

Fast Policy Extragradient Methods for Competitive Games with Entropy Regularization

Shicong Cen*
CMU

Yuting Wei†
CMU

Yuejie Chi‡
CMU

May 29, 2021

Abstract

This paper investigates the problem of computing the equilibrium of competitive games, which is often modeled as a constrained saddle-point optimization problem with probability simplex constraints. Despite recent efforts in understanding the last-iterate convergence of extragradient methods in the unconstrained setting, the theoretical underpinnings of these methods in the constrained settings, especially those using multiplicative updates, remain highly inadequate, even when the objective function is bilinear. Motivated by the algorithmic role of entropy regularization in single-agent reinforcement learning and game theory, we develop provably efficient extragradient methods to find the quantal response equilibrium (QRE)—which are solutions to zero-sum two-player matrix games with entropy regularization—at a linear rate. The proposed algorithms can be implemented in a decentralized manner, where each player executes symmetric and multiplicative updates iteratively using its own payoff without observing the opponent’s actions directly. In addition, by controlling the knob of entropy regularization, the proposed algorithms can locate an approximate Nash equilibrium of the unregularized matrix game at a sublinear rate without assuming the Nash equilibrium to be unique. Our methods also lead to efficient policy extragradient algorithms for solving entropy-regularized zero-sum Markov games at a linear rate. All of our convergence rates are nearly dimension-free, which are independent of the size of the state and action spaces up to logarithm factors, highlighting the positive role of entropy regularization for accelerating convergence.

Keywords: zero-sum Markov game, matrix game, entropy regularization, global convergence, multiplicative updates, extragradient methods

Contents

1	Introduction	2
1.1	Last-iterate convergence in competitive games	2
1.2	Our contributions	3
1.3	Related works	4
1.4	Notation	5
2	Zero-sum matrix games with entropy regularization	5
2.1	Background and problem formulation	5
2.2	Proposed extragradient methods: PU and OMWU	6
2.3	Performance guarantees	7
3	Zero-sum Markov games with entropy regularization	10
3.1	Background and problem formulation	10
3.2	From value iteration to policy extragradient methods	11

*Department of Electrical and Computer Engineering, Carnegie Mellon University; email: shicongc@andrew.cmu.edu.

†Department of Statistics and Data Science, Carnegie Mellon University; email: ytwei@cmu.edu.

‡Department of Electrical and Computer Engineering, Carnegie Mellon University; email: yuejiechi@cmu.edu.

4	Conclusions	13
A	Analysis for entropy-regularized matrix games	13
A.1	Proof of Proposition 1	14
A.2	Proof of Theorem 1	15
B	Analysis for entropy-regularized Markov games	21
B.1	Proof of Proposition 2	21
B.2	Proof of Theorem 2	22
C	Proof of auxiliary lemmas	23
C.1	Proof of Lemma 1	23
C.2	Proof of Lemma 2	24
C.3	Proof of Lemma 3	24
C.4	Proof of Lemma 4	24

1 Introduction

Finding the equilibrium of competitive games, which can be viewed as constrained saddle-point optimization problems with probability simplex constraints, lies at the heart of modern machine learning and decision making paradigms such as Generative Adversarial Networks (GANs) (Goodfellow et al., 2014), competitive reinforcement learning (RL) (Littman, 1994), game theory (Shapley, 1953), adversarial training (Mertikopoulos et al., 2018b), to name a few.

In this paper, we study one of the most basic forms of competitive games, namely two-player zero-sum games, in both the matrix setting and the Markov setting. Our goal is to find the equilibrium policies of both players in an *independent* and *decentralized* manner (Daskalakis et al., 2020; Wei et al., 2021a) with guaranteed *last-iterate convergence*. Namely, each player will execute symmetric and independent updates iteratively using its own payoff without observing the opponent’s actions directly, and the final policies of the iterative process should be a close approximation to the equilibrium up to any prescribed precision. This kind of algorithms is more advantageous and versatile especially in federated environments, as it requires neither prior coordination between the players like two-timescale algorithms, nor a central controller to collect and disseminate the policies of all the players, which are often unavailable due to privacy constraints.

1.1 Last-iterate convergence in competitive games

In recent years, there have been significant progresses in understanding the last-iterate convergence of simple iterative algorithms for *unconstrained* saddle-point optimization, where one is interested in bounding the sub-optimality of the last iterate of the algorithm, rather than say, the ergodic iterate — which is the average of all the iterations — that are commonly studied in the earlier literature. This shift of focus is motivated, for example, by the infeasibility of averaging large machine learning models in training GANs (Goodfellow et al., 2014). While vanilla Gradient Descent / Ascent (GDA) may diverge or cycle even for bilinear matrix games (Daskalakis et al., 2018), quite remarkably, small modifications lead to guaranteed last-iterate convergence to the equilibrium in a non-asymptotic fashion. A flurry of algorithms is proposed, including Optimistic Gradient Descent Ascent (OGDA) (Rakhlin and Sridharan, 2013; Daskalakis and Panageas, 2018b; Wei et al., 2021b), predictive updates (Yadav et al., 2017), implicit updates (Liang and Stokes, 2019), and more. Several unified analyses of these algorithms have been carried out (see, e.g. Mokhtari et al. (2020a); Liang and Stokes (2019) and references therein), where these methods in principle all make clever extrapolation of the local curvature in a predictive manner to accelerate convergence. With slight abuse of terminology, in this paper, we refer to this ensemble of algorithms as extragradient methods (Korpelevich, 1976; Tseng, 1995; Mertikopoulos et al., 2018a; Harker and Pang, 1990).

However, saddle-point optimization in the *constrained setting*, which includes competitive games as a special case, remains largely under-explored even for bilinear matrix games. While it is possible to reformulate constrained bilinear games to unconstrained ones using softmax parameterization of the probability simplex, this approach falls short of preserving the bilinear structure and convex-concave properties

in the original problem, which are crucial to the convergence of gradient methods. Therefore, there is a strong necessity of understanding and developing improved extragradient methods in the constrained setting. [Daskalakis and Panageas \(2018a\)](#) proposed the optimistic variant of the multiplicative weight updates (MWU) method ([Arora et al., 2012](#))—which is extremely natural and popular for optimizing over probability simplexes—called Optimistic Multiplicative Weight Updates (OMWU), and established the asymptotic last-iterate convergence of OMWU for matrix games. Very recently, [Wei et al. \(2021b\)](#) established non-asymptotic last-iterate convergences of OMWU. However, these last-iterate convergence results require the Nash equilibrium to be unique, and cannot be applied to problems with multiple Nash equilibria.

1.2 Our contributions

Motivated by the algorithmic role of entropy regularization in single-agent RL ([Neu et al., 2017](#); [Geist et al., 2019](#); [Cen et al., 2020](#)) as well as its wide use in game theory to account for imperfect and noisy information ([McKelvey and Palfrey, 1995](#); [Savas et al., 2019](#)), we initiate the design and analysis of extragradient algorithms using *multiplicative updates* for finding the so-called quantal response equilibrium (QRE), which are solutions to competitive games with entropy regularization ([McKelvey and Palfrey, 1995](#)). While finding QRE is of interest in its own right, by controlling the knob of entropy regularization, the QRE provides a close approximation to the Nash equilibrium (NE), and in turn acts as a smoothing scheme for finding the NE. Our contributions are summarized below.

- *Near dimension-free last-iterate convergence to QRE of entropy-regularized matrix games.* We propose two policy extragradient algorithms to solve entropy-regularized matrix games, namely the Predictive Update (PU) and OMWU methods, where both players execute symmetric and multiplicative updates without knowing the entire payoff matrix nor the opponent’s actions. Encouragingly, we show that the last iterate of the proposed algorithms converges to the unique QRE at a linear rate that is almost independent of the size of the action spaces. Roughly speaking, to find an ϵ -optimal QRE in terms of Kullback-Leibler (KL) divergence, it takes no more than

$$\tilde{O}\left(\frac{1}{\eta\tau} \log\left(\frac{1}{\epsilon}\right)\right)$$

iterations, where $\tilde{O}(\cdot)$ hides logarithmic dependencies. Here, τ is the regularization parameter, and η is the learning rate of both players. Optimizing the learning rate, the iteration complexity is bounded by $\tilde{O}((1 + \|A\|_\infty/\tau) \log(1/\epsilon))$, where $\|A\|_\infty = \max_{i,j} |A_{i,j}|$ is the ℓ_∞ norm of the payoff matrix A .

- *Last-iterate convergence to ϵ -NE of unregularized matrix games without uniqueness assumption.* The QRE provides an accurate approximation to the NE by setting the entropy regularization τ sufficiently small, therefore our result directly translates to finding a NE with last-iterate convergence guarantee. Roughly speaking, to find an ϵ -NE ([Zhang et al., 2020](#), Definition 2.1), it takes no more than

$$\tilde{O}\left(\frac{\|A\|_\infty}{\epsilon}\right)$$

iterations with optimized learning rates, again independent of the size of the action spaces up to logarithmic factors. Unlike prior literature ([Daskalakis and Panageas, 2018a](#); [Wei et al., 2021b](#)), our last-iterate convergence guarantee does not require the NE to be unique.

- *Extensions to two-player zero-sum Markov games.* By connecting value iteration with matrix games, we propose a policy extragradient method for solving infinite-horizon discounted entropy-regularized zero-sum Markov games, which finds an ϵ -optimal minimax soft Q-function—in terms of ℓ_∞ error—in at most $\tilde{O}\left(\frac{1}{\tau(1-\gamma)^2} \log^2\left(\frac{1}{\epsilon}\right)\right)$ iterations, where $\gamma \in (0, 1)$ is the discount factor.

To the best of our knowledge, our paper is the first one that develops policy extragradient algorithms for solving entropy-regularized competitive games with multiplicative updates and dimension-free linear last-iterate convergence, and demonstrates entropy regularization as a smoothing technique to find ϵ -NE without the uniqueness assumption. Table 1 provides a detailed comparison of the proposed entropy-regularized PU

Equilibrium type	Method	Convergence rate	Dimension-free	Require unique NE
ϵ -QRE	PU & OMWU (this work)	linear	yes	n/a
ϵ -NE	OMWU (Daskalakis and Panageas, 2018a)	asymptotic	no	yes
	OMWU (Wei et al., 2021b)	sublinear + linear	no	yes
	PU & OMWU (this work)	sublinear	yes	no

Table 1: Comparisons of last-iterate convergence of the proposed entropy-regularized PU and OMWU methods with prior results for finding ϵ -QRE or ϵ -NE of competitive matrix games. We note that the convergence rates of unregularized OMWU established in Wei et al. (2021b) are problem-dependent, and scale at least polynomially on the size of the action spaces. Desirable features in the last two columns are highlighted in blue.

and OMWU methods with prior last-iterate convergence guarantees of unregularized OMWU. Our results highlight the positive role of entropy regularization for accelerating convergence and safeguarding against imperfect information in competitive games.

1.3 Related works

Our work lies at the intersection of saddle-point optimization, game theory, and reinforcement learning. In what follows, we discuss a few topics that are closely related to ours.

Unregularized matrix game. Freund and Schapire (1999) showed that many standard methods such as GDA and MWU have a converging average duality gap at the rate of $O(1/\sqrt{T})$, which is improved to $O(1/T)$ by considering optimistic variants of these methods, such as OGDA and OMWU (Rakhlin and Sridharan, 2013; Daskalakis et al., 2011; Syrgkanis et al., 2015). However, the last-iterate convergence of these methods are less understood until recently (Daskalakis and Panageas, 2018a; Wei et al., 2021b). In particular, under the assumption that the NE is unique for the unregularized matrix game, Daskalakis and Panageas (2018a) showed the asymptotic convergence of the last iterate of OMWU to the unique equilibrium, and Wei et al. (2021b) showed the last iterate of OMWU achieves a linear rate of convergence after an initial phase of sublinear convergence, however the rates therein can be highly pessimistic in terms of the problem dimension, while our rate for entropy-regularized OMWU is dimension-free up to logarithmic factors.

Saddle-point optimization. Considerable progress has been made towards understanding OGDA and extragradient (EG) methods in the unconstrained convex-concave saddle-point optimization with general objective functions (Mokhtari et al., 2020a,b; Nemirovski, 2004; Liang and Stokes, 2019). However, the last-iterate convergence of constrained convex-concave saddle-point optimization still lacks theoretical understanding in general and most works fall short of characterizing a finite-time convergence result. In particular, Mertikopoulos et al. (2018a) demonstrated the asymptotic last-iterate convergence of EG, and Hsieh et al. (2019) investigated similar questions for single-call EG algorithms. Lei et al. (2021) showed that OMWU converges to the equilibrium locally without an explicit rate. Wei et al. (2021b) showed that the last-iterate of OGDA converges linearly for strongly-convex strongly-concave constrained saddle-point optimization with an explicit rate.

Entropy regularization in RL and games. In single-agent RL, the role of entropy regularization as an algorithmic mechanism to encourage exploration and accelerate convergence has been investigated extensively (Neu et al., 2017; Geist et al., 2019; Mei et al., 2020; Cen et al., 2020; Lan, 2021; Zhan et al., 2021). Turning

to the game setting, entropy regularization is used to account for imperfect information in the seminal work of [McKelvey and Palfrey \(1995\)](#) that introduced the QRE, and a few representative works on entropy and more general regularizations in games include [Savas et al. \(2019\)](#); [Hofbauer and Sandholm \(2002\)](#); [Mertikopoulos and Sandholm \(2016\)](#).

Zero-sum Markov games. There have been a significant recent interest in developing provably efficient self-play algorithms for Markov games, including model-based algorithms ([Perolat et al., 2015](#); [Zhang et al., 2020](#)), value-based algorithms ([Bai and Jin, 2020](#); [Xie et al., 2020](#)), and policy-based algorithms ([Daskalakis et al., 2020](#); [Wei et al., 2021a](#); [Zhao et al., 2021](#)). Our approach can be regarded as a policy-based algorithm to approximate value iteration, which can be implemented in a decentralized manner with symmetric and multiplicative updates from both players, and the iteration complexity is almost independent of the size of the state-action space.

1.4 Notation

We denote by $\Delta(\mathcal{A})$ the probability simplex over the set \mathcal{A} . We overload the functions such as $\log(\cdot)$ and $\exp(\cdot)$ to take vector inputs with the understanding that the function is applied in an entrywise manner. For instance, given any vector $z = [z_i]_{1 \leq i \leq n} \in \mathbb{R}^n$, the notation $\exp(z)$ denotes $\exp(z) := [\exp(z_i)]_{1 \leq i \leq n}$; other functions are defined analogously. Given two probability distributions μ and μ' over \mathcal{A} , the KL divergence from μ' to μ is defined by $\text{KL}(\mu \parallel \mu') := \sum_{a \in \mathcal{A}} \mu(a) \log \frac{\mu(a)}{\mu'(a)}$. Given a matrix A , $\|A\|_\infty$ is used to denote entrywise maximum norm, namely, $\|A\|_\infty = \max_{i,j} |A_{i,j}|$. The all-one vector is denoted as $\mathbf{1}$.

2 Zero-sum matrix games with entropy regularization

In this section, we consider a two-player zero-sum game with bilinear objective and probability simplex constraints, and demonstrate the positive role of entropy regularization in solving this problem. Throughout this paper, let $\mathcal{A} = \{1, \dots, m\}$ and $\mathcal{B} = \{1, \dots, n\}$ be the action spaces of each player. The proofs for this section are collected in Appendix [A](#).

2.1 Background and problem formulation

Zero-sum two-player matrix game. The focal point of this subsection is a constrained two-player zero-sum matrix game, which can be formulated as the following min-max problem (or saddle point optimization problem):

$$\max_{\mu \in \Delta(\mathcal{A})} \min_{\nu \in \Delta(\mathcal{B})} f(\mu, \nu) := \mu^\top A \nu, \quad (1)$$

where $A \in \mathbb{R}^{m \times n}$ denotes the payoff matrix, $\mu \in \Delta(\mathcal{A})$ and $\nu \in \Delta(\mathcal{B})$ stand for the mixed/randomized policies of each player, defined respectively as distributions over the probability simplex $\Delta(\mathcal{A})$ and $\Delta(\mathcal{B})$. It is well known since [Neumann \(1928\)](#) that the max and min operators in (1) can be exchanged without affecting the solution. A pair of policies (μ^*, ν^*) is said to be a *Nash equilibrium (NE)* of (1) if

$$f(\mu^*, \nu) \geq f(\mu^*, \nu^*) \geq f(\mu, \nu^*) \quad \text{for all } (\mu, \nu) \in \Delta(\mathcal{A}) \times \Delta(\mathcal{B}). \quad (2)$$

In words, the NE corresponds to when both players play their best-response strategies against their respective opponents.

Entropy-regularized zero-sum two-player matrix game. There is no shortage of scenarios where the payoff matrix A might not be known perfectly. In an attempt to accommodate imperfect knowledge of A , [McKelvey and Palfrey \(1995\)](#) proposed a seminal extension to the Nash equilibrium called the *quantal response equilibrium (QRE)* when the payoffs are perturbed by Gumbel-distributed noise. Formally, this amounts to solving the following matrix game with entropy regularization ([Mertikopoulos and Sandholm, 2016](#)):

$$\max_{\mu \in \Delta(\mathcal{A})} \min_{\nu \in \Delta(\mathcal{B})} f_\tau(\mu, \nu) := \mu^\top A \nu + \tau \mathcal{H}(\mu) - \tau \mathcal{H}(\nu), \quad (3)$$

where $\mathcal{H}(\pi) = -\sum_i \pi_i \log(\pi_i)$ denotes the Shannon entropy of a distribution π , and $\tau \geq 0$ is the regularization parameter. As is well known, the optimal solution (μ_τ^*, ν_τ^*) to (3), dubbed as the QRE, is unique whenever $\tau > 0$ (due to the presence of strong concavity/convexity), which satisfies the following fixed point equations:

$$\begin{cases} \mu_\tau^*(a) = \frac{\exp([A\nu_\tau^*]_a/\tau)}{\sum_{a=1}^m \exp([A\nu_\tau^*]_a/\tau)} \propto \exp([A\nu_\tau^*]_a/\tau), & \text{for all } a \in \mathcal{A}, \\ \nu_\tau^*(b) = \frac{\exp(-[A^\top \mu_\tau^*]_b/\tau)}{\sum_{b=1}^n \exp(-[A^\top \mu_\tau^*]_b/\tau)} \propto \exp(-[A^\top \mu_\tau^*]_b/\tau), & \text{for all } b \in \mathcal{B}. \end{cases} \quad (4)$$

Goal. We aim to efficiently compute the QRE of the entropy-regularized matrix game in a decentralized manner, and investigate how an efficient solver of QRE can be leveraged to find a NE of the unregularized matrix game (1). Namely, we only assume access to “first-order information” as opposed to full knowledge of the payoff matrix A or the actions of the opponent. The information received by each player is formally described in the following sampling oracle.

Definition 1 (Sampling oracle for matrix games). *For any policy pair (μ, ν) and payoff matrix A , the sampling oracle returns the exact values of $\mu^\top A$ and $A\nu$.*

Additional notation. For notational convenience, we let ζ represent the concatenation of $\mu \in \mathbb{R}^{|\mathcal{A}|}$ and $\nu \in \mathbb{R}^{|\mathcal{B}|}$, namely, $\zeta = (\mu, \nu)$. The solution to (3), which is specified in (4), is denoted by $\zeta_\tau^* = (\mu_\tau^*, \nu_\tau^*)$. For any $\zeta = (\mu, \nu)$ and $\zeta' = (\mu', \nu')$, we shall often abuse the notation and let

$$\text{KL}(\zeta \parallel \zeta') = \text{KL}(\mu \parallel \mu') + \text{KL}(\nu \parallel \nu').$$

The duality gap of the entropy-regularized matrix game (3) at $\zeta = (\mu, \nu)$ is defined as

$$\text{DualGap}_\tau(\zeta) = \max_{\mu' \in \Delta(\mathcal{A})} f_\tau(\mu', \nu) - \min_{\nu' \in \Delta(\mathcal{B})} f_\tau(\mu, \nu') \quad (5)$$

which is clearly nonnegative and $\text{DualGap}_\tau(\zeta_\tau^*) = 0$. Similarly, let the optimality gap of the entropy-regularized matrix game (3) at $\zeta = (\mu, \nu)$ be $\text{OptGap}(\zeta) = |f_\tau(\mu, \nu) - f_\tau(\mu_\tau^*, \nu_\tau^*)|$.

2.2 Proposed extragradient methods: PU and OMWU

To begin, assume we are given a pair of policies $z_1 \in \Delta(\mathcal{A})$, $z_2 \in \Delta(\mathcal{B})$ employed by each player respectively. If we proceed with fictitious play, i.e. player 1 (resp. player 2) aims to optimize its own policy by assuming the opponent’s policy is fixed as z_2 (resp. z_1), the saddle-point optimization problem (3) is then decoupled into two independent min/max optimization problems:

$$\max_{\mu \in \Delta(\mathcal{A})} \mu^\top A z_2 + \tau \mathcal{H}(\mu) - \tau \mathcal{H}(z_1) \quad \text{and} \quad \min_{\nu \in \Delta(\mathcal{B})} z_1^\top A \nu + \tau \mathcal{H}(z_2) - \tau \mathcal{H}(\nu),$$

which are naturally solved via mirror descent/ascent with KL divergence. Specifically, one step of mirror descent/ascent takes the form

$$\begin{cases} \mu^{(t+1)} = \arg \max_{\mu \in \Delta(\mathcal{A})} (A z_2 - \tau \log \mu^{(t)})^\top \mu - \frac{1}{\eta} \text{KL}(\mu \parallel \mu^{(t)}) \\ \nu^{(t+1)} = \arg \min_{\nu \in \Delta(\mathcal{B})} (A^\top z_1 + \tau \log \nu^{(t)})^\top \nu + \frac{1}{\eta} \text{KL}(\nu \parallel \nu^{(t)}) \end{cases},$$

where η is the learning rate, or equivalently

$$\begin{cases} \mu^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A z_2]_a), & \text{for all } a \in \mathcal{A}, \\ \nu^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top z_1]_b), & \text{for all } b \in \mathcal{B}. \end{cases} \quad (6)$$

The above update rule forms the basis of our algorithm design.

Motivation: a form of implicit updates with linear convergence. It turns out, if we could select the policy pair $(z_1, z_2) = \zeta^{(t+1)} := (\mu^{(t+1)}, \nu^{(t+1)})$ as the ones to be taken in the future, and call the resulting update rule as the Implicit Update (IU) method:

$$\text{Implicit Update: } \begin{cases} \mu^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\nu^{(t+1)}]_a), & \text{for all } a \in \mathcal{A}, \\ \nu^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \mu^{(t+1)}]_b), & \text{for all } b \in \mathcal{B}. \end{cases} \quad (7)$$

Though unrealistic — since it uses the future updates — it leads to a one-step convergence to the QRE when $\eta = 1/\tau$ (see the optimality condition in (4)). Encouragingly, we have the following linear convergence guarantee of IU when adopting a general learning rate.

Proposition 1 (Linear convergence of IU). *Assume $0 < \eta \leq 1/\tau$, then for all $t \geq 0$, the iterates $\zeta^{(t)} := (\mu^{(t)}, \nu^{(t)})$ of the IU method in (7) satisfy*

$$\text{KL}(\zeta_t^* \parallel \zeta^{(t)}) \leq (1 - \eta\tau)^t \text{KL}(\zeta_t^* \parallel \zeta^{(0)}).$$

In words, the IU method achieves an appealing linear rate of convergence that is independent of the problem dimension. Motivated by this observation, we seek to design algorithms where the policies (z_1, z_2) employed in (6) serve as good predictions of $(\mu^{(t+1)}, \nu^{(t+1)})$, such that the resulting algorithms are both practical and retain the appealing convergence rate of IU.

Proposed algorithms. We propose two extragradient algorithms for solving the entropy-regularized matrix game, namely the *Predictive Update (PU)* method and the *Optimistic Multiplicative Weights Update (OMWU)* method, the latter adapted from Rakhlin and Sridharan (2013); Daskalakis et al. (2011). Detailed procedures can be found in Algorithm 1 and Algorithm 2, respectively. On a high level, both algorithms maintain two intertwined sequences $\{(\mu^{(t)}, \nu^{(t)})\}_{t \geq 0}$ and $\{(\bar{\mu}^{(t)}, \bar{\nu}^{(t)})\}_{t \geq 0}$, and in each iteration $t = 0, 1, \dots$, proceed in two steps:

- The midpoint $(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})$ serves as a prediction of $(\mu^{(t+1)}, \nu^{(t+1)})$ by running one step of mirror descent / ascent (cf. (6)) from either $(z_1, z_2) = (\mu^{(t)}, \nu^{(t)})$ (for PU) or $(z_1, z_2) = (\bar{\mu}^{(t)}, \bar{\nu}^{(t)})$ (for OMWU).
- The update of $(\mu^{(t+1)}, \nu^{(t+1)})$ then mimics the implicit update (7) using the prediction $(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})$ obtained above.

When the proposed algorithms converge, both $(\mu^{(t)}, \nu^{(t)})$ and $(\bar{\mu}^{(t)}, \bar{\nu}^{(t)})$ converge to the same point. The two players are completely symmetric and adopt the same learning rate, and require *only* first-order information provided by the sampling oracle. While the two algorithms resemble each other in many aspects, a key difference lies in the query and use of the sampling oracle: in each iteration, OMWU makes a single call to the sampling oracle for gradient evaluation, while PU calls the sampling oracle twice. It is worth noting that, when $\tau = 0$ (i.e., no entropy regularization is enforced), the OMWU method in Algorithm 2 reduces to the method analyzed in Rakhlin and Sridharan (2013); Daskalakis and Panageas (2018a); Wei et al. (2021b) without entropy regularization.

2.3 Performance guarantees

We are now positioned to present our main theorem concerning the last-iterate convergence of PU and OMWU for solving (3).

Theorem 1 (Last-iterate convergence of PU and OMWU). *Suppose that the learning rates $\eta = \eta_{\text{PU}}$ of PU in Algorithm 1 and $\eta = \eta_{\text{OMWU}}$ of OMWU in Algorithm 2 satisfy*

$$0 < \eta_{\text{PU}} \leq \frac{1}{\tau + 2\|A\|_\infty}, \quad \text{and} \quad 0 < \eta_{\text{OMWU}} \leq \min \left\{ \frac{1}{2\tau + 2\|A\|_\infty}, \frac{1}{4\|A\|_\infty} \right\}. \quad (8)$$

Then for any $t \geq 0$, the iterates $\zeta^{(t)} = (\mu^{(t)}, \nu^{(t)})$ and $\bar{\zeta}^{(t)} = (\bar{\mu}^{(t)}, \bar{\nu}^{(t)})$ of both PU and OMWU achieve

Algorithm 1: The PU method	Algorithm 2: The OMWU method
1 initialization: $\mu^{(0)}, \nu^{(0)}$. 2 for $t = 0, 1, 2, \dots$ do 3 Update $\bar{\mu}$ and $\bar{\nu}$ according to $\begin{cases} \bar{\mu}^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\nu^{(t)}]_a), \\ \bar{\nu}^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \mu^{(t)}]_b). \end{cases}$ 4 Update μ and ν according to $\begin{cases} \mu^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\bar{\nu}^{(t+1)}]_a), \\ \nu^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \bar{\mu}^{(t+1)}]_b). \end{cases}$	1 initialization: $\bar{\mu}^{(0)}, \bar{\nu}^{(0)}$. 2 for $t = 0, 1, 2, \dots$ do 3 Update $\bar{\mu}$ and $\bar{\nu}$ according to $\begin{cases} \bar{\mu}^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\bar{\nu}^{(t)}]_a), \\ \bar{\nu}^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \bar{\mu}^{(t)}]_b). \end{cases}$ 4 Update μ and ν according to $\begin{cases} \mu^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\bar{\nu}^{(t+1)}]_a), \\ \nu^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \bar{\mu}^{(t+1)}]_b). \end{cases}$

- **Linear convergence of policies in KL divergence and entrywise log-ratios:**

$$\max \left\{ \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}), \frac{1}{2} \text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) \right\} \leq (1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}), \quad (9a)$$

$$\left\| \log \frac{\zeta^{(t)}}{\zeta_\tau^*} \right\|_\infty \leq 2(1 - \eta\tau)^t \left\| \log \frac{\zeta^{(0)}}{\zeta_\tau^*} \right\|_\infty + \frac{8 \|A\|_\infty}{\tau} (1 - \eta\tau)^{t/2} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})^{1/2}. \quad (9b)$$

- **Linear convergence of values in optimality and duality gaps:**

$$\text{OptGap}_\tau(\bar{\zeta}^{(t)}) \leq \eta^{-1} \cdot \frac{1}{1 - (\tau + \|A\|_\infty)\eta} \cdot \frac{(1 - \eta\tau)^t}{1 - (1 - \eta\tau)^t} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}), \quad (9c)$$

$$\text{DualGap}_\tau(\bar{\zeta}^{(t)}) \leq \left(\eta^{-1} + 2\tau^{-1} \|A\|_\infty^2 \right) (1 - \eta\tau)^{t-1} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}). \quad (9d)$$

Remark 1. Setting $\mu^{(0)}$ and $\nu^{(0)}$ to be uniform policies leads to a universal bound

$$\text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}) = \log |\mathcal{A}| + \log |\mathcal{B}| - \mathcal{H}(\mu_\tau^*) - \mathcal{H}(\nu_\tau^*) \leq \log |\mathcal{A}| + \log |\mathcal{B}|$$

regardless of $\zeta_\tau^* = (\mu_\tau^*, \nu_\tau^*)$.

Remark 2. Similar results continue to hold even when the two players use different regularization parameters $\tau_\mu, \tau_\nu > 0$ in (3), as long as the regularization parameter τ is replaced by $\max\{\tau_\mu, \tau_\nu\}$ in the upper bounds of the learning rate, and the contraction parameter is replaced by $1 - \min\{\tau_\mu, \tau_\nu\}\eta$.

Theorem 1 characterizes the convergence of the *last-iterates* $\zeta^{(t)}$ and $\bar{\zeta}^{(t)}$ of PU and OMWU as long as the learning rate lies within the specified ranges. While PU doubles the number of calls to the sampling oracle, it also allows roughly as large as twice the learning rate compared with OMWU (cf. (8)). Compared with the vast literature analyzing the average-iterate performance of variants of extragradient methods, our results contribute towards characterizing the last-iterate convergence of multiplicative update methods in the presence of entropy regularization and simplex constraints, which to the best of our knowledge, are the first of its kind. Several remarks are in order.

- **Linear convergence to QRE.** To achieve an ϵ -accurate estimate of the QRE in terms of the KL divergence, the bound (9a) tells that it is sufficient to take

$$\frac{1}{\eta\tau} \log \left(\frac{\log |\mathcal{A}| + \log |\mathcal{B}|}{\epsilon} \right)$$

iterations using either PU or OMWU. Notably, this iteration complexity does not depend on any hidden constants and only depends double logarithmically on the cardinality of action spaces, which is almost dimension-free. Maximizing the learning rate, the iteration complexity is bounded by $(1 + \|A\|_\infty/\tau) \log(1/\epsilon)$ (modulo log factors), which only depends on the ratio $\|A\|_\infty/\tau$.

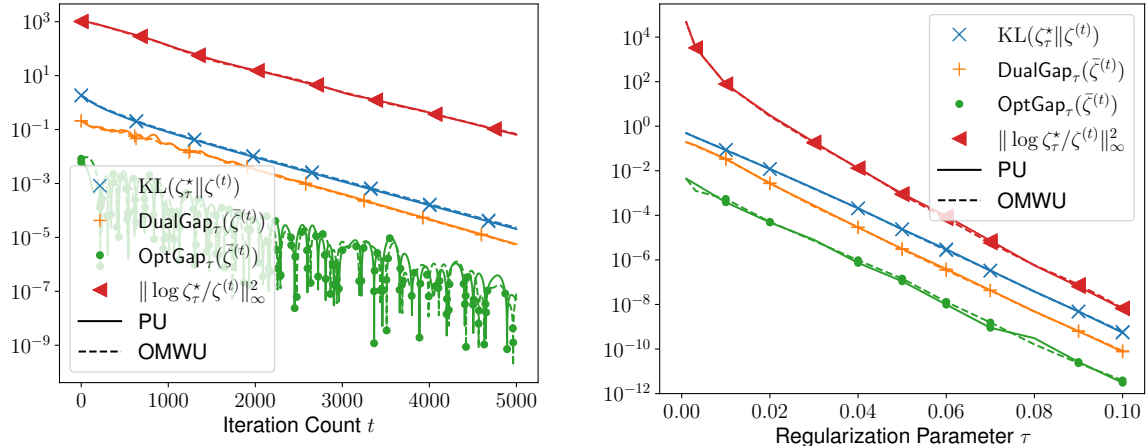


Figure 1: Performance illustration of the PU and OMWU methods for solving entropy-regularized matrix games with $|\mathcal{A}| = |\mathcal{B}| = 100$, where the entries of the payoff matrix A is generated independently from the uniform distribution on $[-1, 1]$. The learning rates are fixed as $\eta = 0.1$. The left panel plots various error metrics of convergence w.r.t. the iteration count with the entropy regularization parameter $\tau = 0.01$, while the right panel plots these error metrics at 1000-th iteration with different choices of τ .

- **Entrywise error of the policy log-ratios.** Both PU and OMWU enjoy strong entrywise guarantees in the sense we can guarantee the convergence of the ℓ_∞ norm of the log-ratios between the learned policy pair and the QRE at the same dimension-free linear rate (cf. (9b)), which suggests the policy pair converges in a somewhat uniform manner across the entire action space.
- **Linear convergence of optimality and duality gaps.** Our theorem also establishes the last-iterate convergence of the game values in terms of the optimality gap (cf. (9c)) and the duality gap (cf. (9d)) for both PU and OMWU. In particular, as will be seen, bounding the optimality gap of matrix games turns out to be the key enabler for generalizing our algorithms to Markov games, and bounding the duality gap allows to directly translate our results to finding a NE of unregularized matrix games.

Figure 1 illustrates the performance of the proposed PU and OMWU methods for solving randomly generated entropy-regularized matrix games. It is evident that both algorithms converge linearly, and achieve faster convergence rates when the regularization parameter increases.

Last-iterate convergence to approximate NE. The entropy-regularized matrix game can be thought as a smooth surrogate of the unregularized matrix game (1); in particular, it is possible to find an ϵ -NE by setting τ sufficiently small in (3). According to (Zhang et al., 2020, Definition 2.1), a policy pair $\zeta = (\mu, \nu)$ is an ϵ -NE if it satisfies

$$\text{DualGap}(\zeta) := \max_{\mu' \in \Delta(\mathcal{A})} f(\mu', \nu) - \min_{\nu' \in \Delta(\mathcal{B})} f(\mu, \nu') \leq \epsilon.$$

Observe that setting $\tau = \frac{\epsilon/4}{\log |\mathcal{A}| + \log |\mathcal{B}|}$ guarantees

$$|f_\tau(\mu, \nu) - f(\mu, \nu)| < \epsilon/4 \quad \text{for all } (\mu, \nu) \in \Delta(\mathcal{A}) \times \Delta(\mathcal{B})$$

in view of the boundedness of the Shannon entropy $\mathcal{H}(\cdot)$. Theorem 9 (cf. (9d)) also ensures that our proposed algorithms find an approximate QRE $\bar{\zeta}^{(T)}$ such that $\text{DualGap}_\tau(\bar{\zeta}^{(T)}) \leq \epsilon/2$ after taking $T = \tilde{O}\left(\frac{1}{\eta\epsilon}\right)$ iterations, which is no more than

$$\tilde{O}\left(\frac{\|A\|_\infty}{\epsilon}\right)$$

iterations with optimized learning rates. It follows immediately that

$$\text{DualGap}(\bar{\zeta}^{(T)}) \leq \text{DualGap}_\tau(\bar{\zeta}^{(T)}) + \max_{\mu', \nu'} \left| f_\tau(\mu', \bar{\nu}^{(T)}) - f_\tau(\bar{\mu}^{(T)}, \nu') - (f(\mu', \bar{\nu}^{(T)}) - f(\bar{\mu}^{(T)}, \nu')) \right| \leq \epsilon, \quad (10)$$

and therefore $\bar{\zeta}^{(T)}$ is an ϵ -NE. Intriguingly, unlike prior work (Daskalakis and Panageas, 2018a; Wei et al., 2021b) that analyzed the last-iterate convergence of OMWU in the unregularized setting ($\tau = 0$), our last-iterate convergence does not require the NE of (1) to be unique.

Rationality. Another attractive feature of the algorithms developed above is being *rational* (as introduced in Bowling and Veloso (2001)) in the sense that the algorithm returns the best-response policy of one player when the opponent takes any *fixed* stationary policy. More specially, in terms of matrix games, when player 2 sticks to a stationary policy ν , the update of player 1 reduces to

$$\mu^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\nu]_a). \quad (11)$$

In this case, Theorem 1 can be established in exactly the same fashion by restricting attention only to the updates of $\mu^{(t)}$.

3 Zero-sum Markov games with entropy regularization

Leveraging the success of PU and OMWU in solving the entropy-regularized matrix games, this section extends our current analysis to solve the zero-sum two-player Markov game, which is again formulated as finding the equilibrium of a saddle-point optimization problem. We start by introducing its basic setup, along with the entropy-regularized Markov game, which will be followed by the proposed policy extragradient method with its theoretical guarantees. The proofs for this section are collected in Appendix B.

3.1 Background and problem formulation

Zero-sum two-player Markov game. We consider a discounted Markov Game (MG) which is defined as $\mathcal{M} = \{\mathcal{S}, \mathcal{A}, \mathcal{B}, P, r, \gamma\}$, with discrete state space \mathcal{S} , action spaces of two players \mathcal{A} and \mathcal{B} , transition probability P , reward function $r : \mathcal{S} \times \mathcal{A} \times \mathcal{B} \rightarrow [0, 1]$ and discount factor $\gamma \in [0, 1)$. A policy $\mu : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ (resp. $\nu : \mathcal{S} \rightarrow \Delta(\mathcal{B})$) defines how player 1 (resp. player 2) reacts to a given state s , where the probability of taking action $a \in \mathcal{A}$ (resp. $b \in \mathcal{B}$) is $\mu(a|s)$ (resp. $\nu(b|s)$). The transition probability kernel $P : \mathcal{S} \times \mathcal{A} \times \mathcal{B} \rightarrow \Delta(\mathcal{S})$ defines the dynamics of the Markov game, where $P(s'|s, a, b)$ specifies the probability of transiting to state s' from state s when the players take actions a and b respectively. The state value of a given policy pair (μ, ν) is evaluated by the expected discounted cumulative reward:

$$V^{\mu, \nu}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t) \mid s_0 = s \right],$$

where the trajectory $(s_0, a_0, b_0, s_1, \dots)$ is generated by the MG \mathcal{M} under the policy pair (μ, ν) , starting from the state s_0 . Similarly, the Q-function captures the expected discounted cumulative reward with an initial state s and initial action pair (a, b) for a given policy pair (μ, ν) :

$$Q^{\mu, \nu}(s, a, b) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t) \mid s_0 = s, a_0 = a, b_0 = b \right].$$

In a zero-sum game, one player seeks to maximize the value function while the other player wants to minimize it. The minimax game value on state s is defined by

$$V^*(s) = \max_{\mu} \min_{\nu} V^{\mu, \nu}(s) = \min_{\nu} \max_{\mu} V^{\mu, \nu}(s).$$

Similarly, the minimax Q-function $Q^*(s, a, b)$ is defined by

$$Q^*(s, a, b) = r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} V^*(s'). \quad (12)$$

It is proved by Shapley (1953) that a pair of stationary policy (μ^*, ν^*) attaining the minimax value on state s attains the minimax value on all states as well (Filar and Vrieze, 2012), and is called the NE of the MG.

Entropy-regularized zero-sum two-player Markov game. Motivated by entropy regularization in Markov decision processes (MDP) (Geist et al., 2019; Cen et al., 2020), we consider an entropy-regularized variant of MG, where the value function is modified as

$$V_{\tau}^{\mu,\nu}(s) := \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t, b_t) - \tau \log \mu(a_t|s_t) + \tau \log \nu(b_t|s_t)) \mid s_0 = s \right], \quad (13)$$

where the quantity $\tau \geq 0$ denotes the regularization parameter, and the expectation is evaluated over the randomness of the transition kernel as well as the policies. The regularized Q-function $Q_{\tau}^{\mu,\nu}$ of a policy pair (μ, ν) is related to $V_{\tau}^{\mu,\nu}$ as

$$Q_{\tau}^{\mu,\nu}(s, a, b) = r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a, b)} [V_{\tau}^{\mu,\nu}(s')]. \quad (14)$$

We will call $V_{\tau}^{\mu,\nu}$ and $Q_{\tau}^{\mu,\nu}$ the *soft value function* and *soft Q-function*, respectively. A policy pair $(\mu_{\tau}^*, \nu_{\tau}^*)$ is said to be the quantal response equilibrium (QRE) of the entropy-regularized MG, if its value attains the minimax value of the entropy-regularized MG over all states $s \in \mathcal{S}$, i.e.

$$V_{\tau}^*(s) = \max_{\mu} \min_{\nu} V_{\tau}^{\mu,\nu}(s) = \min_{\nu} \max_{\mu} V_{\tau}^{\mu,\nu}(s) := V_{\tau}^{\mu_{\tau}^*, \nu_{\tau}^*}(s),$$

where V_{τ}^* is called the optimal minimax soft value function, and similarly $Q_{\tau}^* := Q_{\tau}^{\mu_{\tau}^*, \nu_{\tau}^*}$ is called the optimal minimax soft Q-function.

Goal. Our goal is to find the QRE of the entropy-regularized MG in a decentralized manner where the players only observe its own reward without accessing the opponent's actions.

Remark 3. For any policy pair (μ, ν) , it is straightforward to show that

$$\|V_{\tau}^{\mu,\nu} - V^{\mu,\nu}\|_{\infty} \leq \tau(\log |\mathcal{A}| + \log |\mathcal{B}|).$$

Hence, similar to the case of matrix games, setting the regularization parameter sufficiently small τ , solving the entropy-regularized MG also allows us to find an approximate NE of the unregularized MG. We omit the details for conciseness.

3.2 From value iteration to policy extragradient methods

Entropy-regularized value iteration. It is known that classical dynamic programming approaches such as value iteration can be extended to solve MG (Perolat et al., 2015), where each iteration amounts to solving a series of matrix games for each state. Similar to the single-agent case (Cen et al., 2020), we can extend these approaches to solve the entropy-regularized MG. Setting the stage, let us introduce the per-state Q-value matrix $Q(s) := Q(s, \cdot, \cdot) \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{B}|}$ for every $s \in \mathcal{S}$, where the element indexed by the action pair (a, b) is $Q(s, a, b)$. Similarly, we define the per-state policies $\mu(s) := \mu(\cdot|s) \in \Delta(\mathcal{A})$ and $\nu(s) := \nu(\cdot|s) \in \Delta(\mathcal{B})$ for both players.

In parallel to the original Bellman operator, we denote the *soft Bellman operator* \mathcal{T}_{τ} as

$$\mathcal{T}_{\tau}(Q)(s, a, b) := r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a, b)} \left[\max_{\mu(s') \in \Delta(\mathcal{A})} \min_{\nu(s') \in \Delta(\mathcal{B})} f_{\tau}(Q(s'); \mu(s'), \nu(s')) \right], \quad (15)$$

where for each per-state Q-value matrix $Q(s)$, we introduce an entropy-regularized matrix game in the form of

$$\max_{\mu \in \Delta(\mathcal{A})} \min_{\nu \in \Delta(\mathcal{B})} f_{\tau}(Q(s); \mu(s), \nu(s)) := \mu(s)^{\top} Q(s) \nu(s) - \tau \mathcal{H}(\mu(s)) + \tau \mathcal{H}(\nu(s)).$$

The entropy-regularized value iteration then proceeds as

$$Q^{(t+1)} = \mathcal{T}_{\tau}(Q^{(t)}), \quad (16)$$

where $Q^{(0)}$ is an initialization. By definition, the optimal minimax soft Q-function obeys $\mathcal{T}_{\tau}(Q_{\tau}^*) = Q_{\tau}^*$ and therefore corresponds to the fix point of the soft Bellman operator. Given the above entropy-regularized value iteration, the following lemma states its iterates contract linearly to the optimal minimax soft Q-function at a rate of the discount factor γ .

Algorithm 3: Policy Extragradient Method for Entropy-regularized Markov Game

- 1 **initialization:** $Q^{(0)} = 0$.
 - 2 **for** $t = 0, 1, 2, \dots, T_{\text{main}}$ **do**
 - 3 Let $Q^{(t)}$ denote

$$Q^{(t)}(s, a, b) = r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a, b)} V^{(t)}(s'). \quad (18)$$
 - 4 Invoke PU (Algorithm 1) or OMWU (Algorithm 2) for T_{sub} iterations to solve the following entropy-regularized matrix game for every state s , where the initialization is set as uniform distributions:

$$\max_{\mu(s) \in \Delta(\mathcal{A})} \min_{\nu(s) \in \Delta(\mathcal{B})} f_{\tau}(Q^{(t)}(s); \mu(s), \nu(s)).$$
 - 5 Return the last iterate $\bar{\mu}^{(t, T_{\text{sub}})}(s), \bar{\nu}^{(t, T_{\text{sub}})}(s)$.
 - 6 Set $V^{(t+1)}(s) = f_{\tau}(Q^{(t)}(s); \bar{\mu}^{(t, T_{\text{sub}})}(s), \bar{\nu}^{(t, T_{\text{sub}})}(s))$.
-

Proposition 2. *The entropy-regularized value iteration (16) converges at a linear rate, i.e.*

$$\|Q^{(t)} - Q_{\tau}^*\|_{\infty} \leq \gamma^t \|Q^{(0)} - Q_{\tau}^*\|_{\infty}. \quad (17)$$

Approximate value iteration via policy extragradient methods. Proposition 2 suggests that the optimal minimax soft Q-function of the entropy-regularized MG can be found by solving a series of entropy-regularized matrix games induced by $\{Q^{(t)}\}_{t \geq 0}$ in (16), a task that can be accomplished by adopting the fast extragradient methods developed earlier. To proceed, we first define the following sampling oracle, which makes it rigorous that the proposed algorithm does not require access to the Q-function of the entire MG, but only its own single-agent Q-function when playing against the opponent's policy.

Definition 2 (Sampling oracle for Markov games). *Given any policy pair $\mu(s), \nu(s)$ and Q-value matrix $Q(s)$ for any $s \in \mathcal{S}$, the sampling oracle returns*

$$[Q(s)\nu(s)]_a = \mathbb{E}_{b \sim \nu(s)} [Q(s, a, b)], \quad \text{and} \quad [Q(s)^{\top} \mu(s)]_b = \mathbb{E}_{a \sim \mu(s)} [Q(s, a, b)]$$

for any $a \in \mathcal{A}$ and $b \in \mathcal{B}$.

Algorithm 3 describes the proposed policy extragradient method. Encouragingly, by judiciously setting the number of iterations in both the outer loop (for updating the Q-value matrices) and the inner loop (for updating the QRE of the corresponding Q-value matrix), we are guaranteed to find the QRE of the entropy-regularized MG in a small number of iterations, as dictated by the following theorem.

Theorem 2. *Assume $|\mathcal{A}| \geq |\mathcal{B}|$ and $\tau \leq 1$. Setting $\eta = \frac{1-\gamma}{2(1+\tau(\log|\mathcal{A}|+1-\gamma))}$, the total iterations (namely, the product $T_{\text{main}} \cdot T_{\text{sub}}$) required for Algorithm 3 to achieve $\|Q^{(T_{\text{main}})} - Q_{\tau}^*\|_{\infty} \leq \epsilon$ is at most*

$$O\left(\frac{(\log|\mathcal{A}| + 1/\tau)}{(1-\gamma)^2} \left(\log \frac{\log|\mathcal{A}|}{(1-\gamma)\epsilon}\right)^2\right).$$

Theorem 2 ensures that within $\tilde{O}\left(\frac{1}{\tau(1-\gamma)^2} \log^2\left(\frac{1}{\epsilon}\right)\right)$ iterations, Algorithm 3 finds a pair of policies whose value is close to the optimal minimax soft Q-function Q_{τ}^* in an entrywise manner to a prescribed accuracy ϵ . Remarkably, the iteration complexity is independent of the dimensions of the state space and the action space (up to log factors).

Figure 2 illustrates the performance of Algorithm 3 for solving a randomly generated entropy-regularized Markov game with $|\mathcal{A}| = |\mathcal{B}| = 20$, $|\mathcal{S}| = 100$ and $\gamma = 0.99$ with varying choices of T_{main} , T_{sub} and τ . Here, the transition probability kernel and the reward function are generated as follows. For each state-action pair (s, a, b) , we randomly select 10 states to form a set $\mathcal{S}_{s, a, b}$, and set $P(s'|s, a, b) \propto U_{s, a, b, s'}$ if $s' \in \mathcal{S}_{s, a, b}$, and 0 otherwise, where $\{U_{s, a, b, s'} \sim U[0, 1]\}$ are drawn independently from the uniform distribution over $[0, 1]$. The reward function is generated by $r(s, a, b) \sim U_{s, a, b} \cdot U_s$, where $U_{s, a, b}$ and U_s are drawn independently from the uniform distribution over $[0, 1]$. It is seen that the convergence speed of the ℓ_{∞} error on $\|Q^{(T_{\text{main}})} - Q_{\tau}^*\|_{\infty}$ improves as we increase the regularization parameter τ .

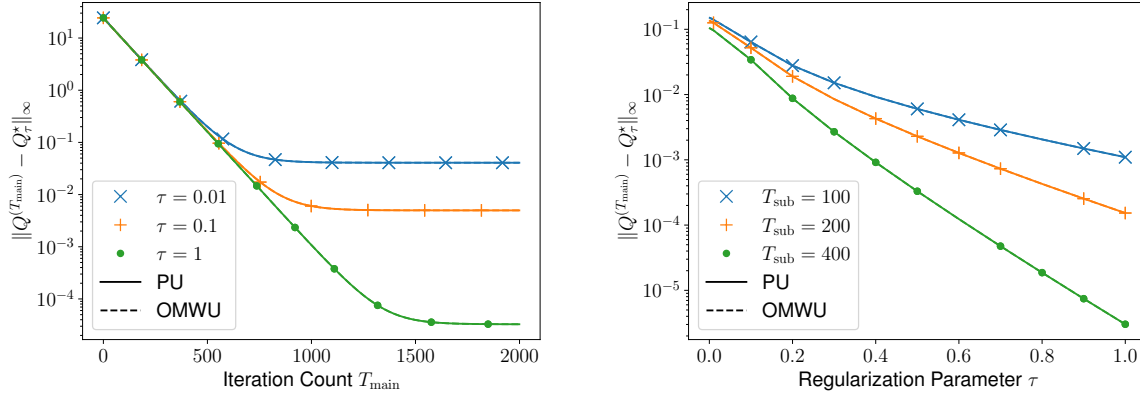


Figure 2: Performance illustration of Algorithm 3 for solving a random generated entropy-regularized Markov game with $|\mathcal{A}| = |\mathcal{B}| = 20$, $|\mathcal{S}| = 100$ and $\gamma = 0.99$. The learning rates of both players are fixed as $\eta = 0.005$. The left panel plots $\|Q^{(T_{\text{main}})} - Q_{\tau}^*\|_{\infty}$ w.r.t. the iteration count T_{main} when $T_{\text{sub}} = 400$ under various entropy regularization parameters, while the right panel plots $\|Q^{(T_{\text{main}})} - Q_{\tau}^*\|_{\infty}$ w.r.t. the regularization parameter τ when $T_{\text{main}} = 2000$ with different choices of T_{sub} .

4 Conclusions

This paper develops provably efficient policy extragradient methods (PU and OMWU) for entropy-regularized matrix games and Markov games, whose last iterates are guaranteed to converge linearly to the quantal response equilibrium at a linear rate. Encouragingly, the rate of convergence is independent of the dimension of the problem, i.e. the sizes of the space space and the action space. In addition, the last iterates of the proposed algorithms can also be used to locate Nash equilibria for the unregularized competitive games without assuming the uniqueness of the Nash equilibria by judiciously tuning the amount of regularization. This work opens up interesting opportunities for further investigations of policy extragradient methods for solving competitive games. For example, can we develop a two-time-scale policy extragradient algorithms for Markov games where the Q-function is updated simultaneously with the policy but potentially at a different time scale, using samples, such as in an actor-critic algorithm (Konda and Tsitsiklis, 2000)? Can we generalize the proposed algorithms to handle more general regularization terms, similar to what has been accomplished in the single-agent setting (Lan, 2021; Zhan et al., 2021)? We leave the answers to future work.

Acknowledgments

S. Cen and Y. Chi are supported in part by the grants ONR N00014-18-1-2142 and N00014-19-1-2404, ARO W911NF-18-1-0303, NSF CCF-1901199, CCF-2007911 and CCF-2106778. Y. Wei is supported in part by the NSF grants CCF-2007911, DMS-2015447 and CCF-2106778.

A Analysis for entropy-regularized matrix games

Before embarking on the main proof, it is useful to first consider the update rule (6) that underlies both PU and OMWU, which is reproduced below for convenience:

$$\begin{cases} \mu^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[Az_2]_a), & \text{for all } a \in \mathcal{A}, \\ \nu^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top z_1]_b), & \text{for all } b \in \mathcal{B}, \end{cases} \quad (19)$$

where $z_1 \in \Delta(\mathcal{A})$ and $z_2 \in \Delta(\mathcal{B})$. These updates satisfy the following property, whose proof is provided in Appendix C.1.

Lemma 1. Denote $\zeta^{(t)} = (\mu^{(t)}, \nu^{(t)})$ and $\zeta(z) = (z_1, z_2)$. The update rule (19) satisfies:

$$\langle \log \mu^{(t+1)} - (1 - \eta\tau) \log \mu^{(t)} - \eta\tau \log \mu_\tau^*, z_1 - \mu_\tau^* \rangle = \eta(\mu_\tau^* - z_1)^\top A(\nu_\tau^* - z_2), \quad (20a)$$

$$\langle \log \nu^{(t+1)} - (1 - \eta\tau) \log \nu^{(t)} - \eta\tau \log \nu_\tau^*, z_2 - \nu_\tau^* \rangle = -\eta(\nu_\tau^* - z_1)^\top A(\nu_\tau^* - z_2), \quad (20b)$$

and

$$\langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \zeta(z) - \zeta_\tau^* \rangle = 0. \quad (21)$$

As we shall see, the above lemma plays a crucial role in establishing the claimed convergence results. The next lemma gives some basic decompositions related to the game values that are helpful.

Lemma 2. For every $(\mu, \nu) \in \Delta(\mathcal{A}) \times \Delta(\mathcal{B})$, the following relations hold

$$f_\tau(\mu_\tau^*, \nu) - f_\tau(\mu, \nu_\tau^*) = \tau \text{KL}(\zeta \parallel \zeta_\tau^*), \quad (22a)$$

$$f_\tau(\mu, \nu) - f_\tau(\mu_\tau^*, \nu_\tau^*) = (\mu_\tau^* - \mu)^\top A(\nu_\tau^* - \nu) + \tau \text{KL}(\nu \parallel \nu_\tau^*) - \tau \text{KL}(\mu \parallel \mu_\tau^*). \quad (22b)$$

In addition, we also make record of the following elementary lemma that is used frequently.

Lemma 3. For any $\mu_1, \mu_2 \in \Delta(\mathcal{A})$ satisfying

$$\mu_1(a) \propto \exp(x_1(a)) \quad \text{and} \quad \mu_2(a) \propto \exp(x_2(a))$$

for some $x_1, x_2 \in \mathbb{R}^{|\mathcal{A}|}$, we have

$$\|\log \mu_1 - \log \mu_2\|_\infty \leq 2 \|x_1 - x_2\|_\infty.$$

A.1 Proof of Proposition 1

Setting $\zeta(z) = \zeta^{(t+1)}$ in Lemma 1, we have

$$\langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \zeta^{(t+1)} - \zeta_\tau^* \rangle = 0. \quad (23)$$

By the definition of the KL divergence, one has

$$\begin{aligned} & - \langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \zeta_\tau^* \rangle \\ &= -(1 - \eta\tau) \langle \log \zeta_\tau^* - \log \zeta^{(t)}, \zeta_\tau^* \rangle + \langle \log \zeta_\tau^* - \log \zeta^{(t+1)}, \zeta_\tau^* \rangle \\ &= -(1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) + \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}), \end{aligned} \quad (24)$$

and similarly,

$$\begin{aligned} & \langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \zeta^{(t+1)} \rangle \\ &= (1 - \eta\tau) \langle \log \zeta^{(t+1)} - \log \zeta^{(t)}, \zeta^{(t+1)} \rangle + \eta\tau \langle \log \zeta^{(t+1)} - \log \zeta_\tau^*, \zeta^{(t+1)} \rangle \\ &= (1 - \eta\tau) \text{KL}(\zeta^{(t+1)} \parallel \zeta^{(t)}) + \eta\tau \text{KL}(\zeta^{(t+1)} \parallel \zeta_\tau^*). \end{aligned}$$

Combining the above two equalities with (23), we arrive at

$$\text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) + \eta\tau \text{KL}(\zeta^{(t+1)} \parallel \zeta_\tau^*) + (1 - \eta\tau) \text{KL}(\zeta^{(t+1)} \parallel \zeta^{(t)}) = (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}). \quad (25)$$

This immediately leads to $\text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) \leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)})$ by the nonnegativity of the KL divergence, as long as $1 - \eta\tau \geq 0$. Therefore

$$\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) \leq (1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}) \quad \text{for all } t \geq 0.$$

A.2 Proof of Theorem 1

A.2.1 Proof of policy convergence in KL divergence (9a)

First noticing that both PU and OMWU share the same update rule for $\mu^{(t+1)}$ and $\nu^{(t+1)}$, which takes the form

$$\begin{cases} \mu^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\bar{\nu}^{(t+1)}]_a), \\ \nu^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \bar{\mu}^{(t+1)}]_b). \end{cases}$$

Regarding this sequence, Lemma 1 (cf. (21)) gives

$$\langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \bar{\zeta}^{(t+1)} - \zeta_\tau^* \rangle = 0. \quad (26)$$

With the optimism that $\bar{\zeta}^{(t+1)}$ approximates $\zeta^{(t+1)}$ well, we can expect similar convergence guarantees to that of the implicit updates established in Proposition 1. Following the same argument as (24), we have

$$- \langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \zeta_\tau^* \rangle = -(1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) + \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}). \quad (27)$$

On the other hand, it is easily seen that

$$\begin{aligned} & \langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \bar{\zeta}^{(t+1)} \rangle \\ &= \langle \log \bar{\zeta}^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \bar{\zeta}^{(t+1)} \rangle + \langle \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)}, \zeta^{(t+1)} \rangle \\ &\quad - \langle \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)}, \bar{\zeta}^{(t+1)} - \zeta^{(t+1)} \rangle \\ &= (1 - \eta\tau) \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) + \eta\tau \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) + \text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) \\ &\quad - \langle \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)}, \bar{\zeta}^{(t+1)} - \zeta^{(t+1)} \rangle. \end{aligned} \quad (28)$$

Combining inequalities (27), (28) with (26), we are left with the following relation pertaining to bounding $\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)})$:

$$\begin{aligned} (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) &= (1 - \eta\tau) \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) + \eta\tau \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) + \text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) \\ &\quad - \langle \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)}, \bar{\zeta}^{(t+1)} - \zeta^{(t+1)} \rangle + \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}). \end{aligned} \quad (29)$$

In addition, to bound $\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)})$, we will resort to the following three-point equality, which reads

$$\begin{aligned} \text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) &= \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) - \langle \zeta_\tau^*, \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)} \rangle \\ &= \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) - \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t+1)}) - \langle \zeta_\tau^* - \bar{\zeta}^{(t+1)}, \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)} \rangle, \end{aligned} \quad (30)$$

which can be checked directly using the definition of the KL divergence.

To proceed, we need to control $\langle \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)}, \bar{\zeta}^{(t+1)} - \zeta^{(t+1)} \rangle$ on the right-hand side of inequality (29), and $\langle \zeta_\tau^* - \bar{\zeta}^{(t+1)}, \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)} \rangle$ on the right-hand side of inequality (30), for which we continue the proofs for PU and OMWU separately as follows.

Bounding $\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)})$ for PU. Following the update rule of $\bar{\zeta}^{(t+1)} = (\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})$ in PU, we have

$$\log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)} = \eta A(\nu^{(t)} - \bar{\nu}^{(t+1)}) + c \cdot \mathbf{1} \quad (31)$$

for some normalization constant c . With this relation in place, one has

$$\begin{aligned} \langle \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)}, \bar{\mu}^{(t+1)} - \mu^{(t+1)} \rangle &= \eta (\bar{\mu}^{(t+1)} - \mu^{(t+1)})^\top A(\nu^{(t)} - \bar{\nu}^{(t+1)}) \\ &\leq \eta \|A\|_\infty \left\| \bar{\mu}^{(t+1)} - \mu^{(t+1)} \right\|_1 \left\| \bar{\nu}^{(t+1)} - \nu^{(t)} \right\|_1. \end{aligned}$$

Combined with Pinsker's inequality, it is therefore clear that

$$\langle \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)}, \bar{\mu}^{(t+1)} - \mu^{(t+1)} \rangle \leq \frac{1}{2} \eta \|A\|_\infty \left(\left\| \bar{\mu}^{(t+1)} - \mu^{(t+1)} \right\|_1^2 + \left\| \bar{\nu}^{(t+1)} - \nu^{(t)} \right\|_1^2 \right)$$

$$\leq \eta \|A\|_\infty \left(\text{KL}(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) + \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) \right). \quad (32)$$

Analogously, one can achieve the same bound regarding the quantity $\langle \log \bar{\nu}^{(t+1)} - \log \nu^{(t+1)}, \bar{\nu}^{(t+1)} - \nu^{(t+1)} \rangle$. Summing up these two inequalities, we end up with

$$\langle \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)}, \bar{\zeta}^{(t+1)} - \zeta^{(t+1)} \rangle \leq \eta \|A\|_\infty \left(\text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) \right).$$

Plugging the above inequality into inequality (29) leads to

$$\begin{aligned} \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) &\leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) - (1 - \eta\tau - \eta \|A\|_\infty) \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) - \eta\tau \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\ &\quad - (1 - \eta \|A\|_\infty) \text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}). \end{aligned} \quad (33)$$

Therefore, as long as the learning rate η satisfies $\eta \leq \frac{1}{\tau + \|A\|_\infty}$, we are ensured that

$$\text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) \leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}),$$

which further implies inequality (9a) when applied recursively.

Bounding $\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)})$ for PU. By similar tricks of arriving at (32), we have

$$\begin{aligned} -\langle \mu_\tau^* - \bar{\mu}^{(t+1)}, \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)} \rangle &= -\eta (\mu_\tau^* - \bar{\mu}^{(t+1)})^\top A (\nu^{(t)} - \bar{\nu}^{(t+1)}) \\ &\leq \frac{1}{2} \eta \|A\|_\infty \left(\left\| \mu_\tau^* - \bar{\mu}^{(t+1)} \right\|_1^2 + \left\| \nu^{(t)} - \bar{\nu}^{(t+1)} \right\|_1^2 \right) \\ &\leq \eta \|A\|_\infty \left(\text{KL}(\mu_\tau^* \parallel \bar{\mu}^{(t+1)}) + \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) \right), \end{aligned}$$

following from (31) and Pinsker's inequality. A similar inequality for $-\langle \nu_\tau^* - \bar{\nu}^{(t+1)}, \log \bar{\nu}^{(t+1)} - \log \nu^{(t+1)} \rangle$ can be obtained by symmetry, and summing together the two leads to

$$-\langle \zeta_\tau^* - \bar{\zeta}^{(t+1)}, \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)} \rangle \leq \eta \|A\|_\infty \left(\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) + \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) \right).$$

Plugging the above inequality into (30) and rearranging terms, we reach at

$$(1 - \eta \|A\|_\infty) \text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) \leq \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) + \eta \|A\|_\infty \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}).$$

Along with (33), we have

$$\begin{aligned} (1 - \eta \|A\|_\infty) \text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) &\leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) - (1 - \eta\tau - 2\eta \|A\|_\infty) \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) \\ &\quad - \eta\tau \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) - (1 - \eta \|A\|_\infty) \text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}). \end{aligned} \quad (34)$$

Therefore, with $\eta \leq 1/(\tau + 2\|A\|_\infty)$ we have

$$\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) \leq 2 \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) \leq 2(1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}).$$

Bounding $\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)})$ for OMWU. Following the update rule of $\bar{\zeta}^{(t+1)} = (\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})$ for OMWU, we have

$$\begin{aligned} \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)} &= \eta A (\bar{\nu}^{(t)} - \bar{\nu}^{(t+1)}) + c \cdot \mathbf{1} \\ &= \eta A (\bar{\nu}^{(t)} - \nu^{(t)}) + \eta A (\nu^{(t)} - \bar{\nu}^{(t+1)}) + c \cdot \mathbf{1}, \end{aligned} \quad (35)$$

where c is some normalization constant. Similar to the proof of relation (32), it can be easily demonstrated that

$$\langle \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)}, \bar{\mu}^{(t+1)} - \mu^{(t+1)} \rangle$$

$$\begin{aligned}
&= \eta(\bar{\mu}^{(t+1)} - \mu^{(t+1)})^\top A(\bar{\nu}^{(t)} - \nu^{(t)}) + \eta(\bar{\mu}^{(t+1)} - \mu^{(t+1)})^\top A(\nu^{(t)} - \bar{\nu}^{(t+1)}) \\
&\leq \eta \|A\|_\infty \left(\text{KL}(\nu^{(t)} \parallel \bar{\nu}^{(t)}) + \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) + 2\text{KL}(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) \right). \tag{36}
\end{aligned}$$

By symmetry, we can also establish a similar inequality for $\langle \log \bar{\nu}^{(t+1)} - \log \nu^{(t+1)}, \bar{\nu}^{(t+1)} - \nu^{(t+1)} \rangle$, which in turns yields

$$\begin{aligned}
&\langle \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)}, \bar{\zeta}^{(t+1)} - \zeta^{(t+1)} \rangle \\
&\leq \eta \|A\|_\infty \left(\text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) + \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) + 2\text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) \right).
\end{aligned}$$

Plugging the above inequality into equation (29) and re-organizing terms, we arrive at

$$\begin{aligned}
\text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) &\leq (1 - \eta\tau)\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) - (1 - \eta\tau - \eta \|A\|_\infty)\text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) - \eta\tau\text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\
&\quad - (1 - 2\eta \|A\|_\infty)\text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + \eta \|A\|_\infty \text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}). \tag{37}
\end{aligned}$$

With the choice of the learning rate $\eta \leq \min\{\frac{1}{2\|A\|_\infty + 2\tau}, \frac{1}{4\|A\|_\infty}\}$, it obeys

$$(1 - \eta\tau)(1 - 2\eta \|A\|_\infty) \geq \eta \|A\|_\infty.$$

Combining the above inequality with (37) gives

$$\begin{aligned}
&\text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) + (1 - 2\eta \|A\|_\infty)\text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) \\
&\leq (1 - \eta\tau)\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) + \eta \|A\|_\infty \text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) - \eta\tau\text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\
&\leq (1 - \eta\tau) \left[\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) + (1 - 2\eta \|A\|_\infty)\text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) \right] - \eta\tau\text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*).
\end{aligned}$$

For conciseness, let us introduce the shorthand notation

$$L^{(t)} := \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) + (1 - 2\eta \|A\|_\infty)\text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}). \tag{38}$$

As a result, the above inequality can be restated as

$$L^{(t+1)} \leq (1 - \eta\tau)L^{(t)} - \eta\tau\text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*). \tag{39}$$

Since we initialize OMWU with $\bar{\zeta}^{(0)} = \zeta^{(0)}$, therefore $L^{(0)} = \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})$, which in turn gives

$$\text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) \leq L^{(t)} \leq (1 - \eta\tau)^t L^{(0)} = (1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}).$$

We complete the proof of inequality (9a) for OMWU.

Bounding $\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)})$ for OMWU. By similar tricks of arriving at (36), we have

$$\begin{aligned}
-\langle \mu_\tau^* - \bar{\mu}^{(t+1)}, \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)} \rangle &= \eta(\bar{\mu}^{(t+1)} - \mu_\tau^*)^\top A(\bar{\nu}^{(t)} - \nu^{(t)}) + \eta(\bar{\mu}^{(t+1)} - \mu_\tau^*)^\top A(\nu^{(t)} - \bar{\nu}^{(t+1)}) \\
&\leq \eta \|A\|_\infty \left(\text{KL}(\nu^{(t)} \parallel \bar{\nu}^{(t)}) + \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) + 2\text{KL}(\mu_\tau^* \parallel \bar{\mu}^{(t+1)}) \right),
\end{aligned}$$

where the first line follows from (35). A similar inequality also holds for $-\langle \nu_\tau^* - \bar{\nu}^{(t+1)}, \log \bar{\nu}^{(t+1)} - \log \nu^{(t+1)} \rangle$. Summing the two inequalities leads to

$$-\langle \zeta_\tau^* - \bar{\zeta}^{(t+1)}, \log \bar{\zeta}^{(t+1)} - \log \zeta^{(t+1)} \rangle \leq \eta \|A\|_\infty \left(\text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) + \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) + 2\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) \right).$$

Plugging the above inequality into (30) and rearranging terms, we reach at

$$(1 - 2\eta \|A\|_\infty)\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) \leq \text{KL}(\zeta_\tau^* \parallel \zeta^{(t+1)}) + \eta \|A\|_\infty \left(\text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) + \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) \right).$$

Along with (37), we have

$$\begin{aligned}
(1 - 2\eta \|A\|_\infty) \text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) &\leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) - (1 - \eta\tau - 2\eta \|A\|_\infty) \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) - \eta\tau \text{KL}(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\
&\quad - (1 - 2\eta \|A\|_\infty) \text{KL}(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + 2\eta \|A\|_\infty \text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) \\
&\leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) + 2\eta \|A\|_\infty \text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) \\
&\leq \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}) + (1 - 2\eta \|A\|_\infty) \text{KL}(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) =: L^{(t)},
\end{aligned}$$

where we recall the shorthand notation $L^{(t)}$ in (38). As the learning rate of OMWU satisfies $0 < \eta < \min\left\{\frac{1}{2\|A\|_\infty + 2\tau}, \frac{1}{4\|A\|_\infty}\right\}$, it is clear that

$$\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) \leq 2L^{(t)} \stackrel{(i)}{\leq} 2(1 - \eta\tau)^t L^{(0)} \leq 2(1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}),$$

where (i) follows from the recursive relation $L^{(t+1)} \leq (1 - \eta\tau)L^{(t)}$ shown in inequality (39).

A.2.2 Proof of entrywise convergence of policy log-ratios (9b)

To facilitate the proof, we introduce an auxiliary sequence $\{\xi^{(t)} \in \mathbb{R}^{|\mathcal{A}|}\}$ constructed recursively by

$$\xi^{(0)}(a) = \|\exp(A\nu_\tau^*/\tau)\|_1 \cdot \mu^{(0)}(a), \quad (40a)$$

$$\xi^{(t+1)}(a) = \xi^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\bar{\nu}^{(t+1)}]_a), \quad \forall a \in \mathcal{A}, t \geq 0. \quad (40b)$$

It is easily seen that $\mu^{(t)}(a) \propto \xi^{(t)}(a) = \exp(\log \xi^{(t)}(a))$ for $t \geq 0$. Noticing that $\mu_\tau^* \propto \exp(A\nu_\tau^*)$, one has

$$\left\| \log \mu^{(t+1)} - \log \mu_\tau^* \right\|_\infty \leq 2 \left\| \log \xi^{(t+1)} - A\nu_\tau^*/\tau \right\|_\infty, \quad (41)$$

where we make use of Lemma 3.

Therefore it suffices for us to control the term $\left\| \log \xi^{(t+1)} - A\nu_\tau^*/\tau \right\|_\infty$ on the right-hand side of inequality (41). Taking logarithm on both sides of (40b) yields

$$\begin{aligned}
\log \xi^{(t+1)} - A\nu_\tau^*/\tau &= (1 - \eta\tau) \log \xi^{(t)} + \eta A\bar{\nu}^{(t+1)} - A\nu_\tau^*/\tau \\
&= (1 - \eta\tau) \left(\log \xi^{(t)} - A\nu_\tau^*/\tau \right) + \eta A(\bar{\nu}^{(t+1)} - \nu_\tau^*),
\end{aligned}$$

which, when combined with Pinsker's inequality, implies

$$\begin{aligned}
\left\| \log \xi^{(t+1)} - A\nu_\tau^*/\tau \right\|_\infty &\leq (1 - \eta\tau) \left\| \log \xi^{(t)} - A\nu_\tau^*/\tau \right\|_\infty + \eta \|A\|_\infty \left\| \bar{\nu}^{(t+1)} - \nu_\tau^* \right\|_1 \\
&\leq (1 - \eta\tau) \left\| \log \xi^{(t)} - A\nu_\tau^*/\tau \right\|_\infty + \eta \|A\|_\infty \left[2\text{KL}(\nu_\tau^* \parallel \bar{\nu}^{(t+1)}) \right]^{1/2} \\
&\leq (1 - \eta\tau) \left\| \log \xi^{(t)} - A\nu_\tau^*/\tau \right\|_\infty + \eta \|A\|_\infty \left[2\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) \right]^{1/2}. \quad (42)
\end{aligned}$$

Plugging the bound of $\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)})$ from relation (9a) into (42) and invoking the inequality recursively leads to

$$\begin{aligned}
&\left\| \log \xi^{(t+1)} - A\nu_\tau^*/\tau \right\|_\infty \\
&\leq (1 - \eta\tau)^{t+1} \left\| \log \xi^{(0)} - A\nu_\tau^*/\tau \right\|_\infty + 2\eta \|A\|_\infty \sum_{s=1}^{t+1} (1 - \eta\tau)^{t+1-s/2} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})^{1/2} \\
&\leq (1 - \eta\tau)^{t+1} \left\| \log \xi^{(0)} - A\nu_\tau^*/\tau \right\|_\infty + 2\eta \|A\|_\infty (1 - \eta\tau)^{(t+1)/2} \frac{1}{1 - (1 - \eta\tau)^{1/2}} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})^{1/2} \\
&\leq (1 - \eta\tau)^{t+1} \left\| \log \xi^{(0)} - A\nu_\tau^*/\tau \right\|_\infty + 4\tau^{-1} \|A\|_\infty (1 - \eta\tau)^{(t+1)/2} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})^{1/2},
\end{aligned}$$

where the last line results from the fact that $(1 - \eta\tau)^{1/2} \leq 1 - \eta\tau/2$. Combining pieces together, we end up with

$$\begin{aligned} \left\| \log \mu^{(t+1)} - \log \mu_\tau^* \right\|_\infty &\leq 2 \left\| \log \xi^{(t+1)} - A\nu_\tau^*/\tau \right\|_\infty \\ &\leq 2(1 - \eta\tau)^{t+1} \left\| \log \xi^{(0)} - A\nu_\tau^*/\tau \right\|_\infty + 8\tau^{-1} \|A\|_\infty (1 - \eta\tau)^{(t+1)/2} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})^{1/2} \\ &\leq 2(1 - \eta\tau)^{t+1} \left\| \log \mu^{(0)} - \log \mu_\tau^* \right\|_\infty + 8\tau^{-1} \|A\|_\infty (1 - \eta\tau)^{(t+1)/2} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})^{1/2}. \end{aligned}$$

Similarly, one can establish the corresponding inequality for $\left\| \log \nu^{(t+1)} - \log \nu_\tau^* \right\|_\infty$, therefore completing the proof of inequality (9b).

A.2.3 Proof of convergence of optimality gap (9c)

To streamline our discussions, we only provide the proof of inequality (9c) concerning upper bounding $f_\tau(\bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_\tau(\mu_\tau^*, \nu_\tau^*)$ without taking the absolute value; the other direction of the inequality can be established in the similar manner and hence is omitted.

We first make note of an important relation that holds both for PU and OMWU. Consider the update rule of $(\mu^{(t+1)}, \nu^{(t+1)})$, which is the same in PU and OMWU. Lemma 1 inequality (20a) gives

$$\langle \log \mu^{(t+1)} - (1 - \eta\tau) \log \mu^{(t)} - \eta\tau \log \mu_\tau^*, \bar{\mu}^{(t+1)} - \mu_\tau^* \rangle = \eta(\mu_\tau^* - \bar{\mu}^{(t+1)})^\top A(\nu_\tau^* - \bar{\nu}^{(t+1)}). \quad (43)$$

Similar to what we have done in the proof of (9a) (cf. (29)), based on the above relation, we can therefore rearrange terms and conclude that

$$\begin{aligned} &\eta \left(\tau \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu_\tau^*) - (\mu_\tau^* - \bar{\mu}^{(t+1)})^\top A(\nu_\tau^* - \bar{\nu}^{(t+1)}) \right) \\ &= (1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) - (1 - \eta\tau) \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) - \text{KL}(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) \\ &\quad + \langle \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)}, \bar{\mu}^{(t+1)} - \mu^{(t+1)} \rangle - \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}). \end{aligned}$$

In conjunction with Lemma 2 (cf. (22b)), we can further derive

$$\begin{aligned} \eta(f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})) &\leq \eta \left(\tau \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu_\tau^*) - (\mu_\tau^* - \bar{\mu}^{(t+1)})^\top A(\nu_\tau^* - \bar{\nu}^{(t+1)}) \right) \\ &= (1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) - (1 - \eta\tau) \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) - \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}) \\ &\quad - \text{KL}(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) + \langle \log \bar{\mu}^{(t+1)} - \log \mu^{(t+1)}, \bar{\mu}^{(t+1)} - \mu^{(t+1)} \rangle, \end{aligned} \quad (44)$$

where the second line follows from (44). From this point, we shall continue the proofs for PU and OMWU separately but follow similar strategies.

Remaining steps for PU. Plugging relation (32) into (44), we arrive at

$$\begin{aligned} &\eta(f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})) \\ &\leq (1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) - (1 - \eta\tau) \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) - \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}) \\ &\quad - (1 - \eta\tau) \|A\|_\infty \text{KL}(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) + \eta \|A\|_\infty \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) \\ &\leq (1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) - \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}) - (1 - \eta\tau) \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) + \eta \|A\|_\infty \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}), \end{aligned} \quad (45)$$

where the last line holds since $\eta(\tau + \|A\|_\infty) \leq 1$. Similarly, from Lemma 1 inequality (20b), one can establish the following inequality in parallel

$$\begin{aligned} &\eta(f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)}) - f_\tau(\mu_\tau^*, \nu_\tau^*)) \\ &\leq (1 - \eta\tau) \text{KL}(\nu_\tau^* \parallel \nu^{(t)}) - \text{KL}(\nu_\tau^* \parallel \nu^{(t+1)}) - (1 - \eta\tau) \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) + \eta \|A\|_\infty \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}). \end{aligned} \quad (46)$$

We are ready to establish inequality (9c) for PU. Computing $(1 - \eta\tau) \cdot (45) + \eta \|A\|_\infty \cdot (46)$ gives

$$\eta[1 - (\|A\|_\infty + \tau)\eta](f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)}))$$

$$\begin{aligned}
&\leq (1 - \eta\tau) \left[(1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t)}) \right] \\
&\quad - \left[(1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t+1)}) \right] - \left[(1 - \eta\tau)^2 - \eta^2 \|A\|_\infty^2 \right] \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) \\
&\leq (1 - \eta\tau) \left[(1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t)}) \right] \\
&\quad - \left[(1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t+1)}) \right]. \tag{47}
\end{aligned}$$

Here, the last step is due to the fact that $1 - \eta\tau \geq \eta \|A\|_\infty \geq 0$ when $0 < \eta \leq \frac{1}{\tau + \|A\|_\infty}$. As a direct consequence, the difference $f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t)}, \bar{\nu}^{(t)})$ satisfies

$$\begin{aligned}
&\eta \left[1 - (\|A\|_\infty + \tau)\eta \right] (f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t)}, \bar{\nu}^{(t)})) \\
&\leq (1 - \eta\tau) \left[(1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t-1)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t-1)}) \right] \\
&\leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t-1)}) \leq (1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}).
\end{aligned}$$

We conclude by noting that the other side of (9c) can be shown by considering $\eta \|A\|_\infty \cdot (45) + (1 - \eta\tau) \cdot (46)$ combined with similar arguments, and are therefore omitted.

Remaining steps for OMWU. Similar to the case of PU, plugging (36) into (44) gives

$$\begin{aligned}
&\eta (f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})) \\
&\leq (1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) - (1 - \eta\tau) \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) - \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}) - (1 - 2\eta \|A\|_\infty) \text{KL}(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) \\
&\quad + \eta \|A\|_\infty \left[\text{KL}(\nu^{(t)} \parallel \bar{\nu}^{(t)}) + \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) \right]. \tag{48}
\end{aligned}$$

Similarly, one can establish a symmetric inequality as follows

$$\begin{aligned}
&\eta (f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)}) - f_\tau(\mu_\tau^*, \nu_\tau^*)) \\
&\leq (1 - \eta\tau) \text{KL}(\nu_\tau^* \parallel \nu^{(t)}) - (1 - \eta\tau) \text{KL}(\bar{\nu}^{(t+1)} \parallel \nu^{(t)}) - \text{KL}(\nu_\tau^* \parallel \nu^{(t+1)}) - (1 - 2\eta \|A\|_\infty) \text{KL}(\nu^{(t+1)} \parallel \bar{\nu}^{(t+1)}) \\
&\quad + \eta \|A\|_\infty \left[\text{KL}(\mu^{(t)} \parallel \bar{\mu}^{(t)}) + \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) \right]. \tag{49}
\end{aligned}$$

Directly computing $(1 - \eta\tau) \cdot (48) + \eta \|A\|_\infty \cdot (49)$ gives

$$\begin{aligned}
&\eta (1 - (\|A\|_\infty + \tau)\eta) (f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})) \\
&\leq (1 - \eta\tau) \left[(1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t)}) \right] \\
&\quad - \left[(1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t+1)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t+1)}) \right] - \left[(1 - \eta\tau)^2 - \eta^2 \|A\|_\infty^2 \right] \text{KL}(\bar{\mu}^{(t+1)} \parallel \mu^{(t)}) \\
&\quad + \eta \|A\|_\infty \left[(1 - \eta\tau) \text{KL}(\mu^{(t)} \parallel \bar{\mu}^{(t)}) + \eta \|A\|_\infty \text{KL}(\nu^{(t)} \parallel \bar{\nu}^{(t)}) \right] \\
&\quad - (1 - 2\eta \|A\|_\infty) \left[(1 - \eta\tau) \text{KL}(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) + \eta \|A\|_\infty \text{KL}(\nu^{(t+1)} \parallel \bar{\nu}^{(t+1)}) \right]. \tag{50}
\end{aligned}$$

With our choice of the learning rate $\eta \leq \min\{\frac{1}{2\|A\|_\infty + 2\tau}, \frac{1}{4\|A\|_\infty}\}$, it is guaranteed that

$$(1 - \eta\tau)^2 - \eta^2 \|A\|_\infty^2 \geq 0 \quad \text{and} \quad (1 - \eta\tau)(1 - 2\eta \|A\|_\infty) \geq \eta \|A\|_\infty.$$

To proceed, let us introduce the shorthand notation

$$\begin{aligned}
G^{(t)} &:= (1 - \eta\tau) \text{KL}(\mu_\tau^* \parallel \mu^{(t)}) + \eta \|A\|_\infty \text{KL}(\nu_\tau^* \parallel \nu^{(t)}) \\
&\quad + (1 - 2\eta \|A\|_\infty) \left[(1 - \eta\tau) \text{KL}(\mu^{(t)} \parallel \bar{\mu}^{(t)}) + \eta \|A\|_\infty \text{KL}(\nu^{(t)} \parallel \bar{\nu}^{(t)}) \right].
\end{aligned}$$

With this piece of notation, we can write inequality (50) as

$$\eta (1 - (\|A\|_\infty + \tau)\eta) (f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})) \leq (1 - \eta\tau) G^{(t)} - G^{(t+1)}, \tag{51}$$

which in turn implies

$$\begin{aligned} & \eta(1 - (\|A\|_\infty + \tau)\eta)(f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})) \\ & \leq (1 - \eta\tau)G^{(t)} \leq L^{(t)} \leq (1 - \eta\tau)^t L^{(0)} = (1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}), \end{aligned}$$

with $L^{(t)}$ defined in (38). This finishes the proof of (9c) for OMWU.

A.2.4 Proof of convergence of duality gap (9d)

The proof of inequality (9d) is built upon the following lemma whose proof is deferred to Appendix C.4.

Lemma 4. *The duality gap at $\zeta = (\mu, \nu)$ can be bounded as*

$$\max_{\mu' \in \Delta(\mathcal{A})} f_\tau(\mu', \nu) - \min_{\nu' \in \Delta(\mathcal{B})} f_\tau(\mu, \nu') \leq \tau \text{KL}(\zeta \parallel \zeta_\tau^*) + \tau^{-1} \|A\|_\infty^2 \text{KL}(\zeta_\tau^* \parallel \zeta).$$

Applying Lemma 4 to $\bar{\zeta}^{(t)} = (\bar{\mu}^{(t)}, \bar{\nu}^{(t)})$ yields

$$\begin{aligned} \text{D}_\tau(\bar{\zeta}^{(t)}) & \leq \tau \text{KL}(\bar{\zeta}^{(t)} \parallel \zeta_\tau^*) + \tau^{-1} \|A\|_\infty^2 \text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t)}) \\ & \leq \tau \text{KL}(\bar{\zeta}^{(t)} \parallel \zeta_\tau^*) + 2\tau^{-1} \|A\|_\infty^2 (1 - \eta\tau)^{t-1} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}), \end{aligned} \quad (52)$$

where the second step results from (9a). It remains to bound $\tau \text{KL}(\bar{\zeta}^{(t)} \parallel \zeta_\tau^*)$, which we proceed separately for PU and OMWU.

Remaining steps for PU. From inequality (33), we are ensured that

$$\eta\tau \text{KL}(\bar{\zeta}^{(t)} \parallel \zeta_\tau^*) \leq (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t-1)}) - \text{KL}(\zeta_\tau^* \parallel \zeta^{(t)}).$$

It thus follows that

$$\tau \text{KL}(\bar{\zeta}^{(t)} \parallel \zeta_\tau^*) \leq \eta^{-1} (1 - \eta\tau) \text{KL}(\zeta_\tau^* \parallel \zeta^{(t-1)}) \leq \eta^{-1} (1 - \eta\tau)^{t-1} \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}),$$

where the last inequality is due to inequality (9a). Plugging the above inequality into (52) completes the proof of inequality (9d) for PU.

Remaining steps for OMWU. From inequality (39), we are ensured that

$$\tau \text{KL}(\bar{\zeta}^{(t)} \parallel \zeta_\tau^*) \leq \eta^{-1} (1 - \eta\tau) L^{(t-1)} \leq \eta^{-1} (1 - \eta\tau)^t L^{(0)} = \eta^{-1} (1 - \eta\tau)^t \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)}),$$

where the last equality follows from $L^{(0)} = \text{KL}(\zeta_\tau^* \parallel \zeta^{(0)})$. Plugging the above inequality into (52) finishes the proof of inequality (9d) for OMWU.

B Analysis for entropy-regularized Markov games

B.1 Proof of Proposition 2

For each t , let

$$V^{(t)}(s) := \max_{\mu(s) \in \Delta(\mathcal{A})} \min_{\nu(s) \in \Delta(\mathcal{B})} f_\tau(Q^{(t)}(s); \mu(s), \nu(s)),$$

which is, in other words, the minimax value of the associated matrix game using a payoff matrix $Q^{(t)}(s)$. We start by making a simple observation that for $\mu(s) \in \Delta(\mathcal{A}), \nu(s) \in \Delta(\mathcal{B})$,

$$\begin{aligned} |f_{Q^{(t)}(s)}(\mu(s), \nu(s)) - f_{Q_\tau^*(s)}(\mu(s), \nu(s))| & = \left| \mu^\top (Q^{(t)}(s) - Q_\tau^*(s)) \nu \right| \\ & \leq \left\| Q^{(t)}(s) - Q_\tau^*(s) \right\|_\infty \leq \left\| Q^{(t)} - Q_\tau^* \right\|_\infty. \end{aligned}$$

As a direct consequence, we can control $V^{(t)}(s) - V_\tau^*(s)$ by

$$\begin{aligned} \left| V^{(t)}(s) - V_\tau^*(s) \right| &= \left| \max_{\mu(s) \in \Delta(\mathcal{A})} \min_{\nu(s) \in \Delta(\mathcal{B})} f_{Q^{(t)}(s)}(\mu(s), \nu(s)) - \max_{\mu(s) \in \Delta(\mathcal{A})} \min_{\nu(s) \in \Delta(\mathcal{B})} f_{Q_\tau^*(s)}(\mu(s), \nu(s)) \right| \\ &\leq \left\| Q^{(t)} - Q_\tau^* \right\|_\infty. \end{aligned}$$

Recalling the definition of the soft Bellman operator \mathcal{T}_τ in (15), it then follows that

$$\begin{aligned} \left\| Q^{(t+1)} - Q_\tau^* \right\|_\infty &= \left\| \mathcal{T}_\tau(Q^{(t)}) - \mathcal{T}_\tau(Q_\tau^*) \right\| = \gamma \left\| \mathbb{E}_{s' \sim P(\cdot|s, a, b)} \left[V^{(t)}(s') - V_\tau^*(s') \right] \right\| \\ &\leq \gamma \left\| V^{(t)} - V_\tau^* \right\|_\infty \leq \gamma \left\| Q^{(t)} - Q_\tau^* \right\|_\infty. \end{aligned}$$

Recursively invoking the above inequality proves inequality (17).

B.2 Proof of Theorem 2

The inner loop of Algorithm 3 aims to solve an entropy-regularized matrix game indexed by $Q^{(t)}(s)$, which is done by running the proposed PU or OMWU methods. To analyze the efficacy of the inner loop, let us denote the exact minimax game value on state s at t -th iteration by

$$\check{V}^{(t+1)}(s) := \max_{\mu(s) \in \Delta(\mathcal{A})} \min_{\nu(s) \in \Delta(\mathcal{A})} f_\tau(Q^{(t)}(s); \mu(s), \nu(s)), \quad (53)$$

which is adopted in the exact value iteration analyzed in Proposition 2, and achieved by the equilibrium $\zeta^{*(t)} = (\mu^{*(t)}, \nu^{*(t)})$ of (53).

- Denote the output of the inner loop as $\bar{\zeta}^{(t, T_{\text{sub}})} = (\bar{\mu}^{(t, T_{\text{sub}})}(s), \bar{\nu}^{(t, T_{\text{sub}})}(s))$, which the entropy-regularized matrix game (53) is approximately solved by executing PU / OMWU for T_{sub} iterations. Theorem 1 (cf. (9c)) guarantees that for every $s \in \mathcal{S}$, one has

$$\begin{aligned} \left| V^{(t+1)}(s) - \check{V}^{(t+1)}(s) \right| &= \left| f_{Q^{(t)}(s)}^\tau(\bar{\mu}^{(t, T_{\text{sub}})}(s), \bar{\nu}^{(t, T_{\text{sub}})}(s)) - f_{Q^{(t)}(s)}^\tau(\mu^{*(t)}(s), \nu^{*(t)}(s)) \right| \\ &\leq \eta^{-1} \frac{1}{1 - (\tau + \|Q^{(t)}(s)\|_\infty)\eta} \cdot \frac{(1 - \eta\tau)^{T_{\text{sub}}}}{1 - (1 - \eta\tau)^{T_{\text{sub}}}} \cdot \text{KL}(\zeta^{*(t)} \parallel \zeta^{(0)}) \\ &\leq 2\eta^{-1} \cdot \frac{(1 - \eta\tau)^{T_{\text{sub}}}}{1 - (1 - \eta\tau)^{T_{\text{sub}}}} \cdot 2 \log |\mathcal{A}|, \end{aligned}$$

where the last step makes use of the choice of the learning rate

$$\eta = \frac{1 - \gamma}{2(1 + \tau(\log |\mathcal{A}| + 1))} \leq \frac{1}{2(\tau + \|Q^{(t)}(s)\|_\infty)}$$

and $\text{KL}(\zeta^{*(t)} \parallel \zeta^{(0)}) \leq \log |\mathcal{A}| + \log |\mathcal{B}| \leq 2 \log |\mathcal{A}|$. As a consequence, setting

$$T_{\text{sub}} = O\left(\frac{1}{\eta\tau} \left(\log \frac{1}{\epsilon} + \log \frac{1}{1 - \gamma} + \log \log |\mathcal{A}| + \log \frac{1}{\eta}\right)\right) \quad (54)$$

yields

$$\left| V^{(t+1)}(s) - \check{V}^{(t+1)}(s) \right| \leq (1 - \gamma)\epsilon, \quad \text{for all } s \in \mathcal{S}. \quad (55)$$

- We now move to monitor the progress of the outer loop. Combining (55) with some basic calculations, we arrive at

$$\begin{aligned} \left\| Q^{(t+1)} - Q_\tau^* \right\|_\infty &\leq \gamma \left\| V^{(t+1)} - V_\tau^* \right\|_\infty \leq \gamma \left\| V^{(t+1)} - \check{V}^{(t+1)} \right\|_\infty + \gamma \left\| \check{V}^{(t+1)} - V_\tau^* \right\|_\infty \\ &\leq (1 - \gamma)\epsilon + \gamma \left\| Q^{(t)} - Q_\tau^* \right\|_\infty. \end{aligned}$$

Now invoking the above relation recursively, it is ensured that

$$\left\| Q^{(t+1)} - Q_\tau^* \right\|_\infty \leq \epsilon + \gamma^{t+1} \left\| Q^{(0)} - Q_\tau^* \right\|_\infty.$$

In view of the above relation, if one takes

$$T_{\text{main}} = O\left(\frac{1}{1-\gamma} \left(\log \frac{1}{\epsilon} + \log \frac{1+\tau \log |\mathcal{A}|}{1-\gamma}\right)\right) \quad (56)$$

iterations of the outer loop in Algorithm 3, we have $\|Q^{(T_{\text{main}})} - Q_\tau^*\|_\infty \leq 2\epsilon$ as desired.

Putting things together, the total iteration complexity sufficient to achieve ϵ -accuracy equals to

$$T_{\text{main}} T_{\text{sub}} = O\left(\frac{1}{\eta\tau(1-\gamma)} \left(\log \frac{1}{\epsilon} + \log \log |\mathcal{A}| + \log \frac{1}{1-\gamma} + \log \tau\right) \left(\log \frac{1}{\epsilon} + \log \log |\mathcal{A}| + \log \frac{1}{1-\gamma} + \log \frac{1}{\eta}\right)\right).$$

Therefore the advertised iteration complexity in Theorem 2 holds true by simply noticing that $\eta = \frac{1-\gamma}{2(1+\tau(\log(|\mathcal{A}|)+1))}$ and $\tau < 1$, and hence

$$\log\left(\frac{1}{\eta}\right) \leq \log\left(\frac{2 \log |\mathcal{A}| + 4}{1-\gamma}\right), \quad \text{and} \quad \log(\tau) \leq \log\left(\frac{1}{1-\gamma}\right).$$

C Proof of auxiliary lemmas

C.1 Proof of Lemma 1

Lemma 1 follows directly from the update sequence (6) and the form of the optimal solution pair (μ_τ^*, ν_τ^*) , provided in (4). Given the update sequence (6), taking logarithm of both sides of the first equation gives

$$\log \mu^{(t+1)} = (1 - \eta\tau) \log \mu^{(t)} + \eta A z_2 + c \cdot \mathbf{1},$$

where c is the corresponding normalization constant. By rearranging terms and taking the inner product with $z_1 - \mu_\tau^*$, we have

$$\langle \log \mu^{(t+1)} - (1 - \eta\tau) \log \mu^{(t)}, z_1 - \mu_\tau^* \rangle = \eta z_1^\top A z_2 - \eta \mu_\tau^{*\top} A z_2, \quad (57)$$

Similarly, one can derive

$$\langle \log \nu^{(t+1)} - (1 - \eta\tau) \log \nu^{(t)}, z_2 - \nu_\tau^* \rangle = -\eta z_1^\top A z_2 + \eta z_1^\top A \nu_\tau^*. \quad (58)$$

By summing up equations (57) and (58), it is guaranteed that

$$\langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)}, \zeta_z - \zeta_\tau^* \rangle = -\eta \mu_\tau^{*\top} A z_2 + \eta z_1^\top A \nu_\tau^*, \quad (59)$$

where $\zeta(z) = (z_1, z_2)$.

On the other hand, recall the optimal policy pair (μ_τ^*, ν_τ^*) satisfies the following fixed point equation

$$\begin{cases} \mu_\tau^*(a) \propto \exp([A \nu_\tau^*]_a / \tau), & \forall a \in \mathcal{A}, \\ \nu_\tau^*(b) \propto \exp(-[A^\top \mu_\tau^*]_b / \tau), & \forall b \in \mathcal{B}. \end{cases}$$

Taking logarithm of both sides of the first relation gives

$$\eta\tau \log \mu_\tau^* = \eta A \nu_\tau^* + c \cdot \mathbf{1}, \quad (60)$$

for some normalization constant c . Again, by taking the inner product with $z_1 - \mu_\tau^*$, we have

$$\langle \eta\tau \log \mu_\tau^*, z_1 - \mu_\tau^* \rangle = \eta (z_1 - \mu_\tau^*)^\top A \nu_\tau^*, \quad (61)$$

and similarly

$$\langle \eta\tau \log \nu_\tau^*, z_2 - \nu_\tau^* \rangle = \eta \mu_\tau^{*\top} A (z_2 - \nu_\tau^*). \quad (62)$$

Combining inequalities (57) and (61), we arrive at inequality (20a); combining inequalities (58) and (62) gives inequality (20b). Moreover, putting together inequalities (59), (61) and (62) leads to

$$\langle \log \zeta^{(t+1)} - (1 - \eta\tau) \log \zeta^{(t)} - \eta\tau \log \zeta_\tau^*, \zeta(z) - \zeta_\tau^* \rangle = 0.$$

C.2 Proof of Lemma 2

We begin with establishing (22a). By the definition of $f_\tau(\mu, \nu)$, direct calculations yield

$$\begin{aligned} f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\mu, \nu_\tau^*) &= (\mu_\tau^* - \mu)^\top A \nu_\tau^* + \tau \mu^\top \log \mu - \tau \mu_\tau^{*\top} \log \mu_\tau^* \\ &= \tau \left(\langle \mu_\tau^* - \mu, \log \mu_\tau^* \rangle + \mu^\top \log \mu - \mu_\tau^{*\top} \log \mu_\tau^* \right) = \tau \text{KL}(\mu \| \mu_\tau^*). \end{aligned} \quad (63)$$

Here, the second equality is obtained by plugging in (60). Similarly, we have

$$f_\tau(\mu_\tau^*, \nu) - f_\tau(\mu_\tau^*, \nu_\tau^*) = \tau \text{KL}(\nu \| \nu_\tau^*). \quad (64)$$

Summing these two equalities completes the proof of (22a).

Turning to (22b), we first write

$$\begin{aligned} f_\tau(\mu, \nu) + f_\tau(\mu_\tau^*, \nu_\tau^*) &= \mu^\top A \nu + \mu_\tau^{*\top} A \nu_\tau^* + \tau \mathcal{H}(\mu) - \tau \mathcal{H}(\nu) + \tau \mathcal{H}(\mu_\tau^*) - \tau \mathcal{H}(\nu_\tau^*), \\ f_\tau(\mu_\tau^*, \nu) + f_\tau(\mu, \nu_\tau^*) &= \mu_\tau^{*\top} A \nu + \mu^\top A \nu_\tau^* + \tau \mathcal{H}(\mu_\tau^*) - \tau \mathcal{H}(\nu) + \tau \mathcal{H}(\mu) - \tau \mathcal{H}(\nu_\tau^*). \end{aligned}$$

As a consequence, taking the difference of the above two equations leads to

$$f_\tau(\mu, \nu) + f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\mu_\tau^*, \nu) - f_\tau(\mu, \nu_\tau^*) = (\mu_\tau^* - \mu)^\top A (\nu_\tau^* - \nu).$$

This in turn allows us to write $f_\tau(\mu, \nu) - f_\tau(\mu_\tau^*, \nu_\tau^*)$ as follows

$$f_\tau(\mu, \nu) - f_\tau(\mu_\tau^*, \nu_\tau^*) = (\mu_\tau^* - \mu)^\top A (\nu_\tau^* - \nu) + f_\tau(\mu_\tau^*, \nu) + f_\tau(\mu, \nu_\tau^*) - 2f_\tau(\mu_\tau^*, \nu_\tau^*). \quad (65)$$

Finally, plugging (63) and (64) into (65) reveals the desired relation (22b).

C.3 Proof of Lemma 3

By straightforward calculations, the gradient of the function $\log(\|\exp(x)\|_1)$ is given by

$$\nabla_x \log(\|\exp(x)\|_1) = \exp(x) / \|\exp(x)\|_1,$$

which implies $\|\nabla_x \log(\|\exp(x)\|_1)\|_1 = 1, \forall x \in \mathbb{R}^{|\mathcal{A}|}$. Therefore, we have

$$\begin{aligned} \|\log \mu_1 - \log \mu_2\|_\infty &= \|x_1 - x_2 - \log(\|\exp(x_1)\|_1) \cdot \mathbf{1} + \log(\|\exp(x_2)\|_1) \cdot \mathbf{1}\|_\infty \\ &\leq \|x_1 - x_2\|_\infty + \left| -\log(\|\exp(x_1)\|_1) + \log(\|\exp(x_2)\|_1) \right| \\ &= \|x_1 - x_2\|_\infty + \left| \langle x_1 - x_2, \nabla_x \log(\|\exp(x)\|_1)|_{x=x_c} \rangle \right| \\ &\leq \|x_1 - x_2\|_\infty + \left| \|x_1 - x_2\|_\infty \|\nabla_x \log(\|\exp(x)\|_1)|_{x=x_c}\|_1 \right| \\ &= 2 \|x_1 - x_2\|_\infty, \end{aligned}$$

where x_c is a certain convex combination of x_1 and x_2 .

C.4 Proof of Lemma 4

Since

$$\max_{\mu' \in \Delta(\mathcal{A})} f_\tau(\mu', \nu) - \min_{\nu' \in \Delta(\mathcal{B})} f_\tau(\mu, \nu') = \max_{\mu' \in \Delta(\mathcal{A}), \nu' \in \Delta(\mathcal{B})} f_\tau(\mu', \nu) - f_\tau(\mu, \nu'),$$

it boils down to control $f_\tau(\mu', \nu) - f_\tau(\mu, \nu')$ for any $(\mu', \nu') \in \Delta(\mathcal{A}) \times \Delta(\mathcal{B})$. Towards this, we have

$$\begin{aligned} f_\tau(\mu', \nu) - f_\tau(\mu, \nu') &= (f_\tau(\mu', \nu) - f_\tau(\mu', \nu_\tau^*) - f_\tau(\mu, \nu') + f_\tau(\mu_\tau^*, \nu')) - (f_\tau(\mu_\tau^*, \nu') - f_\tau(\mu', \nu_\tau^*)) \\ &= (f_\tau(\mu', \nu) - f_\tau(\mu', \nu_\tau^*) - f_\tau(\mu, \nu') + f_\tau(\mu_\tau^*, \nu')) - \tau \text{KL}(\zeta' \| \zeta_\tau^*), \end{aligned} \quad (66)$$

where the last step is due to $f_\tau(\mu, \nu_\tau^*) - f_\tau(\mu_\tau^*, \nu) = \tau \text{KL}(\zeta \| \zeta_\tau^*)$, as revealed in Lemma 2 (cf. (22a)).

To continue, observe that

$$\begin{aligned} f_\tau(\mu', \nu) - f_\tau(\mu', \nu_\tau^*) &= \mu'^\top A(\nu - \nu_\tau^*) + \nu^\top \log \nu - \nu_\tau^{*\top} \log \nu_\tau^* \\ &= (\mu' - \mu_\tau^*)^\top A(\nu - \nu_\tau^*) + f_\tau(\mu_\tau^*, \nu) - f_\tau(\mu_\tau^*, \nu_\tau^*). \end{aligned}$$

Similarly, we have

$$-f_\tau(\mu, \nu') + f_\tau(\mu_\tau^*, \nu') = -(\mu - \mu_\tau^*)^\top A(\nu' - \nu_\tau^*) + f_\tau(\mu_\tau^*, \nu_\tau^*) - f_\tau(\mu, \nu_\tau^*).$$

Plugging the above two equalities into (66) gives

$$\begin{aligned} f_\tau(\mu', \nu) - f_\tau(\mu, \nu') &= (\mu' - \mu_\tau^*)^\top A(\nu - \nu_\tau^*) - (\mu - \mu_\tau^*)^\top A(\nu' - \nu_\tau^*) + f_\tau(\mu_\tau^*, \nu) - f_\tau(\mu, \nu_\tau^*) - \tau \text{KL}(\zeta' \parallel \zeta_\tau^*) \\ &= (\mu' - \mu_\tau^*)^\top A(\nu - \nu_\tau^*) - (\mu - \mu_\tau^*)^\top A(\nu' - \nu_\tau^*) + \tau \text{KL}(\zeta \parallel \zeta_\tau^*) - \tau \text{KL}(\zeta' \parallel \zeta_\tau^*) \\ &\leq \|A\|_\infty (\|\mu' - \mu_\tau^*\|_1 \|\nu - \nu_\tau^*\|_1 + \|\nu' - \nu_\tau^*\|_1 \|\mu - \mu_\tau^*\|_1) + \tau \text{KL}(\zeta \parallel \zeta_\tau^*) - \tau \text{KL}(\zeta' \parallel \zeta_\tau^*) \\ &\stackrel{(i)}{\leq} \frac{1}{2} \|A\|_\infty \left[\frac{\tau}{\|A\|_\infty} (\|\mu' - \mu_\tau^*\|_1^2 + \|\nu' - \nu_\tau^*\|_1^2) + \frac{\|A\|_\infty}{\tau} (\|\mu - \mu_\tau^*\|_1^2 + \|\nu - \nu_\tau^*\|_1^2) \right] \\ &\quad + \tau \text{KL}(\zeta \parallel \zeta_\tau^*) - \tau \text{KL}(\zeta' \parallel \zeta_\tau^*) \\ &\stackrel{(ii)}{\leq} \tau \text{KL}(\zeta' \parallel \zeta_\tau^*) + \frac{\|A\|_\infty^2}{\tau} \text{KL}(\zeta_\tau^* \parallel \zeta) + \tau \text{KL}(\zeta \parallel \zeta_\tau^*) - \tau \text{KL}(\zeta' \parallel \zeta_\tau^*) \\ &= \frac{\|A\|_\infty^2}{\tau} \text{KL}(\zeta_\tau^* \parallel \zeta) + \tau \text{KL}(\zeta \parallel \zeta_\tau^*), \end{aligned}$$

where the second step invokes Lemma 2 (cf. (22a)), (i) follows from Young's inequality, namely $ab \leq \frac{a^2}{2\varepsilon} + \frac{\varepsilon b^2}{2}$ with $\varepsilon = \frac{\|A\|_\infty}{\tau}$, and (ii) results from Pinsker's inequality. Taking maximum over μ', ν' finishes the proof.

References

- Arora, S., Hazan, E., and Kale, S. (2012). The multiplicative weights update method: a meta-algorithm and applications. *Theory of Computing*, 8(1):121–164.
- Bai, Y. and Jin, C. (2020). Provable self-play algorithms for competitive reinforcement learning. In *International Conference on Machine Learning*, pages 551–560. PMLR.
- Bowling, M. and Veloso, M. (2001). Rational and convergent learning in stochastic games. In *Proceedings of the 17th international joint conference on Artificial intelligence-Volume 2*, pages 1021–1026.
- Cen, S., Cheng, C., Chen, Y., Wei, Y., and Chi, Y. (2020). Fast global convergence of natural policy gradient methods with entropy regularization. *arXiv preprint arXiv:2007.06558*.
- Daskalakis, C., Deckelbaum, A., and Kim, A. (2011). Near-optimal no-regret algorithms for zero-sum games. In *Proceedings of the twenty-second annual ACM-SIAM symposium on Discrete Algorithms*, pages 235–254. SIAM.
- Daskalakis, C., Foster, D. J., and Golowich, N. (2020). Independent policy gradient methods for competitive reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 5527–5540.
- Daskalakis, C., Ilyas, A., Syrgkanis, V., and Zeng, H. (2018). Training GANs with optimism. In *International Conference on Learning Representations (ICLR 2018)*.
- Daskalakis, C. and Panageas, I. (2018a). Last-iterate convergence: Zero-sum games and constrained min-max optimization. *arXiv preprint arXiv:1807.04252*.
- Daskalakis, C. and Panageas, I. (2018b). The limit points of (optimistic) gradient descent in min-max optimization. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 9256–9266.

- Filar, J. and Vrieze, K. (2012). *Competitive Markov decision processes*. Springer Science & Business Media.
- Freund, Y. and Schapire, R. E. (1999). Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1-2):79–103.
- Geist, M., Scherrer, B., and Pietquin, O. (2019). A theory of regularized Markov decision processes. In *International Conference on Machine Learning*, pages 2160–2169.
- Goodfellow, I. J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. (2014). Generative adversarial networks. *arXiv preprint arXiv:1406.2661*.
- Harker, P. T. and Pang, J.-S. (1990). Finite-dimensional variational inequality and nonlinear complementarity problems: a survey of theory, algorithms and applications. *Mathematical programming*, 48(1):161–220.
- Hofbauer, J. and Sandholm, W. H. (2002). On the global convergence of stochastic fictitious play. *Econometrica*, 70(6):2265–2294.
- Hsieh, Y.-G., Iutzeler, F., Malick, J., and Mertikopoulos, P. (2019). On the convergence of single-call stochastic extra-gradient methods. *arXiv preprint arXiv:1908.08465*.
- Konda, V. R. and Tsitsiklis, J. N. (2000). Actor-critic algorithms. In *Advances in neural information processing systems*, pages 1008–1014. Citeseer.
- Korpelevich, G. M. (1976). The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756.
- Lan, G. (2021). Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *arXiv preprint arXiv:2102.00135*.
- Lei, Q., Nagarajan, S. G., Panageas, I., and Wang, X. (2021). Last iterate convergence in no-regret learning: constrained min-max optimization for convex-concave landscapes. In *International Conference on Artificial Intelligence and Statistics*, pages 1441–1449. PMLR.
- Liang, T. and Stokes, J. (2019). Interaction matters: A note on non-asymptotic local convergence of generative adversarial networks. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 907–915. PMLR.
- Littman, M. L. (1994). Markov games as a framework for multi-agent reinforcement learning. In *Machine Learning Proceedings*, pages 157–163. Elsevier.
- McKelvey, R. D. and Palfrey, T. R. (1995). Quantal response equilibria for normal form games. *Games and economic behavior*, 10(1):6–38.
- Mei, J., Xiao, C., Szepesvari, C., and Schuurmans, D. (2020). On the global convergence rates of softmax policy gradient methods. *arXiv preprint arXiv:2005.06392*.
- Mertikopoulos, P., Lecouat, B., Zenati, H., Foo, C.-S., Chandrasekhar, V., and Piliouras, G. (2018a). Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. In *International Conference on Learning Representations*.
- Mertikopoulos, P., Papadimitriou, C., and Piliouras, G. (2018b). Cycles in adversarial regularized learning. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2703–2717. SIAM.
- Mertikopoulos, P. and Sandholm, W. H. (2016). Learning in games via reinforcement and regularization. *Mathematics of Operations Research*, 41(4):1297–1324.
- Mokhtari, A., Ozdaglar, A., and Pattathil, S. (2020a). A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In *International Conference on Artificial Intelligence and Statistics*, pages 1497–1507. PMLR.

- Mokhtari, A., Ozdaglar, A. E., and Pattathil, S. (2020b). Convergence rate of $O(1/k)$ for optimistic gradient and extragradient methods in smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 30(4):3230–3251.
- Nemirovski, A. (2004). Prox-method with rate of convergence $O(1/t)$ for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251.
- Neu, G., Jonsson, A., and Gómez, V. (2017). A unified view of entropy-regularized Markov decision processes. *arXiv preprint arXiv:1705.07798*.
- Neumann, J. V. (1928). Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320.
- Perolat, J., Scherrer, B., Piot, B., and Pietquin, O. (2015). Approximate dynamic programming for two-player zero-sum Markov games. In *International Conference on Machine Learning*, pages 1321–1329. PMLR.
- Rakhlin, A. and Sridharan, K. (2013). Optimization, learning, and games with predictable sequences. *arXiv preprint arXiv:1311.1869*.
- Savas, Y., Ahmadi, M., Tanaka, T., and Topcu, U. (2019). Entropy-regularized stochastic games. In *2019 IEEE 58th Conference on Decision and Control (CDC)*, pages 5955–5962. IEEE.
- Shapley, L. S. (1953). Stochastic games. *Proceedings of the National Academy of Sciences*, 39(10):1095–1100.
- Syrkkanis, V., Agarwal, A., Luo, H., and Schapire, R. E. (2015). Fast convergence of regularized learning in games. In *Proceedings of the 28th International Conference on Neural Information Processing Systems—Volume 2*, pages 2989–2997.
- Tseng, P. (1995). On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics*, 60(1-2):237–252.
- Wei, C.-Y., Lee, C.-W., Zhang, M., and Luo, H. (2021a). Last-iterate convergence of decentralized optimistic gradient descent/ascent in infinite-horizon competitive Markov games. *arXiv preprint arXiv:2102.04540*.
- Wei, C.-Y., Lee, C.-W., Zhang, M., and Luo, H. (2021b). Linear last-iterate convergence in constrained saddle-point optimization. In *International Conference on Learning Representations (ICLR)*.
- Xie, Q., Chen, Y., Wang, Z., and Yang, Z. (2020). Learning zero-sum simultaneous-move Markov games using function approximation and correlated equilibrium. In *Conference on Learning Theory*, pages 3674–3682. PMLR.
- Yadav, A., Shah, S., Xu, Z., Jacobs, D., and Goldstein, T. (2017). Stabilizing adversarial nets with prediction methods. *arXiv preprint arXiv:1705.07364*.
- Zhan, W., Cen, S., Huang, B., Chen, Y., Lee, J. D., and Chi, Y. (2021). Policy mirror descent for regularized reinforcement learning: A generalized framework with linear convergence. *arXiv preprint arXiv:2105.11066*.
- Zhang, K., Kakade, S., Basar, T., and Yang, L. (2020). Model-based multi-agent RL in zero-sum Markov games with near-optimal sample complexity. *Advances in Neural Information Processing Systems*, 33.
- Zhao, Y., Tian, Y., Lee, J. D., and Du, S. S. (2021). Provably efficient policy gradient methods for two-player zero-sum Markov games. *arXiv preprint arXiv:2102.08903*.