# Provably Efficient Generalized Lagrangian Policy Optimization for Safe Multi-Agent Reinforcement Learning

**Dongsheng Ding**                                DONGSHED@SEAS.UPENN.EDU
*University of Pennsylvania, Philadelphia, PA 19104, USA*

**Xiaohan Wei**                                       UBIMETEOR@FB.COM
*Meta, Menlo Park, CA 94065 USA*

**Zhuoran Yang**                                ZHUORAN.YANG@YALE.EDU
*Yale University, New Haven, CT 06511, USA*

**Zhaoran Wang**                              ZHAORANWANG@GMAIL.COM
*Northwestern University, Evanston, IL 60208, USA*

**Mihailo R. Jovanović**                             MIHAILO@USC.EDU
*University of Southern California, Los Angeles, CA 90089, USA*

**Editors:** N. Matni, M. Morari, G. J. Pappas

## Abstract

We examine online safe multi-agent reinforcement learning using constrained Markov games in which agents compete by maximizing their expected total rewards under a constraint on expected total utilities. Our focus is confined to an episodic two-player zero-sum constrained Markov game with independent transition functions that are unknown to agents, adversarial reward functions, and stochastic utility functions. For such a Markov game, we employ an approach based on the occupancy measure to formulate it as an online constrained saddle-point problem with an explicit constraint. We extend the Lagrange multiplier method in constrained optimization to handle the constraint by creating a generalized Lagrangian with minimax decision primal variables and a dual variable. Next, we develop an upper confidence reinforcement learning algorithm to solve this Lagrangian problem while balancing exploration and exploitation. Our algorithm updates the minimax decision primal variables via online mirror descent and the dual variable via projected gradient step and we prove that it enjoys sublinear rate $O((|X| + |Y|)L\sqrt{T(|A| + |B|)})$ for both regret and constraint violation after playing $T$ episodes of the game. Here, $L$ is the horizon of each episode, $(|X|, |A|)$ and $(|Y|, |B|)$ are the state/action space sizes of the min-player and the max-player, respectively. To the best of our knowledge, we provide the first provably efficient online safe reinforcement learning algorithm in constrained Markov games.

**Keywords:** safe multi-agent reinforcement learning, constrained Markov game, upper confidence reinforcement learning, generalized Lagrange multiplier method, online mirror descent

## 1. Introduction

Safe Reinforcement Learning (RL) studies how a single agent learns to maximize its expected total reward subject to safety-concerned constraints by interacting with an unknown environment over time (García and Fernández, 2015; Thomas, 2015; Amodei et al., 2016). The constrained Markov decision processes (MDPs) provide a standard class of constraint critical environment models (Altman, 1999) that are utilized in autonomous robots (Feyzabadi, 2017; Fisac et al., 2018),

personalized medicine (Girard, 2018), online advertising (Boutilier and Lu, 2016), and financial management (Abe et al., 2010). General constrained MDPs for two or more agents are often formulated as constrained Markov games (MGs) in which agents compete under constraints (Altman and Shwartz, 2000; Altman et al., 2005, 2008), providing an effective model for safe multi-agent RL (Nguyen et al., 2014; Shalev-Shwartz et al., 2016; Zhang et al., 2021).

Considerable recent progress has been made in single-agent safe RL, especially for solving constrained MDP problems with constraint satisfaction guarantees (Efroni et al., 2020; Brantley et al., 2020; Bai et al., 2020a; Ding et al., 2021; Chen et al., 2021; Singh et al., 2022; Ding et al., 2022b). In these references, Lagrangian-based methods have been combined with the optimistic exploration to address exploration-exploitation trade-off under constraints. These constrained MDP learning algorithms are sample-efficient (in achieving both low regret and low constraint violation) and they effectively enhance classical RL methods to attain safety requirements. However, most of these algorithms are limited to the single-agent setting and it is an open question how to balance the exploration-exploitation trade-off under constraints for multiple agents. Another motivation for our work comes from recent advances on the efficient competitive RL algorithms in MGs (Wei et al., 2017; Bai and Jin, 2020; Bai et al., 2020b; Xie et al., 2020).

In this work, we take initial steps towards developing provably efficient safe multi-agent RL algorithms. We examine perhaps the most basic safe multi-agent RL setup that involves a two-player zero-sum constrained MG with independent state transitions (Altman and Shwartz, 2000; Altman et al., 2005, 2008; Singh and Hemachandra, 2014). This problem represents a generalization of constrained MDPs to the two-player case with coupled constraints. In such a constrained MG, two players follow their own state transitions independently, take actions simultaneously, and observe the reward and utility functions while competing against each other by maximizing/minimizing the reward while both are restrained by the constraint regarding some utility for safety reasons. The decision-coupling that arises from the constraint is often encountered in multi-agent systems (Rosen, 1965; Li and Marden, 2014; Kulkarni, 2011, 2017; De Nijs, 2019). More specifically, we aim to design an online RL algorithm for solving episodic two-player zero-sum constrained MGs. Here, two players do not know the transition models and have no access to a generative model, but can play the game for multiple episodes using arbitrary policies. The goal is to find an approximate constrained Nash equilibrium of the game in hindsight, a generalization of Nash equilibrium to characterize violating constraints if any unilateral deviations occur. We utilize a notion of regret to quantify the approximation error of the constrained Nash equilibrium and employ a constraint dissatisfaction (which results from violation of any utility constraints) to evaluate the constraint violation.

**Contribution**. We develop the first provably efficient algorithm for a constrained Markov game (MG) with $O(\sqrt{T})$ regret and $O(\sqrt{T})$ constraint violation. Specifically, we introduce an episodic constrained MG with unknown independent transition functions and decision-couplings that come from both adversarial reward functions and coupled stochastic constraints on utility functions. We use the occupancy measure approach to formulate such a MG as a constrained saddle-point problem with an explicit constraint. We extend the Lagrange method in constrained optimization to deal with the constraint by creating a generalized Lagrangian with minimax decision primal variables and a dual variable. We develop an upper confidence reinforcement learning algorithm – an Upper Confidence Bound Constrained SAddle-Point Optimization (UCB-CSAPO) algorithm – to solve this Lagrangian problem while balancing exploration and exploitation. Our algorithm updates the minimax decision primal variables via optimistic mirror descent and the dual variable via projected gradient step and we

prove that it enjoys sublinear rate $O((|X| + |Y|)L\sqrt{T(|A| + |B|)}))$ for both regret and constraint violation after playing $T$ episodes. Here, $L$ is the horizon of each episode, $(|X|, |A|)$ and $(|Y|, |B|)$ are the state/action space sizes of the min-player and max-player, respectively.

**Related Work.** We briefly review the most-related work; see Appendix 6 for details. Our work is closely related to safe multi-agent RL in constrained MGs. The Nash equilibrium for constrained MGs have been studied in Altman and Shwartz (2000); Gómez-Ramırez et al. (2003); Altman et al. (2005); Alvarez-Mena and Hernández-Lerma (2006); Altman et al. (2007, 2008); Altman and Solan (2009); Singh and Hemachandra (2014) using the notion of *constrained Nash equilibrium* (which generalizes the concept of *generalized Nash equilibrium* in static games (Arrow and Debreu, 1954) to MGs); see more studies in Yaji and Bhatnagar (2015); Zhang (2019); Wei (2020, 2021); Zhang and Zou (2021). These results are not applicable to the RL setting that assumes unknown models. Recently, asymptotic convergence in learning constrained MGs was examined in Hakami and Dehghan (2015); Jiang et al. (2020) but sample efficiency and exploration were not fully addressed, except for a concurrent work on learning correlated equilibria (Chen et al., 2022b). Our work fills this gap by adding built-in exploration mechanisms under constraints and proving the first non-asymptotic convergence for learning constrained Nash equilibria.

Our work is also pertinent to a rich RL literature on learning constrained MDPs (Zheng and Ratliff, 2020; Qiu et al., 2020; Kalagarla et al., 2020; Bai et al., 2020a; Chow et al., 2017; Tessler et al., 2019; Ding et al., 2020, 2021, 2022b; Wachi and Sui, 2020; Efroni et al., 2020; Brantley et al., 2020; Chen et al., 2021; Liu et al., 2021a; Ying et al., 2022; Liu et al., 2021b; Bai et al., 2022; Zhao and You, 2021; Li et al., 2021; Chen et al., 2022a). While these results provide provably efficient algorithms regarding regret and constraint satisfaction in the single-agent setting, they are not applicable to our multi-agent game being played under constraints, because of the *non-convexity* nauture of constrained multi-agent policy optimization and the *non-stationary* environment each agent is facing. An extended line of work on constrained MDPs focuses on cooperative multi-agent learning under constraints and most efforts study the case where multiple agents have independent MDPs with a coupled budget/resource constraint (Meuleau et al., 1998; Boutilier and Lu, 2016; Wei et al., 2018; de Nijs and Stuckey, 2020; Gagrani and Nayyar, 2020). All these results assume knowing transition models or system dynamics. Only a few studies considered the shared MDP case (Diddigi et al., 2019; Lu et al., 2020; Parnika et al., 2021; Gu et al., 2021), but they lack theoretical guarantees and do not handle exploration. In contrast, our work focuses on the MG setting with unknown models and attacks the exploration challenge directly.

## 2. Problem Setup

In this section, we introduce zero-sum Markov games (MGs) with constraints, which are categorized as constrained Markov/stochastic games (Altman and Shwartz, 2000; Altman et al., 2005, 2008).

In an episodic constrained MG there are two players; a *min-player* $- (X, A, P_1, r, g, T)$, which minimizes the reward, and a *max-player* $- (Y, B, P_2, r, h, T)$, which maximizes the reward, while adhering to a coupled utility constraint. Here, $T$ is the number of episodes, $X$ and $Y$ are finite state spaces, $A$ and $B$ are finite action spaces, $P_1$ and $P_2$ are transition probability measures where $P_1(\cdot \mid x, a)$ is a distribution over $X$ if the min-player takes action $a$ in state $x$ and $P_2(\cdot \mid y, b)$ is a distribution over $Y$ if the max-player takes action $b$ in state $y$, $r := \{r^t\}_{t=1}^T$ is a collection of players' reward functions $r^t \colon X \times Y \times A \times B \to [0, 1]$, whereas $g := \{g^t\}_{t=1}^T$ and $h := \{h^t\}_{t=1}^T$ are

collections of players' utility functions $g^t\colon X \times A \to [0,1]$, $h^t\colon Y \times B \to [0,1]$. For two independent transitions, players are coupled via the reward function and a constraint on their utility functions.

We utilize layered Markov decision processes to model the environment dynamics. For each player, e.g., the min-player, we assume that the state space $X$ has $L+1$ layers and that it satisfies the loop-free property: (i) $X := X_0 \cup \cdots \cup X_L$ and $X_{\ell_1} \cap X_{\ell_2} = \emptyset$ for $\ell_1 \neq \ell_2$; (ii) $X_0 = \{x_0\}$ and $X_L = \{x_L\}$; (iii) if $P_1(x' \mid x, a) > 0$, then $x' \in X_{\ell+1}$ and $x \in X_\ell$ for some $\ell \in \{0, 1, \cdots, L\}$. This assumption is common in loop-free stochastic shortest path problems (György et al., 2007; Jaksch et al., 2010; Neu et al., 2010; Rosenberg and Mansour, 2019; Jin et al., 2020); it is often used to simplify notation/analysis since any episodic MDPs can be reduced to be loop-free.

The min/max players interact with the environment in episode $t$ as follows. At the beginning, the environment determines the reward function $r^t$ and the utility functions $g^t$ and $h^t$. Meanwhile, two players decide their policies $\pi^t\colon X \times A \to [0,1]$ and $\mu^t\colon Y \times B \to [0,1]$, where $\pi^t(\cdot \mid x)$ and $\mu^t(\cdot \mid y)$ are probability distributions over their action spaces $A$ and $B$, respectively. Then, given initial states $x_0$ and $y_0$, both players execute their own policies $\pi^t$ or $\mu^t$ for $L$ steps. At step $\ell \in \{0, \ldots, L-1\}$, each player only observes its own state $x_\ell \in X$ or $y_\ell \in Y$, takes action $a_\ell$ or $b_\ell$ following its own policy $\pi^t$ or $\mu^t$, transits to next state $x_{\ell+1}$ or $y_{\ell+1}$ according to its own transition $P_1(\cdot \mid x_\ell, a_\ell)$ or $P_2(\cdot \mid y_\ell, b_\ell)$, and observes reward $r^t$ and local utility $g^t$ or $h^t$. Assume there is no dependence between functions $r^t$, $g^t$, and $h^t$ and they are independent of the underlying MDPs.

To define the learning objective, for the min-player in episode $t$ we introduce the occupancy measure $q_1^t\colon X \times A \times X \to [0,1]$ by $q_1^t(x, a, x') := \mathrm{Prob}(x_\ell = x, a_\ell = a, x_{\ell+1} = x')$ for $x \in X_\ell$, describing the marginal probability of visiting $(x, a, x')$ when executing policy $\pi^t$ under the transition $P_1$. Similarly, we introduce the occupancy measure $q_2^t\colon Y \times B \times Y \to [0,1]$ for the max-player. We recall that a function $q\colon X \times A \times X \to [0,1]$ is an occupancy measure associated with policy $\pi$ and transition $P$ if and only if it satisfies two conditions (Altman, 1999): (i) $\sum_{x \in X_\ell} \sum_{a \in A} \sum_{x' \in X_{\ell+1}} q(x, a, x') = 1$ for $\ell \in \{0, \ldots, L-1\}$; (ii) $\sum_{x \in X_{\ell-1}} \sum_{a \in A} q(x, a, x') = \sum_{a \in A} \sum_{x'' \in X_{\ell+1}} q(x', a, x'')$ for $x' \in X_\ell$ and $\ell \in \{1, \ldots, L-1\}$. We denote by $\Delta(P)$ a set of valid occupancy measures under $P$,

$$\Delta(P) := \big\{ q\colon X \times A \times X \to [0,1] \mid q \text{ satisfies (i) and (ii) as shown above} \big\}.$$

It is worth noting that the occupancy measure set is convex and compact for finite MDPs (Altman, 1999). Using an occupancy measure $q$, we can express associated transition $P$ and policy $\pi$ as

$$P(x' \mid x, a) = \frac{q(x, a, x')}{\sum_{x'' \in X_{\ell+1}} q(x, a, x'')} \text{ and } \pi(a \mid x) = \frac{\sum_{x' \in X_{\ell+1}} q(x, a, x')}{\sum_{a \in A} \sum_{x'' \in X_{\ell+1}} q(x, a, x'')} \quad (1)$$

where $x \in X_\ell$. Slightly extending the notation $q$, we use it to represent the probability of visiting $(x, a)$, i.e., $q(x, a) = \sum_{x' \in X_{\ell+1}} q(x, a, x')$ for $x \neq x_L$. These properties imply that the problem of learning a policy equals learning the associated occupancy measure (Zimin and Neu, 2013).

In episode $t$, given a min-policy $\pi^t$ and a max-policy $\mu^t$, we introduce the expected total reward,

$$\mathbb{E}_{P_1, P_2, \pi^t, \mu^t} \left[ \sum_{\ell=0}^{L-1} r^t(x_\ell, y_\ell, a_\ell, b_\ell) \right] = \sum_{\ell=0}^{L-1} \sum_{x \in X_\ell, y \in Y_\ell} \sum_{a \in A, b \in B} q_1^t(x, a) q_2^t(y, b) r^t(x, y, a, b)$$
$$:= \big\langle q_1^t \cdot q_2^t, r^t \big\rangle$$

$$(2)$$

where the expectation $\mathbb{E}$ is taken over the random state-action sequence $\{(x_\ell, y_\ell, a_\ell, b_\ell)\}_{\ell=0}^{L-1}$; the action $a_\ell$ follows the policy $\pi^t(\cdot \,|\, x_\ell)$ in the state $x_\ell$ and the next state $x_{\ell+1}$ follows the transition $P_1(\cdot \,|\, x_\ell, a_\ell)$; the action $b_\ell$ follows the policy $\mu^t(\cdot \,|\, y_\ell)$ in the state $y_\ell$ and the next state $y_{\ell+1}$ follows the transition $P_2(\cdot \,|\, y_\ell, b_\ell)$. Similarly, we can define the expected total utilities as

$$\mathbb{E}_{P_1, \pi^t} \left[ \sum_{\ell=0}^{L-1} g_x^t(x_\ell, a_\ell) \right] = \sum_{\ell=0}^{L-1} \sum_{x \in X_\ell} \sum_{a \in A} q_1^t(x, a) g^t(x, a) := \left\langle q_1^t, g^t \right\rangle \tag{3a}$$

$$\mathbb{E}_{P_2, \mu^t} \left[ \sum_{\ell=0}^{L-1} h^t(y_\ell, a_\ell) \right] = \sum_{\ell=0}^{L-1} \sum_{y \in Y_\ell} \sum_{b \in B} q_2^t(y, b) h^t(y, b) := \left\langle q_2^t, h^t \right\rangle. \tag{3b}$$

In general, reward function $r^t$ and utility functions $g^t$ and $h^t$ all can change arbitrarily, i.e., being adversarial. However, even if we fix the opponent's policy, there is no algorithm for the player to achieve sublinear regret and constraint violation at the same time when the constraints are changing adversarially (Mannor et al., 2009). Hence, we restrict the utility functions to be stochastic: $g^t(x, a) := g(x, a; \xi^t)$, $h^t(y, b) := h(y, b; \xi^t)$ with $\mathbb{E}\left[g^t(x, a)\right] = g(x, a)$ and $\mathbb{E}\left[h^t(y, b)\right] = h(y, b)$, for any $x \in X$, $a \in A$ and $y \in Y$, $b \in B$, where $\xi^t$ is an independent random variable.

**Learning Performance**. We now define the underlying constrained optimization problem and the solution concept for learning constrained MGs. Using the notion of occupancy measure, we formulate a constrained minimax problem in which the objective function is a sum of the expected total rewards over $T$ episodes and the constraint is on a sum of two agent' expected total utilities,

$$\underset{q_1 \in \Delta(P_1)}{\text{minimize}} \ \underset{q_2 \in \Delta(P_2)}{\text{maximize}} \ \sum_{t=0}^{T-1} \left\langle q_1 \cdot q_2, r^t \right\rangle \quad \text{subject to} \quad \left\langle q_1, g \right\rangle + \left\langle q_2, h \right\rangle \leq b \tag{4}$$

where we take $b \in (0, 2L]$ to avoid trivial cases since we note that $\langle q_1, g \rangle$, $\langle q_2, h \rangle \in [0, L]$. The coupled constraint is used to model the limited use of budget/resource for two players; multi-agent problems with a common constraint are often called *weakly-coupled* or *non-orthogonal* in the literature on CMDPs (Meuleau et al., 1998; Boutilier and Lu, 2016; Wei et al., 2018; Salemi Parizi, 2018; Gagrani and Nayyar, 2020) and constrained MGs (Altman et al., 2008; Altman and Solan, 2009; Kulkarni, 2011; Singh and Hemachandra, 2014; Kulkarni, 2017). We can generalize it to multiple or local side constraints, e.g., $\langle q_1, g \rangle \leq b_1$ or $\langle q_2, h \rangle \leq b_2$. When transitions $P_1$ and $P_2$ are known, the occupancy measure sets $\Delta(P_1)$ and $\Delta(P_2)$ define convex polytopes on $q_1$ and $q_2$.

Let $(q_1^\star, q_2^\star)$ be a solution to Problem (4) in hindsight. The existence of $(q_1^\star, q_2^\star)$ follows from compactness of the constraint sets (Neumann, 1928; Rosen, 1965). It is standard to define an intuitive solution – constrained Nash equilibrium – via two conditions (Altman and Shwartz, 2000; Daskalakis et al., 2021):

(i) $\displaystyle\sum_{t=0}^{T-1} \langle q_1^\star \cdot q_2^\star, r^t \rangle \ \leq \ \sum_{t=0}^{T-1} \langle q_1 \cdot q_2^\star, r^t \rangle$ for any $q_1 \in \Delta(P_1)$ satisfying $\langle q_1, g \rangle + \langle q_2^\star, h \rangle \leq b$;

(ii) $\displaystyle\sum_{t=0}^{T-1} \langle q_1^\star \cdot q_2, r^t \rangle \ \leq \ \sum_{t=0}^{T-1} \langle q_1^\star \cdot q_2^\star, r^t \rangle$ for any $q_2 \in \Delta(P_2)$ satisfying $\langle q_1^\star, g \rangle + \langle q_2, h \rangle \leq b$.

Any unilateral deviation from the constrained Nash equilibrium will either break the constraint, or if it is not, then there is no benefit for this player. With this solution concept, we define the regret for any algorithm that plays the game for $T$ episodes by

$$\text{Regret}(T) \; = \; \sum_{t=0}^{T-1} \left( \langle q_1^t \cdot q_2^\star, r^t \rangle - \langle q_1^\star \cdot q_2^t, r^t \rangle \right) \tag{5}$$

which adds two side optimality gaps, $\sum_{t=0}^{T-1} \langle q_1^t \cdot q_2^\star, r^t \rangle - \langle q_1^\star \cdot q_2^\star, r^t \rangle$ for the min-player and $\sum_{t=0}^{T-1} \langle q_1^\star \cdot q_2^\star, r^t \rangle - \langle q_1^\star \cdot q_2^t, r^t \rangle$ for the max-player, and two players take policies $\pi^t$ and $\mu^t$ in episode $t$ and they define occupancy measures $q_1^t$ and $q_2^t$ under the true transitions $P_1$ and $P_2$. This regret works in a notion of weak regret (Brafman and Tennenholtz, 2002; Bai and Jin, 2020; Xie et al., 2020) instead of the single-agent type regret (Tian et al., 2020; Bai et al., 2020b) which is statistically and computationally hard to bound sublinearly.

To measure the constraint satisfaction, we introduce the violation as a non-negative part of accumulated constraint violations $\langle q_1^t, g \rangle + \langle q_2^t, h \rangle - b$ over $T$ episodes,

$$\text{Violation}(T) \; = \; \left[ \sum_{t=0}^{T-1} \left( \langle q_1^t, g^t \rangle + \langle q_2^t, h^t \rangle - b \right) \right]_+. \tag{6}$$

We next assume feasibility that ensures the existence of constrained Nash equilibrium (Altman and Shwartz, 2000). Feasibility can be verified by a priori knowledge on feasible policies.

**Assumption 1 (Feasibility)** *There exists a joint policy $(\bar{\pi}, \bar{\mu})$ associated to the occupancy measure $(\bar{q}_1, \bar{q}_2)$ and a constant $\xi > 0$ such that $\langle \bar{q}_1, g \rangle + \langle \bar{q}_2, h \rangle + \xi \leq b$.*

Having defined the learning performance, we will work with the occupancy measure in the online learning setting where the two players do not know the transition functions, only observe reward/utility functions at the end of each episode, repeatedly play the game for a fixed number of episodes to learn the constrained Nash equilibrium in hindsight.

## 3. Proposed Algorithm

We present a variant of upper confidence reinforcement learning in Algorithm 1 – an Upper Confidence Bound Constrained SAddle-Point Optimization (UCB-CSAPO) algorithm – for learning constrained MGs. Conceptually, the algorithm works as the primal-dual policy optimization (Efroni et al., 2020; Ding et al., 2021; Chen et al., 2021) in the Lagrangian-based framework, which makes it a simple policy optimization algorithm. However, our primal update exploits the structure of constrained MGs to maintain two players' occupancy measures. The domain set of occupancy measures builds on the upper confidence bound exploration or optimism (Jaksch et al., 2010) regarding the estimated transition models using past trajectories. The dual update determines the penalty weight by collecting the possible constraint violation already acquired. In each episode, our algorithm has two key stages: (i) The generalized Lagrangian mirror descent step for updating the occupancy measures with optimism; (ii) The estimation of confidence sets on the occupancy measures.

**Generalized Lagrangian Mirror Descent Step**. The main idea of this step is to apply the online primal-dual mirror descent – an algorithmic generalization of online mirror descent to the constrained problems (Wei et al., 2020) – to the constrained MG setting (Altman and Shwartz, 2000; Altman

et al., 2005, 2008; Singh and Hemachandra, 2014). Let us recall that the occupancy measures $q_1^t$ for the min-player and $q_2^t$ for the max-player are defined over the true transitions $P_1$ and $P_2$ in episode $t$. The primal update of our algorithm maintains two occupancy measures $\widehat{q}_1^t, \widehat{q}_2^t$ to estimate $q_1^t, q_2^t$, separately. Although $\widehat{q}_1^t, \widehat{q}_2^t$ do not necessarily come from the true transitions $P_1, P_2$, they propose a min-policy $\pi^t$ for the min-player and a max-policy $\mu^t$ for the max-player according to the occupancy measure's property (1), i.e., for all $(x, a) \in X \times A$ and $(y, b) \in Y \times B$,

$$\pi^t(a \,|\, x) \;=\; \frac{\sum\limits_{x'} \widehat{q}_1^t(x, a, x')}{\sum\limits_{a, x''} \widehat{q}_1^t(x, a, x'')} \;\text{ and }\; \mu^t(b \,|\, y) \;=\; \frac{\sum\limits_{y'} \widehat{q}_2^t(y, b, y')}{\sum\limits_{b, y''} \widehat{q}_2^t(y, b, y'')}. \tag{7}$$

We describe our Lagrangian-based design to update estimates $\widehat{q}_1^t$ and $\widehat{q}_2^t$ in an online fashion. Assume that the transitions $P_1$ and $P_2$ are known. We consider a one-episode constrained minimax problem based on reward/utility functions: $r^{t-1}, g^{t-1}, h^{t-1}$, revealed at the end of episode $t-1$,

$$\underset{q_1 \in \Delta(P_1)}{\text{minimize}} \; \underset{q_2 \in \Delta(P_2)}{\text{maximize}} \quad \langle q_1 \cdot q_2, r^{t-1} \rangle \quad \text{subject to} \quad \langle q_1, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle \;\le\; b$$

where $\Delta(P_1)$ and $\Delta(P_2)$ are sets of valid occupancy measures under $P_1$ and $P_2$, respectively.

It is standard to use the method of Lagrange multipliers (Bertsekas, 2014) to handle constraints by adding penalty terms, if any constraint violation appears, into the original objective, and formulate an unconstrained problem. This is found in constrained games with separate side constraints (Pearsall, 1976) and multiple MDPs with coupled constraints (Boutilier and Lu, 2016; Wei et al., 2018). However, for constrained MGs either player can contribute to constraint violation $\langle q_1, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle - b$. It is important to specify which player should get such penalty terms (Altman and Solan, 2009; Dai and Zhang, 2020). We employ an attitude that the two players are jointly against the constraint while competing for rewards (Altman and Solan, 2009). As a result, both would sacrifice their rewards to satisfy the constraint if any violation occurs. We approximate the violation for each player as: $\langle q_1, g^{t-1} \rangle + \langle \widehat{q}_2^t, h^{t-1} \rangle - b$ for the min-player, and $\langle \widehat{q}_1^t, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle - b$ for the max-player. We formulate a generalized Lagrangian-type function,

$$\begin{aligned} L^t(q_1, q_2; \lambda) \;&:=\; \langle q_1 \cdot q_2, r^{t-1} \rangle \\ &\quad + \lambda\big(\langle q_1, g^{t-1} \rangle + \langle \widehat{q}_2^t, h^{t-1} \rangle - b\big) \;-\; \lambda\big(\langle \widehat{q}_1^t, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle - b\big) \end{aligned}$$

where $q_1$ is the first primal variable for the min-player, $q_2$ is the second primal variable for the max-player, and $\lambda \ge 0$ works as the Lagrange multiplier or the dual variable in penalizing the min-player/max-player via the first/second $\lambda$-term. Once we update $\lambda = \lambda^{t-1}$ from the last episode, we reach a constrained saddle-point problem, $\text{minimize}_{q_1 \in \Delta(P_1)} \text{maximize}_{q_2 \in \Delta(P_2)} L^t(q_1, q_2; \lambda^{t-1})$.

However, it is not feasible to take the domains $\Delta(P_1)$ and $\Delta(P_2)$ since the true transitions $P_1$ and $P_2$ are unknown. Instead, by the optimism in the face of uncertainty, we use their optimistic estimates $\Delta(k_1^t)$ and $\Delta(k_2^t)$ in sense that $q_1^t \in \Delta(k_1^t)$ and $q_2^t \in \Delta(k_2^t)$ hold with high probability in Lemma 1, where $\Delta(k_1^t)$ and $\Delta(k_2^t)$ are given by (11). Let $\widehat{q}^t := (\widehat{q}_1^t, \widehat{q}_2^t)$ and $D(p \,|\, q) := \sum_i p_i \ln \frac{p_i}{q_i} - \sum_i (p_i - q_i)$ that is the unnormalized Kullback-Leibler (KL) divergence between two distributions $p, q$. By a linear approximation of $L^t(q_1, q_2; \lambda^{t-1})$ at the previous iterate $(q_1^{t-1}, q_2^{t-1})$, we update the primal variable via an online mirror descent step over the domains of $q_1$ and $q_2$,

$$\begin{aligned} \widehat{q}^t \;\leftarrow\; \underset{q_1 \in \Delta(k_1^t)}{\text{argmin}} \; \underset{q_2 \in \Delta(k_2^t)}{\text{argmax}} \; \Big( &V \big\langle q_1 \cdot \widehat{q}_2^{t-1} + \widehat{q}_1^{t-1} \cdot q_2, r^{t-1} \big\rangle \\ &+ \lambda^{t-1}\big(\langle q_1, g^{t-1} \rangle - \langle q_2, h^{t-1} \rangle\big) \;+\; \eta^{-1} D\big(q \,|\, \widetilde{q}^{t-1}\big)\Big) \end{aligned} \tag{8}$$

where $V > 0$ provides the tradeoff between the minimax objective and the constraint, $\eta > 0$ is the learning rate, $D(\cdot \,|\, \cdot)$ is the unnormalized Kullback-Leibler divergence with a slightly abuse in a way that $D(q \,|\, q') := D(q_1 \,|\, q_1') - D(q_2 \,|\, q_2')$, $\widetilde{q}_1^{t-1}$ and $\widetilde{q}_2^{t-1}$ are mixing policies, e.g.,

$$\widetilde{q}_1^{t-1}(x,a) \;=\; (1 - \theta)\,\widehat{q}_1^{t-1}(x,a) \;+\; \theta\,\frac{1}{|X_\ell||A|} \tag{9}$$

for $(x,a) \in X_\ell \times A$, $\ell \in \{0, 1, \ldots, L-1\}$, $\theta \in (0,1]$. The mixing step ensures the uniform boundedness of KL divergence and also adds extra exploration into policy search (Wei et al., 2020). Moreover, we offer an efficient implementation of (8) as solving a convex program in Appendix 8.

Once we obtain $\widehat{q}^t$, we next perform the dual update. If we treat two $\lambda$-related regularization terms in $L^t(\widehat{q}_1^t, \widehat{q}_2^t; \lambda)$ separately, then gradient ascent/descent over either $\lambda$ leads to the same update rule using the constraint violation $\langle \widehat{q}_1^t, g^{t-1} \rangle + \langle \widehat{q}_2^t, h^{t-1} \rangle - b$. Hence, the dual update works in the usual way by adding up all past constraint violations,

$$\lambda^t \;=\; \max\left(\lambda^{t-1} \;+\; (\langle \widehat{q}_1^t, g^{t-1} \rangle + \langle \widehat{q}_2^t, h^{t-1} \rangle - b\,),\, 0\right). \tag{10}$$

The dual update (10) increases $\lambda^{t-1}$ when $\widehat{q}^t$ violates the approximate constraint $\langle q_1, g^{t-1} \rangle + \langle q_2, h^{t-1} \rangle \le b$. It penalizes both players by yielding individual gains to the constraint satisfaction. The dual update finds uses in constrained MDP problems (Efroni et al., 2020; Ding et al., 2021).

**Estimation of Confidence Sets**. To deal with unknown transitions $P_1$ and $P_2$, we employ the upper confidence bound (Jaksch et al., 2010; Neu et al., 2010) to estimate occupancy measure sets $\Delta(P_1)$, $\Delta(P_2)$. We exploit players' history trajectories to estimate their true transitions: $P_1$, $P_2$, and describe estimation uncertainty as confidence sets. The estimation proceeds in epochs as follows.

Let the epoch index for the min-player be $k_1 \in \{1, 2, \ldots\}$ and the epoch index for the max-player be $k_2 \in \{1, 2, \ldots\}$. We may represent them by $k_1^t$ and $k_2^t$ for showing the dependence on episode $t$. The epoch counters work in the following way. For each player, e.g., the min-player, we denote by $N_1^{k_1}(x,a)$ and $M_1^{k_1}(x,a,x')$ the total numbers of visitations to $(x,a)$ and $(x,a,x')$ before epoch $k_1$, respectively; we represent the total numbers of visitations to $(x,a)$ and $(x,a,x')$ in epoch $k_1$ by $n_1^{k_1}(x,a)$ and $m_1^{k_1}(x,a,x')$, respectively; If there exists $(x,a)$ such that $n_1^{k_1}(x,a) \ge N_1^{k_1}(x,a)$, then we set a new epoch by increasing $k_1$ by one. Similarly, we define $N_2^{k_2}(y,b)$, $M_2^{k_2}(y,b,y')$, $n_2^{k_2}(y,b)$, and $m_2^{k_2}(y,b,y')$ for the max-player. Using the defined epoch and visitation counters, we empirically estimate the true transitions $P_1$ or $P_2$ in epoch $k_1$ or $k_2$ by

$$\bar{P}_1^{k_1}(x' \,|\, x,a) \;=\; \frac{M_1^{k_1}(x,a,x')}{\max(1, N_1^{k_1}(x,a))} \quad \text{and} \quad \bar{P}_2^{k_2}(y' \,|\, y,b) \;=\; \frac{M_2^{k_2}(y,b,y')}{\max(1, N_2^{k_2}(y,b))}$$

for all $(x,a,x') \in X \times A \times X$ and $(y,b,y') \in Y \times B \times Y$.

Let the confidence set of epoch $k_1$ for the min-player be $\mathcal{P}_1^{k_1}$ and the confidence set of epoch $k_2$ for the max-player be $\mathcal{P}_2^{k_2}$. We take $\mathcal{P}_1^{k_1}$ and $\mathcal{P}_2^{k_2}$ as collections of transitions that deviate from the empirical ones at most $\epsilon_1^{k_1}$ and $\epsilon_2^{k_2}$,

$$\mathcal{P}_1^{k_1} \;=\; \left\{\widehat{P}_1 \,\middle|\, \|\widehat{P}_1(\cdot \,|\, x,a) - \bar{P}_1^{k_1}(\cdot \,|\, x,a)\|_1 \le \epsilon_1^{k_1}, \forall(x,a)\right\}$$

$$\mathcal{P}_2^{k_2} \;=\; \left\{\widehat{P}_2 \,\middle|\, \|\widehat{P}_2(\cdot \,|\, y,b) - \bar{P}_2^{k_2}(\cdot \,|\, y,b)\|_1 \le \epsilon_2^{k_2}, \forall(y,b)\right\}$$

where we take $\epsilon_1^{k_1}(x,a) = \sqrt{\frac{2|X_{\ell(x)+1}|\log(T|A||X|/\delta)}{\max(1, N_1^{k_1}(x,a))}}$ and $\epsilon_2^{k_2}(y,b) = \sqrt{\frac{2|Y_{\ell(y)+1}|\log(T|B||Y|/\delta)}{\max(1, N_2^{k_2}(y,b))}}$, $\ell(x)$ and $\ell(y)$ are the layers that certain states belong to, and $\delta \in (0,1)$. We recall the occupancy measure

---

**Algorithm 1** <u>U</u>pper <u>C</u>onfidence <u>B</u>ound <u>C</u>onstrained <u>SA</u>ddle-<u>P</u>oint <u>O</u>ptimization (UCB-CSAPO)

---

1: **Input**: State/action spaces $(X, A)$ and $(Y, B)$, episode $T$, parameters $V, \eta, \theta$, and $p \in (0, 1)$.

2: **Initialization**: The min-player: $\widehat{q}_1^0(x, a, x') = \frac{1}{|X^\ell||A||X^{\ell+1}|}, \forall (x, a, x') \in X^\ell \times A \times X^{\ell+1}, \ell \in [0, L-1]; n_1^1(x, a) = N_1^1(x, a) = 0, \forall (x, a); m_1^1(x, a, x') = M_1^1(x, a, x') = \bar{P}_1^1(x' \mid x, a) = 0, \forall (x, a, x')$.

The max-player: $\widehat{q}_2^0(y, b, y') = \frac{1}{|Y^\ell||B||Y^{\ell+1}|}, \forall (y, b, y') \in Y^\ell \times B \times Y^{\ell+1}, \ell \in [0, L-1]; n_2^1(y, b) = N_2^1(y, b) = 0, \forall (y, b); m_2^1(y, b, y') = M_2^1(y, b, y') = \bar{P}_2^1(y' \mid y, b) = 0, \forall (y, b, y')$.
Let $r^0, g^0, h^0$ be zero functions, $\lambda^0$ be zero, and $k_1^1 = k_2^1 = 1$.

3: **for** episode $t = 1, \ldots, T$ **do**

4:     Update the primal variable $\widehat{q}^t$ via (8) and the dual variable $\lambda^t$ via (10).

5:     Compute the min-policy $\pi^t$ and the max-policy $\mu^t$ via (7). Execute them for $L$ steps and record trajectories $(x^0, a^0, x^1, \cdots, a^{L-1}, x^{L-1})$ and $(y^0, b^0, y^1, \cdots, b^{L-1}, y^{L-1})$, and reward/utility functions $r^t, g^t$, and $h^t$.

6:     Update local visitation counters at visited trajectories,

$$n_1^{k_1^t}(x^\ell, a^\ell) \leftarrow n_1^{k_1^t}(x^\ell, a^\ell) + 1 \text{ and } m_1^{k_1^t}(x^\ell, a^\ell, x^{\ell+1}) \leftarrow m_1^{k_1^t}(x^\ell, a^\ell, x^{\ell+1}) + 1$$

$$n_2^{k_2^t}(y^\ell, b^\ell) \leftarrow n_2^{k_2^t}(y^\ell, b^\ell) + 1 \text{ and } m_2^{k_2^t}(y^\ell, b^\ell, y^{\ell+1}) \leftarrow m_2^{k_2^t}(y^\ell, b^\ell, y^{\ell+1}) + 1.$$

7:     **if** $n_1^{k_1^t}(x, a) \geq N_1^{k_1^t}(x, a)$ or $n_2^{k_2^t}(y, b) \geq N_2^{k_2^t}(y, b)$ for some $(x, a) \in X \times A$ or $(y, b) \in Y \times B$ **then**

8:         Increase epoch counter by one, $k_1^{t+1} \leftarrow k_1^t + 1$ or $k_2^{t+1} \leftarrow k_2^t + 1$, and update global visitation counters,

$$N_1^{k_1^{t+1}}(x, a) \leftarrow N_1^{k_1^t}(x, a) + n_1^{k_1^t}(x, a) \text{ or } N_2^{k_2^{t+1}}(y, b) \leftarrow N_2^{k_2^t}(y, b) + n_2^{k_2^t}(y, b)$$

$$M_1^{k_1^{t+1}}(x, a, x') \leftarrow M_1^{k_1^t}(x, a, x') + m_1^{k_1^t}(x, a, x') \text{ or } M_2^{k_2^{t+1}}(y, b, y') \leftarrow M_2^{k_2^t}(y, b, y') + m_2^{k_2^t}(y, b, y').$$

        Update the confidence bounds for $\Delta(k_1^t)$ or $\Delta(k_2^t)$ in (11), and set $n_1^{k_1^{t+1}}(x, a) = m_1^{k_1^{t+1}}(x, a, x') = 0$ for all $(x, a)$ and $(x, a, x')$ or $n_2^{k_2^{t+1}}(y, b) = m_2^{k_2^{t+1}}(y, b, y') = 0$ for all $(y, b)$ and $(y, b, y')$.

9:     **else**

10:         Set either $k_1^{t+1} = k_1^t$ or $k_2^{t+1} = k_2^t$.

11:     **end if**

12: **end for**

---

sets $\Delta(P_1)$ or $\Delta(P_2)$ that are induced by the true transitions $P_1$ or $P_2$. We generalize this notion to define $\Delta(\mathcal{P}_1^{k_1^t})$ or $\Delta(\mathcal{P}_2^{k_2^t})$ as collections of all possible occupancy measures that are induced by the estimated transitions $\widehat{P}_1 \in \mathcal{P}_1^k$ or $\widehat{P}_2 \in \mathcal{P}_2^k$,

$$\Delta(k_1^t) := \Delta(\mathcal{P}_1^{k_1^t}) \text{ or } \Delta(k_2^t) := \Delta(\mathcal{P}_2^{k_2^t}); \text{ see (17) in Appendix 8 for explicit forms.} \quad (11)$$

**Lemma 1** *Fix $\delta \in (0, 1)$. With probability $1 - \delta$, $\Delta(P_1) \subset \Delta(\mathcal{P}_1^{k_1})$ and $\Delta(P_2) \subset \Delta(\mathcal{P}_2^{k_2})$ for all $k_1, k_2 \in \{1, 2, \ldots\}$.*

The proof of Lemma 1 follows the confidence bound construction; we provide it in Appendix 9. For all epoch $k_1^t$ or $k_2^t$ (episode $t$), the true transitions $P_1$ and $P_2$ are contained in $\mathcal{P}_1^{k_1^t}$ and $\mathcal{P}_2^{k_2^t}$, respectively, with high probability. This supports the primal update (8) such that both players are optimistically searching solutions in a large but tractable domain.

## 4. Performance Guarantees

In Theorem 2, we present our main theoretical result on the regret and the constraint violation for Algorithm 1. We recall the total number of games played by the algorithm $T$, the size of state/action spaces of the min-player $|X|, |A|$, and the size of state/action spaces of the max-player $|Y|, |B|$.

**Theorem 2 (Regret Bound and Constraint Violation)**  *Let Assumption 1 hold. Fix $p \in (0, 1)$ and $T \geq \max(|X||A|, |B||Y|)$. In Algorithm 1, we set $V = L\sqrt{T}$, $\eta = 1/(TL)$, and $\theta = 1/T$. Then, with probability $1 - p$, the regret (5) and the constraint violation (6) satisfy*

$$\text{Regret}(T), \ \text{Violation}(T) \ \leq \ \widetilde{O}\big( (|X| + |Y|) L\sqrt{T(|A| + |B|)} \big)$$

*where $\widetilde{O}(\cdot)$ hides the logarithmic factor $\log \frac{1}{p}$.*

In Theorem 2, we prove that UCB-CSAPO enjoys $O(\sqrt{T})$ regret and $O(\sqrt{T})$ constraint violation using appropriate algorithm parameters $\{V, \eta, \theta, p\}$ and Assumption 1; see Appendix 7 for proof. Our bounds have the optimal dependence on the total number of episodes $T$ up to some logarithmic factors. The $\sqrt{|A| + |B|}$ dependence matches the existing lower bound for the single-player case (Bai and Jin, 2020). The only suboptimal dependence comes from $|X|, |Y|$ that also exists in existing unconstrained loop-free stochastic shortest path problems (Rosenberg and Mansour, 2019). It is straightforward to remove knowledge of $T$ by using the doubling trick while not altering our bounds up to logarithmic factors (Rakhlin and Sridharan, 2013).

We see that Assumption 1 does not impose any restrictions on rewards. Hence, UCB-CSAPO is robust against adversarial reward functions. Moreover, Theorem 2 carries to other settings, e.g., constrained MGs with side constraints; see Appendix 14.

## 5. Concluding Remarks

We have examined an episodic two-player zero-sum constrained Markov game (MG) with independent transition functions. In our setup, transition functions are unknown to agents, reward functions are adversarial, and utility functions are stochastic. We have proposed the first provably efficient algorithm for playing constrained MGs with $O(\sqrt{T})$ regret and constraint violation. Our algorithm provides a principled extension of the upper confidence reinforcement learning to deal with coupled constraints in constrained MGs. We also remark that the developed algorithmic framework can be readily applied to learning other constrained MGs, e.g., the ones that involve a single controller.

Our work opens up many interesting directions for future work, such as sharper algorithms with sample complexity lower bounds, constrained rational algorithms, and how to perform safe exploration in other models of constrained MGs.

## Acknowledgments

## References

Naoki Abe, Prem Melville, Cezar Pendus, Chandan K Reddy, David L Jensen, Vince P Thomas, James J Bennett, Gary F Anderson, Brent R Cooley, Melissa Kowalczyk, et al. Optimizing debt collections using constrained reinforcement learning. In *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 75–84, 2010.

Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. Constrained policy optimization. In *International Conference on Machine Learning*, volume 70, pages 22–31, 2017.

Eitan Altman. *Constrained Markov decision processes*, volume 7. CRC Press, 1999.

Eitan Altman and Adam Shwartz. Constrained markov games: Nash equilibria. In *Advances in Dynamic Games and Applications*, pages 213–221. Birkhäuser Boston, 2000.

Eitan Altman and Eilon Solan. Constrained games: The impact of the attitude to adversary's constraints. *IEEE Transactions on Automatic Control*, 54(10):2435–2440, 2009.

Eitan Altman, Konstantin Avrachenkov, Richard Marquez, and Gregory Miller. Zero-sum constrained stochastic games with independent state processes. *Mathematical Methods of Operations Research*, 62(3):375–386, 2005.

Eitan Altman, Saswati Sarkar, and Eilon Solan. Constrained Markov games with transition probabilities controlled by a single player. In *International Conference on Performance Evaluation Methodologies and Tools*, pages 1–6, 2007.

Eitan Altman, Konstantin Avrachenkov, Nicolas Bonneau, Merouane Debbah, Rachid El-Azouzi, and Daniel Sadoc Menasche. Constrained cost-coupled stochastic games with independent state processes. *Operations Research Letters*, 36(2):160–164, 2008.

Jorge Alvarez-Mena and Onésimo Hernández-Lerma. Existence of Nash equilibria for constrained stochastic games. *Mathematical Methods of Operations Research*, 63(2):261–285, 2006.

Dario Amodei, Chris Olah, Jacob Steinhardt, Paul Christiano, John Schulman, and Dan Mané. Concrete problems in AI safety. *arXiv preprint arXiv:1606.06565*, 2016.

Kenneth J Arrow and Gerard Debreu. Existence of an equilibrium for a competitive economy. *Econometrica: Journal of the Econometric Society*, pages 265–290, 1954.

Qinbo Bai, Vaneet Aggarwal, and Ather Gattami. Model-free algorithm and regret analysis for MDPs with long-term constraints. *arXiv preprint arXiv:2006.05961*, 2020a.

Qinbo Bai, Amrit Singh Bedi, Mridul Agarwal, Alec Koppel, and Vaneet Aggarwal. Achieving zero constraint violation for constrained reinforcement learning via primal-dual approach. In *AAAI Conference on Artificial Intelligence*, volume 36, pages 3682–3689, 2022.

Yu Bai and Chi Jin. Provable self-play algorithms for competitive reinforcement learning. In *International Conference on Machine Learning*, pages 551–560, 2020.

Yu Bai, Chi Jin, and Tiancheng Yu. Near-optimal reinforcement learning with self-play. *Advances in Neural Information Processing Systems*, 33, 2020b.

Dimitri P Bertsekas. *Constrained optimization and Lagrange multiplier methods*. Academic Press, 2014.

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

Craig Boutilier and Tyler Lu. Budget allocation using weakly coupled, constrained Markov decision processes. In *Conference on Uncertainty in Artificial Intelligence*, pages 52–61, 2016.

Stephen Boyd, Stephen P Boyd, and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.

Ronen I Brafman and Moshe Tennenholtz. R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3:213–231, 2002.

Kianté Brantley, Miroslav Dudik, Thodoris Lykouris, Sobhan Miryoosefi, Max Simchowitz, Aleksandrs Slivkins, and Wen Sun. Constrained episodic reinforcement learning in concave-convex and knapsack settings. *Advances in Neural Information Processing Systems*, 33:16315–16326, 2020.

Lucian Busoniu, Robert Babuska, and Bart De Schutter. A comprehensive survey of multiagent reinforcement learning. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 38(2):156–172, 2008.

Lucian Buşoniu, Robert Babuška, and Bart De Schutter. Multi-agent reinforcement learning: An overview. *Innovations in multi-agent systems and applications-1*, pages 183–221, 2010.

Liyu Chen, Rahul Jain, and Haipeng Luo. Learning infinite-horizon average-reward Markov decision process with constraints. In *International Conference on Machine Learning*, pages 3246–3270, 2022a.

Yi Chen, Jing Dong, and Zhaoran Wang. A primal-dual approach to constrained Markov decision processes. *arXiv preprint arXiv:2101.10895*, 2021.

Ziyi Chen, Shaocong Ma, and Yi Zhou. Finding correlated equilibrium of constrained Markov game: A primal-dual approach. In *Advances in Neural Information Processing Systems*, 2022b.

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

Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.

Yu-HOng Dai and Liwei Zhang. Optimality conditions for constrained minimax optimization. *arXiv preprint arXiv:2004.09730*, 2020.

Constantinos Daskalakis, Stratis Skoulakis, and Manolis Zampetakis. The complexity of constrained min-max optimization. In *Annual ACM SIGACT Symposium on Theory of Computing*, pages 1466–1478, 2021.

Frits De Nijs. *Resource-constrained multi-agent Markov decision processes*. PhD thesis, Delft University of Technology, 2019.

Frits de Nijs and Peter J Stuckey. Risk-aware conditional replanning for globally constrained multi-agent sequential decision making. In *International Conference on Autonomous Agents and MultiAgent Systems*, pages 303–311, 2020.

Raghuram Bharadwaj Diddigi, Sai Koti Reddy Danda, Shalabh Bhatnagar, et al. Actor-critic algorithms for constrained multi-agent reinforcement learning. *arXiv preprint arXiv:1905.02907*, 2019.

Dongsheng Ding and Mihailo R Jovanović. Policy gradient primal-dual mirror descent for constrained MDPs with large state spaces. In *2022 IEEE 61st Conference on Decision and Control*, pages 4892–4897, 2022.

Dongsheng Ding, Kaiqing Zhang, Tamer Basar, and Mihailo Jovanovic. Natural policy gradient primal-dual method for constrained Markov decision processes. *Advances in Neural Information Processing Systems*, 33:8378–8390, 2020.

Dongsheng Ding, Xiaohan Wei, Zhuoran Yang, Zhaoran Wang, and Mihailo Jovanovic. Provably efficient safe exploration via primal-dual policy optimization. In *International Conference on Artificial Intelligence and Statistics*, pages 3304–3312, 2021.

Dongsheng Ding, Kaiqing Zhang, Tamer Başar, and Mihailo R Jovanović. Convergence and optimality of policy gradient primal-dual method for constrained Markov decision processes. In *2022 American Control Conference*, pages 2851–2856, 2022a.

Dongsheng Ding, Kaiqing Zhang, Jiali Duan, Tamer Başar, and Mihailo R Jovanović. Convergence and sample complexity of natural policy gradient primal-dual methods for constrained MDPs. *arXiv preprint arXiv:2206.02346*, 2022b.

Simon Du, Sham Kakade, Jason Lee, Shachar Lovett, Gaurav Mahajan, Wen Sun, and Ruosong Wang. Bilinear classes: A structural framework for provable generalization in RL. In *International Conference on Machine Learning*, pages 2826–2836, 2021.

Yonathan Efroni, Shie Mannor, and Matteo Pirotta. Exploration-exploitation in constrained MDPs. *arXiv preprint arXiv:2003.02189*, 2020.

Seyedshams Feyzabadi. *Robot Planning with Constrained Markov Decision Processes*. PhD thesis, UC Merced, 2017.

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

Dylan J Foster, Sham M Kakade, Jian Qian, and Alexander Rakhlin. The statistical complexity of interactive decision making. *arXiv preprint arXiv:2112.13487*, 2021.

Mukul Gagrani and Ashutosh Nayyar. Weakly coupled constrained Markov decision processes in Borel spaces. In *2020 American Control Conference*, pages 2790–2795, 2020.

Javier Garcıa and Fernando Fernández. A comprehensive survey on safe reinforcement learning. *Journal of Machine Learning Research*, 16(1):1437–1480, 2015.

Cory Jay Girard. *STRUCTURAL RESULTS FOR CONSTRAINED MARKOV DECISION PROCESSES*. PhD thesis, Cornell University, 2018.

E Gómez-Ramırez, K Najim, and AS Poznyak. Saddle-point calculation for constrained finite Markov chains. *Journal of Economic Dynamics and Control*, 27(10):1833–1853, 2003.

Shangding Gu, Jakub Grudzien Kuba, Munning Wen, Ruiqing Chen, Ziyan Wang, Zheng Tian, Jun Wang, Alois Knoll, and Yaodong Yang. Multi-agent constrained policy optimisation. *arXiv preprint arXiv:2110.02793*, 2021.

András György, Tamás Linder, Gábor Lugosi, and György Ottucsák. The on-line shortest path problem under partial monitoring. *Journal of Machine Learning Research*, 8(10), 2007.

Vesal Hakami and Mehdi Dehghan. Learning stationary correlated equilibria in constrained general-sum stochastic games. *IEEE transactions on cybernetics*, 46(7):1640–1654, 2015.

Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11(4), 2010.

Xiaofeng Jiang, Shuangwu Chen, Jian Yang, Han Hu, and Zhenliang Zhang. Finding the equilibrium for continuous constrained Markov games under the average criteria. *IEEE Transactions on Automatic Control*, 65(12):5399–5406, 2020.

Chi Jin, Tiancheng Jin, Haipeng Luo, Suvrit Sra, and Tiancheng Yu. Learning adversarial Markov decision processes with bandit feedback and unknown transition. In *International Conference on Machine Learning*, pages 4860–4869, 2020.

Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman Eluder dimension: New rich classes of RL problems, and sample-efficient algorithms. *Advances in Neural Information Processing Systems*, 34, 2021.

Chi Jin, Qinghua Liu, Yuanhao Wang, and Tiancheng Yu. V-learning– A simple, efficient, decentralized algorithm for multiagent RL. In *ICLR 2022 Workshop on Gamification and Multiagent Solutions*, 2022a.

Chi Jin, Qinghua Liu, and Tiancheng Yu. The power of exploiter: Provable multi-agent RL in large state spaces. In *International Conference on Machine Learning*, pages 10251–10279, 2022b.

Krishna C Kalagarla, Rahul Jain, and Pierluigi Nuzzo. A sample-efficient algorithm for episodic finite-horizon MDP with constraints. In *AAAI Conference on Artificial Intelligence*, 2020.

Ankur A Kulkarni. *Generalized Nash games with shared constraints: existence, efficiency, refinement and equilibrium constraints*. PhD thesis, University of Illinois at Urbana-Champaign, 2011.

Ankur A Kulkarni. Games and teams with shared constraints. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 375(2100):20160302, 2017.

Na Li and Jason R Marden. Decoupling coupled constraints through utility design. *IEEE Transactions on Automatic Control*, 59(8):2289–2294, 2014.

Tianjiao Li, Ziwei Guan, Shaofeng Zou, Tengyu Xu, Yingbin Liang, and Guanghui Lan. Faster algorithm and sharper analysis for constrained Markov decision process. *arXiv preprint arXiv:2110.10351*, 2021.

Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 157–163, 1994.

Tao Liu, Ruida Zhou, Dileep Kalathil, Panganamala Kumar, and Chao Tian. Learning policies with zero or bounded constraint violation for constrained MDPs. *Advances in Neural Information Processing Systems*, 34:17183–17193, 2021a.

Tao Liu, Ruida Zhou, Dileep Kalathil, PR Kumar, and Chao Tian. Fast global convergence of policy optimization for constrained MDPs. *arXiv preprint arXiv:2111.00552*, 2021b.

Songtao Lu, Kaiqing Zhang, Tianyi Chen, Tamer Basar, and Lior Horesh. Decentralized policy gradient descent ascent for safe multi-agent reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2020.

Shie Mannor, John N Tsitsiklis, and Jia Yuan Yu. Online learning with sample path constraints. *Journal of Machine Learning Research*, 10(3), 2009.

Nicolas Meuleau, Milos Hauskrecht, Kee-Eung Kim, Leonid Peshkin, Leslie Pack Kaelbling, Thomas L Dean, and Craig Boutilier. Solving very large weakly coupled Markov decision processes. In *AAAI/IAAI*, pages 165–172, 1998.

Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.

Gergely Neu, András György, and Csaba Szepesvári. The online loop-free stochastic shortest-path problem. In *COLT*, volume 2010, pages 231–243, 2010.

Gergely Neu, Andras Gyorgy, and Csaba Szepesvári. The adversarial stochastic shortest path problem with unknown transition probabilities. In *Artificial Intelligence and Statistics*, pages 805–813, 2012.

J v Neumann. Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320, 1928.

Duc Thien Nguyen, William Yeoh, Hoong Chuin Lau, Shlomo Zilberstein, and Chongjie Zhang. Decentralized multi-agent reinforcement learning in average-reward dynamic DCOPs. In *AAAI conference on artificial intelligence*, 2014.

Afshin OroojlooyJadid and Davood Hajinezhad. A review of cooperative multi-agent deep reinforcement learning. *arXiv preprint arXiv:1908.03963*, 2019.

P Parnika, Raghuram Bharadwaj Diddigi, Sai Koti Reddy Danda, and Shalabh Bhatnagar. Attention actor-critic algorithm for multi-agent constrained co-operative reinforcement learning. *arXiv preprint arXiv:2101.02349*, 2021.

Edward S Pearsall. A Lagrange multiplier method for certain constrained min-max problems. *Operations Research*, 24(1):70–91, 1976.

Alexei B Piunovskiy and Xuerong Mao. Constrained Markovian decision processes: the dynamic programming approach. *Operations research letters*, 27(3):119–126, 2000.

Shuang Qiu, Xiaohan Wei, Zhuoran Yang, Jieping Ye, and Zhaoran Wang. Upper confidence primal-dual reinforcement learning for CMDP with adversarial loss. *Advances in Neural Information Processing Systems*, 33:15277–15287, 2020.

Alexander Rakhlin and Karthik Sridharan. Online learning with predictable sequences. In *Conference on Learning Theory*, pages 993–1019, 2013.

J Ben Rosen. Existence and uniqueness of equilibrium points for concave n-person games. *Econometrica: Journal of the Econometric Society*, pages 520–534, 1965.

Aviv Rosenberg and Yishay Mansour. Online stochastic shortest path with bandit feedback and unknown transition function. In *Advances in Neural Information Processing Systems*, pages 2212–2221, 2019.

Mahshid Salemi Parizi. *Approximate dynamic programming for weakly coupled Markov decision processes with perfect and imperfect information*. PhD thesis, The University of Washington, 2018.

Lukas M Schmidt, Johanna Brosig, Axel Plinge, Bjoern M Eskofier, and Christopher Mutschler. An introduction to multi-agent reinforcement learning and review of its application to autonomous mobility. *arXiv preprint arXiv:2203.07676*, 2022.

Shai Shalev-Shwartz, Shaked Shammah, and Amnon Shashua. Safe, multi-agent, reinforcement learning for autonomous driving. *arXiv preprint arXiv:1610.03295*, 2016.

Lloyd S Shapley. Stochastic games. *Proceedings of the National Academy of Sciences*, 39(10): 1095–1100, 1953.

Rahul Singh, Abhishek Gupta, and Ness Shroff. Learning in Markov decision processes under constraints. *IEEE Transactions on Control of Network Systems*, 2022.

Vikas Vikram Singh and N Hemachandra. A characterization of stationary Nash equilibria of constrained stochastic games with independent state processes. *Operations Research Letters*, 42 (1):48–52, 2014.

Ziang Song, Song Mei, and Yu Bai. When can we learn general-sum Markov games with a large number of players sample-efficiently? In *International Conference on Learning Representations*, 2021.

Chen Tessler, Daniel J Mankowitz, and Shie Mannor. Reward constrained policy optimization. In *International Conference on Learning Representations*, 2019.

Philip S Thomas. *Safe reinforcement learning*. PhD thesis, University of Massachusetts Libraries, 2015.

Yi Tian, Yuanhao Wang, Tiancheng Yu, and Suvrit Sra. Provably efficient online agnostic learning in Markov games. *arXiv preprint arXiv:2010.15020*, 2020.

Paul Tseng. On accelerated proximal gradient methods for convex-concave optimization. *URL http://www. math. washington. edu/˜ tseng/papers/apgm. pdf*, 2009.

Akifumi Wachi and Yanan Sui. Safe reinforcement learning in constrained Markov decision processes. In *International Conference on Machine Learning*, pages 9797–9806, 2020.

Chen-Yu Wei, Yi-Te Hong, and Chi-Jen Lu. Online reinforcement learning in stochastic games. In *Advances in Neural Information Processing Systems*, pages 4994–5004, 2017.

Qingda Wei. Discrete-time constrained stochastic games with the expected average payoff criteria. *Optimization*, pages 1–32, 2020.

Qingda Wei. Constrained expected average stochastic games for continuous-time jump processes. *Applied Mathematics & Optimization*, 83(3):1277–1309, 2021.

Xiaohan Wei, Hao Yu, and Michael J Neely. Online learning in weakly coupled Markov decision processes: A convergence time study. *ACM on Measurement and Analysis of Computing Systems*, 2(1):1–38, 2018.

Xiaohan Wei, Hao Yu, and Michael J Neely. Online primal-dual mirror descent under stochastic constraints. In *Abstracts of the 2020 SIGMETRICS/Performance Joint International Conference on Measurement and Modeling of Computer Systems*, pages 3–4, 2020.

Qiaomin Xie, Yudong Chen, Zhaoran Wang, and Zhuoran Yang. Learning zero-sum simultaneous-move Markov games using function approximation and correlated equilibrium. In *Conference on Learning Theory*, pages 3674–3682, 2020.

Vinayaka G Yaji and Shalabh Bhatnagar. Necessary and sufficient conditions for optimality in constrained general sum stochastic games. *Systems & Control Letters*, 85:8–15, 2015.

Yaodong Yang and Jun Wang. An overview of multi-agent reinforcement learning from game theoretical perspective. *arXiv preprint arXiv:2011.00583*, 2020.

Donghao Ying, Yuhao Ding, and Javad Lavaei. A dual approach to constrained Markov decision processes with entropy regularization. In *International Conference on Artificial Intelligence and Statistics*, pages 1887–1909, 2022.

Hao Yu, Michael Neely, and Xiaohan Wei. Online convex optimization with stochastic constraints. In *Advances in Neural Information Processing Systems*, pages 1428–1438, 2017.

Kaiqing Zhang, Zhuoran Yang, and Tamer Başar. Multi-agent reinforcement learning: A selective overview of theories and algorithms. *Handbook of Reinforcement Learning and Control*, pages 321–384, 2021.

Wenzhao Zhang. Discrete-time constrained average stochastic games with independent state processes. *Mathematics*, 7(11):1089, 2019.

Wenzhao Zhang and Xiaolong Zou. Constrained average stochastic games with continuous-time independent state processes. *Optimization*, pages 1–24, 2021.

Feiran Zhao and Keyou You. Primal-dual learning for the model-free risk-constrained linear quadratic regulator. In *Learning for Dynamics and Control*, pages 702–714, 2021.

Liyuan Zheng and Lillian Ratliff. Constrained upper confidence reinforcement learning. In *Conference on Learning for Dynamics and Control*, volume 120, pages 620–629, 2020.

Alexander Zimin and Gergely Neu. Online learning in episodic Markovian decision processes by relative entropy policy search. In *Neural Information Processing Systems 26*, 2013.