMDPs. The $d = |\mathcal{S}|^2 |\mathcal{A}|$ dependence relies on the least-squares policy evaluation and it can be improved via other optimistic policy evolution methods if we are limited to the tabular case. We provide such results in Section 5.

4.1 Proof Outline of Theorem 1

We sketch the proof for Theorem 1. We state key lemmas and delay their full versions and proofs to Appendix B. In what follows, we fix $p \in (0,1)$ and use the shorthand w.p. for with probability.

Regret Analysis. We take a regret decomposition,

$$\operatorname{Regret}(K) = \underbrace{\sum_{k=1}^{K} \left(V_{r,1}^{\pi^{*}}(x_{1}) - V_{r,1}^{k}(x_{1}) \right)}_{(R.I)} + \underbrace{\sum_{k=1}^{K} \left(V_{r,1}^{k}(x_{1}) - V_{r,1}^{\pi^{k}}(x_{1}) \right)}_{(R.II)}$$

where π^* is an optimal policy in hindsight, and $V_{r,1}^k(x_1)$ is estimated via our optimistic policy evaluation given by Algorithm 2. Since we use $V_{r,h+1}^k$ to estimate $V_{r,h+1}^{\pi^k}$, it leads a model prediction error in the Bellman equations, $\iota_{r,h}^k := r_h^k + \mathbb{P}_h V_{r,h+1}^k - Q_{r,h}^k$; similarly define $\iota_{g,h}^k$. In Appendix D.3, the UCB optimism of $\iota_{\phi,h}^k$ with $\diamond = r$ or g, shows that or any $(k,h) \in [K] \times [H]$ and $(x,a) \in \mathcal{S} \times \mathcal{A}$, w.p. 1-p/2, we have

$$-2(\Gamma_h^k + \Gamma_{\diamond h}^k)(x,a) \le \iota_{\diamond h}^k(x,a) \le 0.$$

By assumptions of Theorem 1, the policy improvement in line 6 of Algorithm 1 yields Lemma 1, depicting weighted total differences of estimates $V_{r,1}^k(x_1)$, $V_{q,1}^k(x_1)$ to the optimal ones.

Lemma 1 (Policy Improvement: Primal-Dual Mirror Descent Step). Let assumptions of Theorem 1 hold. Then.

$$(R.I) \leq -\sum_{k=1}^{K} Y^{k} \left(V_{g,1}^{\pi^{\star}}(x_{1}) - V_{g,1}^{k}(x_{1}) \right) + \sum_{k=1}^{K} \sum_{h=1}^{H} \mathbb{E}_{\pi^{\star}} \left[\iota_{r,h}^{k}(x_{h}, a_{h}) + Y^{k} \iota_{g,h}^{k}(x_{h}, a_{h}) \right] + O(H^{2.5} \sqrt{T \log |\mathcal{A}|}).$$

Lemma 1 displays coupling between the regret (R.I) and the constraint. This coupling also finds in online convex optimization [49, 75, 68, 41] and CMDP problems [29]. The proof of Lemma 1 takes a primal-dual mirror descent type analysis of line 6 of Algorithm 1, using the performance difference lemma.

Via the dual update in line 9 of Algorithm 1, we can verify that the second total differences $-\sum_{k=1}^{K} Y^k \left(V_{g,1}^{\pi^*}(x_1) - V_{g,1}^k(x_1)\right)$ scales $O(\sqrt{K})$. Together with a decomposition of (R.II),

(R.II) =
$$-\sum_{k=1}^{K} \sum_{h=1}^{H} \iota_{r,h}^{k}(x_{h}^{k}, a_{h}^{k}) + M_{r,H,2}^{K}$$

where $M_{r,H,2}^K$ is a martingale, we now have Lemma 2.

Lemma 2. Let assumptions of Theorem 1 hold. Then,

Regret(K)
$$\leq \sum_{k=1}^{K} \sum_{h=1}^{H} (\mathbb{E}_{\pi^{\star}}[\iota_{r,h}^{k}(x_{h}, a_{h})] - \iota_{r,h}^{k}(x_{h}^{k}, a_{h}^{k})) + M_{r,H,2}^{K} + O(H^{2.5}\sqrt{T\log|\mathcal{A}|}).$$

Finally, we note that $M_{r,H,2}^K$ is a martingale that scales as $O(H\sqrt{T})$ via the Azuma-Hoeffding inequality. For the model prediction error, we use the UCB optimism and apply the elliptical potential lemma.

Lemma 3. Let assumptions of Theorem 1 hold. Then,

$$\begin{split} & \sum_{k=1}^{K} \sum_{h=1}^{H} \left(\mathbb{E}_{\pi^{\star}} [\iota_{r,h}^{k}(x_{h}, a_{h})] - \iota_{r,h}^{k}(x_{h}^{k}, a_{h}^{k}) \right) \\ & \leq O \left(dH^{1.5} \sqrt{T \log (K) \log (dT/p)} \right), \text{ w.p. } 1 - p/2. \end{split}$$

Lemma 4. Let assumptions of Theorem 1 hold. Then,

$$|M_{r,H,2}^K| \le 4H\sqrt{T\log(4/p)}, \text{ w.p. } 1 - p/2.$$

Finally, we apply probability bounds from Lemmas 3 and 4 to Lemma 2 to get our regret bound.

Constraint Violation Analysis. We take a violation decomposition,

Violation(K) =
$$\sum_{k=1}^{K} (b - V_{g,1}^{k}(x_{1}))$$

+ $\sum_{k=1}^{K} (V_{g,1}^{k}(x_{1}) - V_{g,1}^{\pi^{k}}(x_{1}))$.

We begin with the policy improvement in line 6 of Algorithm 1 to refine Lemma 1 as Lemma 5.

Lemma 5 (Policy Improvement: Refined Primal-Dual Mirror Descent Step). Let assumptions of Theorem 1 hold. Then, for any $Y \in [0, \chi]$,

(R.I) +
$$Y \sum_{k=1}^{K} (b - V_{g,1}^{k}(x_1)) \le O(H^{2.5} \sqrt{T \log |\mathcal{A}|}).$$

Lemma 5 removes the dual update Y^k in the second total differences in Lemma 1. We prove Lemma 5 by combining Lemma 1 with the UCB optimism and a change of variable of Y^k for the dual update.

Similar to (R.II), we also have

(V.II) =
$$-\sum_{k=1}^{K} \sum_{h=1}^{H} \iota_{g,h}^{k}(x_{h}^{k}, a_{h}^{k}) + M_{g,H,2}^{K}$$

where $M_{g,H,2}^{K}$ is a martingale. By adding (V.II) to the inequality in Lemma 5 with multiplier $Y \geq 0$, and also adding (R.II) to it,

$$\sum_{k=1}^{K} \left(V_{r,1}^{\pi^{\star}}(x_{1}) - V_{r,1}^{\pi^{k}}(x_{1}) \right) + Y \sum_{k=1}^{K} \left(b - V_{g,1}^{\pi^{k}}(x_{1}) \right)$$

$$\leq - \sum_{k=1}^{K} \sum_{h=1}^{H} \left(\iota_{r,h}^{k}(x_{h}^{k}, a_{h}^{k}) + Y \iota_{g,h}^{k}(x_{h}^{k}, a_{h}^{k}) \right)$$

$$+ O\left(H^{2.5} \sqrt{T \log |\mathcal{A}|} \right) + M_{r,H,2}^{K} + Y M_{g,H,2}^{K}$$

Then, we take Y=0 if $\sum_{k=1}^K \left(b-V_{g,1}^{\pi^k}(x_1)\right) \leq 0$; otherwise $Y=\chi,$ w.p. 1-p, we have,

$$\left(V_{r,1}^{\pi^*}(x_1) - V_{r,1}^{\pi'}(x_1)\right) + \chi \left[b - V_{g,1}^{\pi'}(x_1)\right]_{+} \\
\leq O\left(dH^{2.5}\sqrt{T}\log\left(dT/p\right)/K\right)$$

where $V_{r,1}^{\pi'}(x_1) = \frac{1}{K} \sum_{k=1}^{K} V_{r,1}^{\pi^k}(x_1)$ and $V_{g,1}^{\pi'}(x_1) =$ $\frac{1}{K}\sum_{k=1}^{K}V_{g,1}^{\pi^k}(x_1)$ for some existing policy π' . Here, we bound $\Gamma_h^k + \Gamma_{\diamond,h}^k$ and $M_{\diamond,H,2}^K$ as done in Lemmas 3

Last, by the strong duality in Lemma 1, we apply the constraint violation bound from constrained optimization that is stated in Lemma 10 in Appendix E,

$$[b - V_{q,1}^{\pi'}(x_1)]_+ \le O(dH^{2.5}\sqrt{T}\log(dT/p)/(\chi K))$$

which gives our desired violation bound.

Algorithm 3 Optimistic Policy Evaluation (OPE)

- 1: **Input**: $\{x_h^{\tau}, a_h^{\tau}, r_h(x_h^{\tau}, a_h^{\tau}), g_h(x_h^{\tau}, a_h^{\tau})\}_{h, \tau = 1}^{H, k}$.
- 2: **Initialization**: Set $\{V_{r,H+1}^k, V_{q,H+1}^k\}$ be zero functions, and $\lambda = 1$, $\beta = C_1 H \sqrt{|\mathcal{S}| \log(|\mathcal{S}||\mathcal{A}|T/p)}$.
- 3: **for** step $h=H,H-1,\cdots,1$ **do** $\Rightarrow \Rightarrow =r,g$ 4: Compute counters $n_h^k(x,a,x')$ and $n_h^k(x,a)$ via (7) for all $(x,a,x')\in\mathcal{S}\times\mathcal{A}\times\mathcal{S}$ and $(x,a)\in$
- Estimate reward/utility functions \widehat{r}_h^k , \widehat{g}_h^k via (8) for all $(x, a) \in \mathcal{S} \times \mathcal{A}$.
- Estimate transition $\widehat{\mathbb{P}}_h^k$ via (9) for all $(x, a, x') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, and take bonus $\Gamma_h^k = \beta \left(n_h^k(x, a) + \lambda\right)^{-1/2}$ for all $(x, a) \in \mathcal{S} \times \mathcal{A}$.
- $Q_{h}^{k}(\cdot,\cdot) = \min(H-h+1,\widehat{\diamond}_{h}^{k}(\cdot,\cdot)+$ $\sum_{x' \in \mathcal{S}} \widehat{\mathbb{P}}_h(x' \mid \cdot, \cdot) V_{\diamond, h+1}^k(x') + 2\Gamma_h^k(\cdot, \cdot), \)^+ \cdot V_{\diamond, h}^k(\cdot) = \left\langle Q_{\diamond, h}^k(\cdot, \cdot), \pi_h^k(\cdot \mid \cdot) \right\rangle_{\mathcal{A}}.$
- 9: **Return**: $\{Q_{r,h}^k(\cdot,\cdot), Q_{q,h}^k(\cdot,\cdot)\}_{h=1}^H$.

5 Further Results on Tabular Case

The tabular CMDP($\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r, g$) is a special case of Assumption 2 with $|\mathcal{S}|<\infty$ and $|\mathcal{A}|<\infty$. Let $d_1 = |\mathcal{S}|^2 |\mathcal{A}|$ and $d_2 = |\mathcal{S}| |\mathcal{A}|$. We take the following feature maps $\psi(x, a, x') \in \mathbb{R}^{d_1}$, $\varphi(x, a) \in \mathbb{R}^{d_2}$, and parameter vectors,

$$\psi(x, a, x') = \mathbf{e}_{(x, a, x')}, \ \theta_h = \mathbb{P}_h(\cdot, \cdot, \cdot)
\varphi(x, a) = \mathbf{e}_{(x, a)}, \ \theta_{r, h} = r_h(\cdot, \cdot), \ \theta_{q, h} = g_h(\cdot, \cdot)$$
(6)

where $\mathbf{e}_{(x,a,x')}$ is a canonical basis of \mathbb{R}^{d_1} associated with (x,a,x') and $\theta_h = \mathbb{P}_h(\cdot,\cdot,\cdot)$ reads that for any $(x,a,x') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, the (x,a,x')th entry of θ_h is $\mathbb{P}(x'|x,a)$; similarly we define $\mathbf{e}_{(x,a)}$, $\theta_{r,h}$, and $\theta_{g,h}$. We can verify that $\|\theta_h\| \leq \sqrt{d_1}$, $\|\theta_{r,h}\| \leq \sqrt{d_2}$, $\|\theta_{g,h}\| \leq \sqrt{d_2}$, and for any V: $\mathcal{S} \to [0,H]$ and any $(x,a) \in \mathcal{S} \times \mathcal{A}$, we have $\|\sum_{x' \in \mathcal{S}} \psi(x,a,x')V(x')\| \leq \sqrt{|\mathcal{S}|}H \leq \sqrt{d_1}H$. Therefore, we take $d := \max(d_1,d_2) = |\mathcal{S}|^2|\mathcal{A}|$ in Assumption (2) for the tabular case.

The proof of Theorem 1 is generic, since it is ready to achieve sublinear regret and constraint violation bounds as long as the policy evaluation is sample-efficient, e.g., the UCB design of 'optimism in the face of uncertainty.' In what follows, we introduce another efficient policy evaluation for line 10 of Algorithm 1 in the tabular case. Let us first introduce some notation. For any $(h,k) \in [H] \times [K]$, any $(x,a,x') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, and any $(x,a) \in \mathcal{S} \times \mathcal{A}$, we define two visitation counters $n_h^k(x,a,x')$ and $n_h^k(x,a)$ at step h in episode k,

$$n_h^k(x, a, x') = \sum_{\substack{\tau = 1 \\ k-1}}^{k-1} \mathbf{1}\{(x, a, x') = (x_h^{\tau}, a_h^{\tau}, a_{h+1}^{\tau})\}$$

$$n_h^k(x, a) = \sum_{\tau = 1}^{k-1} \mathbf{1}\{(x, a) = (x_h^{\tau}, a_h^{\tau})\}.$$

$$(7)$$

This allows us to estimate reward function r, utility function g, and transition kernel \mathbb{P}_h for episode k by

$$\widehat{r}_{h}^{k}(x,a) = \sum_{\tau=1}^{k-1} \frac{\mathbf{1}\{(x,a) = (x_{h}^{\tau}, a_{h}^{\tau})\} r_{h}(x_{h}^{\tau}, a_{h}^{\tau})}{n_{h}^{k}(x,a) + \lambda}$$

$$\widehat{g}_{h}^{k}(x,a) = \sum_{\tau=1}^{k-1} \frac{\mathbf{1}\{(x,a) = (x_{h}^{\tau}, a_{h}^{\tau})\} g_{h}(x_{h}^{\tau}, a_{h}^{\tau})}{n_{h}^{k}(x,a) + \lambda}$$
(8)

$$\widehat{\mathbb{P}}_h^k(x'|x,a) = \frac{n_h^k(x,a,x')}{n_h^k(x,a) + \lambda} \tag{9}$$

for all $(x, a, x') \in \mathcal{S} \times \mathcal{A} \times \mathcal{S}$, $(x, a) \in \mathcal{S} \times \mathcal{A}$ where $\lambda > 0$ is the regularization parameter. Moreover, we introduce the bonus term $\Gamma_h^k \colon \mathcal{S} \times \mathcal{A} \to \mathbb{R}$, $\Gamma_h^k(x, a) = \beta \left(n_h^k(x, a) + \lambda\right)^{-1/2}$ which adapts the counter-based bonus terms in the literature [9, 36], where $\beta > 0$ is to be determined later.

Using the estimated transition kernels $\{\widehat{\mathbb{P}}_h^k\}_{h=1}^H$, the estimated reward/utility functions $\{\widehat{r}_h^k, \widehat{g}_h^k\}_{h=1}^H$, and the bonus terms $\{\Gamma_h^k\}_{h=1}^H$, we now can estimate the action-value function via line 7 of Algorithm 3 for any $(x,a) \in \mathcal{S} \times \mathcal{A}$, where $\diamond = r$ or g. Thus,

 $V_{\diamond,h}^k(x) = \langle Q_{\diamond,h}^k(x,\cdot), \pi_h^k(\cdot | x) \rangle_{\mathcal{A}}$. We summarize the above procedure in Algorithm 3. Using already estimated $\{Q_{r,h}^k(\cdot,\cdot),Q_{g,h}^k(\cdot,\cdot)\}_{h=1}^H$, we execute the policy improvement and the dual update in Algorithm 1.

As in Theorem 1, we provide theoretical guarantees in Theorem 2; see Appendix C.2 for the proof. Theorem 2 improves ($|\mathcal{S}|, |\mathcal{A}|$) dependence in Theorem 1 for the tabular case and also matches $|\mathcal{S}|$ dependence in references [29, 55]. It is worthy mentioning our Algorithm 1 is generic in handling an infinite state space.

Theorem 2 (Tabular Case: Regret and Constraint Violation). Let Assumption 1 hold and let Assumption 2 hold with feature maps (6). Fix $p \in (0,1)$. In Algorithm 1, we set $\alpha = \sqrt{\log |\mathcal{A}|}/(H^2K)$, $\beta = C_1H\sqrt{|\mathcal{S}|\log(|\mathcal{S}||\mathcal{A}|T/p)}$, $\eta = 1/\sqrt{K}$, $\theta = 1/K$, and $\lambda = 1$ where C_1 is an absolute constant. Then, with probability 1-p, the regret and the constraint violation in (3) satisfy

$$\begin{aligned} \operatorname{Regret}(K) & \leq & C|\mathcal{S}|\sqrt{|\mathcal{A}|H^5T}\log\left(\frac{|\mathcal{S}||\mathcal{A}|T}{p}\right) \\ \left[\operatorname{Violation}(K)\right]_{+} & \leq & C'|\mathcal{S}|\sqrt{|\mathcal{A}|H^5T}\log\left(\frac{|\mathcal{S}||\mathcal{A}|T}{p}\right) \end{aligned}$$

where C and C' are absolute constants.

6 Concluding Remarks

We have developed a provably efficient safe reinforcement learning algorithm in the linear MDP setting. The algorithm extends the proximal policy optimization to CMDPs by incorporating the UCB exploration. We prove that the proposed algorithm achieves an $\widetilde{O}(\sqrt{T})$ regret and an $\widetilde{O}(\sqrt{T})$ constraint violation under mild conditions, where T is the total number of steps taken by the algorithm. Our algorithm works in the setting where reward/utility functions are given by bandit feedback. To the best of our knowledge, our algorithm is the first provably efficient online policy optimization algorithm for CMDPs in the function approximation setting.

Mathematically, our algorithm framework allows reward/utility functions to be adversarial. We believe that approaches from the adversarial MDP literature allow us to derive similar regret and constraint violation bounds, although we leave it as future work. Beyond linear kernel MDPs, the UCB exploration has previously been applied for other types of MDPs, e.g., factored MDPs, or infinite-horizon MDPs. It remains to be seen if these are extendable for CMDPs in similar settings. In practice, we often encounter general function approximation beyond linear functions, e.g., neural nets. It would be useful to design provably efficient exploration algorithms for CMDPs with general function approximation.

Acknowledgements

D. Ding and M. R. Jovanović were supported by the National Science Foundation under Awards ECCS-1708906 and ECCS-1809833.

References

- Y. Abbasi-Yadkori, D. Pál, and C. Szepesvári. Improved algorithms for linear stochastic bandits. In Advances in Neural Information Processing Systems, volume 24, pages 2312–2320, 2011.
- [2] N. Abe, P. Melville, C. Pendus, C. K. Reddy, D. L. Jensen, V. P. Thomas, J. J. Bennett, G. F. Anderson, B. R. Cooley, M. Kowalczyk, et al. Optimizing debt collections using constrained reinforcement learning. In *International Conference* on Knowledge Discovery and Data Mining, pages 75–84, 2010.
- [3] J. Achiam, D. Held, A. Tamar, and P. Abbeel. Constrained policy optimization. In *International Conference on Machine Learning*, volume 70, pages 22–31, 2017.
- [4] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. Optimality and approximation with policy gradient methods in Markov decision processes. arXiv preprint arXiv:1908.00261, 2019.
- [5] E. Altman. Constrained Markov Decision Processes, volume 7. CRC Press, 1999.
- [6] D. Amodei, C. Olah, J. Steinhardt, P. Christiano, J. Schulman, and D. Mané. Concrete problems in AI safety. arXiv preprint arXiv:1606.06565, 2016.
- [7] P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine learning*, 47(2):235–256, 2002.
- [8] A. Ayoub, Z. Jia, C. Szepesvari, M. Wang, and L. Yang. Model-based reinforcement learning with value-targeted regression. In *International Conference on Machine Learning*, pages 463–474, 2020.
- [9] M. G. Azar, I. Osband, and R. Munos. Minimax regret bounds for reinforcement learning. In International Conference on Machine Learning, pages 263–272, 2017.
- [10] Q. Bai, A. Gattami, and V. Aggarwal. Model-free algorithm and regret analysis for MDPs with peak constraints. arXiv preprint arXiv:2003.05555, 2020.

- [11] Y. Bai, T. Xie, N. Jiang, and Y.-X. Wang. Provably efficient Q-learning with low switching cost. In Advances in Neural Information Processing Systems, volume 32, pages 8002–8011, 2019.
- [12] A. Beck. First-order methods in optimization. SIAM, 2017.
- [13] F. Berkenkamp, M. Turchetta, A. Schoellig, and A. Krause. Safe model-based reinforcement learning with stability guarantees. In Advances in Neural Information Processing Systems, volume 30, pages 908–918, 2017.
- [14] S. Bhatnagar and K. Lakshmanan. An online actor–critic algorithm with function approximation for constrained Markov decision processes. *Journal of Optimization Theory and Applications*, 153(3):688–708, 2012.
- [15] V. S. Borkar. An actor-critic algorithm for constrained Markov decision processes. Systems & Control Letters, 54(3):207–213, 2005.
- [16] J. A. Boyan. Least-squares temporal difference learning. In *International Conference on Machine Learning*, pages 49–56, 1999.
- [17] S. J. Bradtke and A. G. Barto. Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22(1-3):33–57, 1996.
- [18] K. Brantley, M. Dudik, T. Lykouris, S. Miryoosefi, M. Simchowitz, A. Slivkins, and W. Sun. Constrained episodic reinforcement learning in concave-convex and knapsack settings. In Advances in Neural Information Processing Systems, volume 33, 2020.
- [19] S. Bubeck, N. Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. Foundations and Trends® in Machine Learning, 5(1):1–122, 2012.
- [20] Q. Cai, Z. Yang, C. Jin, and Z. Wang. Provably efficient exploration in policy optimization. In *International Conference on Machine Learning*, pages 1283–1294, 2020.
- [21] Y. Chen, J. Dong, and Z. Wang. A primal-dual approach to constrained Markov decision processes. arXiv preprint arXiv:2101.10895, 2021.
- [22] Y. Chow, M. Ghavamzadeh, L. Janson, and M. Pavone. Risk-constrained reinforcement learning with percentile risk criteria. The Journal of Machine Learning Research, 18(1):6070–6120, 2017.

- [23] Y. Chow, O. Nachum, E. Duenez-Guzman, and M. Ghavamzadeh. A lyapunov-based approach to safe reinforcement learning. In *Advances in Neu*ral Information Processing Systems, volume 31, pages 8092–8101, 2018.
- [24] Y. Chow, O. Nachum, A. Faust, M. Ghavamzadeh, and E. Duenez-Guzman. Lyapunov-based safe policy optimization for continuous control. arXiv preprint arXiv:1901.10031, 2019.
- [25] G. Dalal, K. Dvijotham, M. Vecerik, T. Hester, C. Paduraru, and Y. Tassa. Safe exploration in continuous action spaces. arXiv preprint arXiv:1801.08757, 2018.
- [26] D. Ding, K. Zhang, T. Basar, and M. Jovanovic. Natural policy gradient primal-dual method for constrained Markov decision processes. In Advances in Neural Information Processing Systems, volume 33, 2020.
- [27] S. S. Du, S. M. Kakade, R. Wang, and L. F. Yang. Is a good representation sufficient for sample efficient reinforcement learning? In *International Conference on Learning Representations*, 2019.
- [28] G. Dulac-Arnold, D. Mankowitz, and T. Hester. Challenges of real-world reinforcement learning. arXiv preprint arXiv:1904.12901, 2019.
- [29] Y. Efroni, S. Mannor, and M. Pirotta. Exploration-exploitation in constrained MDPs. arXiv preprint arXiv:2003.02189, 2020.
- [30] D. Ernst, P. Geurts, and L. Wehenkel. Treebased batch mode reinforcement learning. *Jour*nal of Machine Learning Research, 6(Apr):503– 556, 2005.
- [31] S. Feyzabadi. Robot Planning with Constrained Markov Decision Processes. PhD thesis, UC Merced, 2017.
- [32] J. F. Fisac, A. K. Akametalu, M. N. Zeilinger, S. Kaynama, J. Gillula, and C. J. Tomlin. A general safety framework for learning-based control in uncertain robotic systems. *IEEE Transactions* on Automatic Control, 64(7):2737–2752, 2018.
- [33] J. Garcia and F. Fernández. A comprehensive survey on safe reinforcement learning. *Journal* of Machine Learning Research, 16(1):1437–1480, 2015.
- [34] C. J. Girard. Structural Results for Constrained Markov Decision Processes. PhD thesis, Cornell University, 2018.

- [35] A. HasanzadeZonuzy, D. Kalathil, and S. Shakkottai. Learning with safety constraints: Sample complexity of reinforcement learning for constrained mdps. arXiv preprint arXiv:2008.00311, 2020.
- [36] C. Jin, Z. Allen-Zhu, S. Bubeck, and M. I. Jordan. Is Q-learning provably efficient? In Advances in Neural Information Processing Systems, volume 31, pages 4863–4873, 2018.
- [37] C. Jin, Z. Yang, Z. Wang, and M. I. Jordan. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143, 2020.
- [38] S. Kakade and J. Langford. Approximately optimal approximate reinforcement learning. In *Proceedings of the Nineteenth International Conference on Machine Learning*, page 267–274, 2002.
- [39] S. M. Kakade. A natural policy gradient. In Advances in Neural Information Processing Systems, volume 15, pages 1531–1538, 2002.
- [40] K. C. Kalagarla, R. Jain, and P. Nuzzo. A sample-efficient algorithm for episodic finitehorizon MDP with constraints. arXiv preprint arXiv:2009.11348, 2020.
- [41] A. Koppel, K. Zhang, H. Zhu, and T. Başar. Projected stochastic primal-dual method for constrained online learning with kernels. *IEEE Transactions on Signal Processing*, 67(10):2528–2542, 2019.
- [42] M. G. Lagoudakis and R. Parr. Least-squares policy iteration. *Journal of Machine Learning Research*, 4(Dec):1107–1149, 2003.
- [43] T. Lattimore, C. Szepesvari, and G. Weisz. Learning with good feature representations in bandits and in RL with a generative model. In *International Conference on Machine Learning*, pages 5662–5670, 2020.
- [44] A. Lazaric, M. Ghavamzadeh, and R. Munos. Finite-sample analysis of LSTD. In *International Conference on Machine Learning*, pages 615–622, 2010.
- [45] H. Le, C. Voloshin, and Y. Yue. Batch policy learning under constraints. In *International Conference on Machine Learning*, pages 3703–3712, 2019.
- [46] Q. Liang, F. Que, and E. Modiano. Accelerated primal-dual policy optimization for safe reinforcement learning. arXiv preprint arXiv:1802.06480, 2018.

- [47] B. Liu, Q. Cai, Z. Yang, and Z. Wang. Neural trust region/proximal policy optimization attains globally optimal policy. In *Advances in Neural In*formation Processing Systems, volume 32, pages 10564–10575, 2019.
- [48] Y. Liu, J. Ding, and X. Liu. IPO: Interior-point policy optimization under constraints. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4940–4947, 2020.
- [49] M. Mahdavi, R. Jin, and T. Yang. Trading regret for efficiency: online convex optimization with long term constraints. *The Journal of Machine Learning Research*, 13(1):2503–2528, 2012.
- [50] A. Modi, N. Jiang, A. Tewari, and S. Singh. Sample complexity of reinforcement learning using linearly combined model ensembles. In *International Conference on Artificial Intelligence and Statistics*, pages 2010–2020, 2020.
- [51] A. Nedić and A. Ozdaglar. Subgradient methods for saddle-point problems. *Journal of optimization theory and applications*, 142(1):205–228, 2009.
- [52] S. Paternain, M. Calvo-Fullana, L. F. Chamon, and A. Ribeiro. Safe policies for reinforcement learning via primal-dual methods. arXiv preprint arXiv:1911.09101, 2019.
- [53] S. Paternain, L. Chamon, M. Calvo-Fullana, and A. Ribeiro. Constrained reinforcement learning has zero duality gap. In Advances in Neural Information Processing Systems, volume 32, pages 7553–7563, 2019.
- [54] B. Á. Pires and C. Szepesvári. Policy error bounds for model-based reinforcement learning with factored linear models. In *Conference on Learning Theory*, pages 121–151, 2016.
- [55] S. Qiu, X. Wei, Z. Yang, J. Ye, and Z. Wang. Upper confidence primal-dual optimization: Stochastically constrained Markov decision processes with adversarial losses and unknown transitions. In Advances in Neural Information Processing Systems, volume 33, 2020.
- [56] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International Conference on Machine Learn*ing, pages 1889–1897, 2015.
- [57] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.

- [58] R. Singh, A. Gupta, and N. B. Shroff. Learning in Markov decision processes under constraints. arXiv preprint arXiv:2002.12435, 2020.
- [59] A. Stooke, J. Achiam, and P. Abbeel. Responsive safety in reinforcement learning by PID Lagrangian methods. In *International Conference on Machine Learning*, pages 9133–9143, 2020.
- [60] R. S. Sutton and A. G. Barto. Reinforcement Learning: An Introduction. MIT press, 2018.
- [61] C. Tessler, D. J. Mankowitz, and S. Mannor. Reward constrained policy optimization. In *Inter*national Conference on Learning Representations, 2019.
- [62] M. Turchetta, F. Berkenkamp, and A. Krause. Safe exploration in finite Markov decision processes with Gaussian processes. In Advances in Neural Information Processing Systems, volume 29, pages 4312–4320, 2016.
- [63] E. Uchibe and K. Doya. Constrained reinforcement learning from intrinsic and extrinsic rewards. In *International Conference on Develop*ment and Learning, pages 163–168, 2007.
- [64] B. Van Roy and S. Dong. Comments on the Du-Kakade-Wang-Yang lower bounds. arXiv preprint arXiv:1911.07910, 2019.
- [65] A. Wachi and Y. Sui. Safe reinforcement learning in constrained Markov decision processes. In International Conference on Machine Learning, pages 9797–9806, 2020.
- [66] A. Wachi, Y. Sui, Y. Yue, and M. Ono. Safe exploration and optimization of constrained MDPs using Gaussian processes. In AAAI Conference on Artificial Intelligence, 2018.
- [67] L. Wang, Q. Cai, Z. Yang, and Z. Wang. Neural policy gradient methods: Global optimality and rates of convergence. In *International Conference* on Learning Representations, 2019.
- [68] X. Wei, H. Yu, and M. J. Neely. Online primaldual mirror descent under stochastic constraints. Proceedings of the ACM on Measurement and Analysis of Computing Systems, 4(2):1–36, 2020.
- [69] T. Xu, Y. Liang, and G. Lan. A primal approach to constrained policy optimization: Global optimality and finite-time analysis. arXiv preprint arXiv:2011.05869, 2020.
- [70] L. Yang and M. Wang. Sample-optimal parametric Q-learning using linearly additive features. In International Conference on Machine Learning, pages 6995–7004, 2019.

- [71] L. Yang and M. Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine* Learning, pages 10746–10756, 2020.
- [72] T.-Y. Yang, J. Rosca, K. Narasimhan, and P. J. Ramadge. Projection-based constrained policy optimization. In *International Conference on Learning Representations*, 2019.
- [73] H. Yu, M. Neely, and X. Wei. Online convex optimization with stochastic constraints. In Advances in Neural Information Processing Systems, volume 30, pages 1428–1438, 2017.
- [74] M. Yu, Z. Yang, M. Kolar, and Z. Wang. Convergent policy optimization for safe reinforcement learning. In Advances in Neural Information Processing Systems, volume 32, pages 3121–3133, 2019.
- [75] J. Yuan and A. Lamperski. Online convex optimization for cumulative constraints. In Advances in Neural Information Processing Systems, volume 31, pages 6137–6146, 2018.
- [76] A. Zanette, A. Lazaric, M. Kochenderfer, and E. Brunskill. Learning near optimal policies with low inherent bellman error. In *International Con*ference on Machine Learning, pages 10978–10989, 2020.
- [77] J. Zhang, A. Koppel, A. S. Bedi, C. Szepesvari, and M. Wang. Variational policy gradient method for reinforcement learning with general utilities. In Advances in Neural Information Processing Systems, volume 33, 2020.
- [78] Y. Zhang, Q. Vuong, and K. Ross. First order constrained optimization in policy space. In Advances in Neural Information Processing Systems, volume 33, 2020.
- [79] D. Zhou, J. He, and Q. Gu. Provably efficient reinforcement learning for discounted MDPs with feature mapping. arXiv preprint arXiv:2006.13165, 2020.