

FASTER LAST-ITERATE CONVERGENCE OF POLICY OPTIMIZATION IN ZERO-SUM MARKOV GAMES

Shicong Cen
Carnegie Mellon University
shicongc@andrew.cmu.edu

Yuejie Chi
Carnegie Mellon University
yuejiechi@cmu.edu

Simon S. Du
University of Washington
ssdu@cs.washington.edu

Lin Xiao
Meta AI Research
linx@fb.com

ABSTRACT

Multi-Agent Reinforcement Learning (MARL)—where multiple agents learn to interact in a shared dynamic environment—permeates across a wide range of critical applications. While there has been substantial progress on understanding the global convergence of policy optimization methods in single-agent RL, designing and analysis of efficient policy optimization algorithms in the MARL setting present significant challenges, which unfortunately, remain highly inadequately addressed by existing theory. In this paper, we focus on the most basic setting of competitive multi-agent RL, namely two-player zero-sum Markov games, and study equilibrium finding algorithms in both the infinite-horizon discounted setting and the finite-horizon episodic setting. We propose a single-loop policy optimization method with symmetric updates from both agents, where the policy is updated via the entropy-regularized optimistic multiplicative weights update (OMWU) method and the value is updated on a slower timescale. We show that, in the full-information tabular setting, the proposed method achieves a finite-time last-iterate linear convergence to the quantal response equilibrium of the regularized problem, which translates to a sublinear last-iterate convergence to the Nash equilibrium by controlling the amount of regularization. Our convergence results improve upon the best known iteration complexities, and lead to a better understanding of policy optimization in competitive Markov games.

Keywords: zero-sum Markov game, entropy regularization, policy optimization, global convergence, multiplicative updates

1 INTRODUCTION

Policy optimization methods (Williams, 1992; Sutton et al., 2000; Kakade, 2002; Peters and Schaal, 2008; Konda and Tsitsiklis, 2000), which cast sequential decision making as value maximization problems with regards to (parameterized) policies, have been instrumental in enabling recent successes of reinforcement learning (RL). See e.g., Schulman et al. (2015; 2017); Silver et al. (2016). Despite its empirical popularity, the theoretical underpinnings of policy optimization methods remain elusive until very recently. For single-agent RL problems, a flurry of recent works has made substantial progress on understanding the global convergence of policy optimization methods under the framework of Markov Decision Processes (MDP) (Agarwal et al., 2020; Bhandari and Russo, 2019; Mei et al., 2020; Cen et al., 2021a; Lan, 2022; Bhandari and Russo, 2020; Zhan et al., 2021; Khodadadian et al., 2021; Xiao, 2022). Despite the nonconcave nature of value maximization, (natural) policy gradient methods are shown to achieve global convergence at a sublinear rate (Agarwal et al., 2020; Mei et al., 2020) or even a linear rate in the presence of regularization (Mei et al., 2020; Cen et al., 2021a; Lan, 2022; Zhan et al., 2021) when the learning rate is constant.

Author are sorted alphabetically.

Moving beyond single-agent RL, Multi-Agent Reinforcement Learning (MARL) is the next frontier—where multiple agents learn to interact in a shared dynamic environment—permeating across critical applications such as multi-agent networked systems, autonomous vehicles, robotics, and so on. Designing and analysis of efficient policy optimization algorithms in the MARL setting present significant challenges and new desiderata, which unfortunately, remain highly inadequately addressed by existing theory.

1.1 POLICY OPTIMIZATION FOR COMPETITIVE RL

In this work, we focus on one of the most basic settings of competitive multi-agent RL, namely two-player zero-sum Markov games (Shapley, 1953), and study equilibrium finding algorithms in both the infinite-horizon discounted setting and the finite-horizon episodic setting. In particular, our designs gravitate around algorithms that are *single-loop*, *symmetric*, with *finite-time last-iterate* convergence to the Nash Equilibrium (NE) or Quantal Response Equilibrium (QRE) under bounded rationality, two prevalent solution concepts in game theory. These design principles naturally come up as a result of pursuing simple yet efficient algorithms: *single-loop* updates preclude sophisticated interleaving of rounds between agents; *symmetric* updates ensure no agent will compromise its rewards in the learning process, which can be otherwise exploited by a faster-updating opponent; in addition, asymmetric updates typically lead to one-sided convergence, i.e., only one of the agents is guaranteed to converge to the minimax equilibrium in a non-asymptotic manner, which is less desirable; moreover, *last-iterate convergence* guarantee absolves the need for agents to switch between learning and deployment; last but not least, it is desirable to converge as fast as possible, where the iteration complexities are *non-asymptotic* with clear dependence on salient problem parameters.

Substantial algorithmic developments have been made for finding equilibria in two-player zero-sum Markov games, where Dynamical Programming (DP) techniques have long been used as a fundamental building block, leading to prototypical iterative schemes such as Value Iteration (VI) (Shapley, 1953) and Policy Iteration (PI) (Van Der Wal, 1978; Patek and Bertsekas, 1999). Different from their single-agent counterparts, these methods require solving a two-player zero-sum matrix game for every state per iteration. A considerable number of recent works (Zhao et al., 2022; Alacaoglu et al., 2022; Cen et al., 2021b; Chen et al., 2021a) are based on these DP iterations, by plugging in various (gradient-based) solvers of two-player zero-sum matrix games. However, these methods are inherently nested-loop, which are less convenient to implement. In addition, PI-based methods are asymmetric and come with only one-sided convergence guarantees (Patek and Bertsekas, 1999; Zhao et al., 2022; Alacaoglu et al., 2022).

Going beyond nested-loop algorithms, single-loop policy gradient methods have been proposed recently for solving two-player zero-sum Markov games. Here, we are interested in finding an ϵ -optimal NE or QRE in terms of the duality gap, i.e. the difference in the value functions when either of the agents deviates from the solution policy.

- For the infinite-horizon discounted setting, Daskalakis et al. (2020) demonstrated that the independent policy gradient method, with direct parameterization and asymmetric learning rates, finds an ϵ -optimal NE within a polynomial number of iterations. Zeng et al. (2022) improved over this rate using an entropy-regularized policy gradient method with softmax parameterization and asymmetric learning rates. On the other end, Wei et al. (2021b) proposed an optimistic gradient descent ascent (OGDA) method (Rakhlin and Sridharan, 2013) with direct parameterization and symmetric learning rates,¹ which achieves a last-iterate convergence at a rather pessimistic iteration complexity.
- For the finite-horizon episodic setting, Zhang et al. (2022); Yang and Ma (2022) showed that the weighted average-iterate of the optimistic Follow-The-Regularized-Leader (FTRL) method, when combined with slow critic updates, finds an ϵ -optimal NE in a polynomial number of iterations.

A more complete summary of prior results can be found in Table 1 and Table 2. In brief, while there have been encouraging progresses in developing computationally efficient policy gradient methods

¹To be precise, Wei et al. (2021b) proved the average-iterate convergence of the duality gap, as well as the last-iterate convergence of the policy in terms of the Euclidean distance to the set of NEs, where it is possible to translate the latter last-iterate convergence to the duality gap (see Appendix G). The resulting iteration complexity, however, is much worse than that of the average-iterate convergence in terms of the duality gap, with a problem-dependent constant that can scale pessimistically with salient problem parameters.

for solving zero-sum Markov games, achieving fast finite-time last-iterate convergence with single-loop and symmetric update rules remains a challenging goal.

1.2 OUR CONTRIBUTIONS

Motivated by the positive role of entropy regularization in enabling faster convergence of policy optimization in single-agent RL (Cen et al., 2021a; Lan, 2022) and two-player zero-sum games (Cen et al., 2021b), we propose a single-loop policy optimization algorithm for two-player zero-sum Markov games in both the infinite-horizon and finite-horizon settings. The proposed algorithm follows the style of actor-critic (Konda and Tsitsiklis, 2000), with the actor updating the policy via the entropy-regularized optimistic multiplicative weights update (OMWU) method (Cen et al., 2021b) and the critic updating the value function on a slower timescale. Both agents execute multiplicative and symmetric policy updates, where the learning rates are carefully selected to ensure a fast last-iterate convergence. In both the infinite-horizon and finite-horizon settings, we prove that the last iterate of the proposed method learns the optimal value function and converges at a linear rate to the unique QRE of the entropy-regularized Markov game, which can be further translated into finding the NE by setting the regularization sufficiently small.

- For the infinite-horizon discounted setting, the last iterate of our method takes at most

$$\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}|}{(1-\gamma)^4\tau}\log\frac{1}{\epsilon}\right)$$

iterations for finding an ϵ -optimal QRE under entropy regularization, where $\tilde{\mathcal{O}}(\cdot)$ hides logarithmic dependencies. Here, $|\mathcal{S}|$ is the size of the state space, γ is the discount factor, and τ is the regularization parameter. Moreover, this implies the last-iterate convergence with an iteration complexity of

$$\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}|}{(1-\gamma)^5\epsilon}\right)$$

for finding an ϵ -optimal NE.

- For the finite-horizon episodic setting, the last iterate of our method takes at most

$$\tilde{\mathcal{O}}\left(\frac{H^2}{\tau}\log\frac{1}{\epsilon}\right)$$

iterations for finding an ϵ -optimal QRE under entropy regularization, where H is the horizon length. Similarly, this implies the last-iterate convergence with an iteration complexity of

$$\tilde{\mathcal{O}}\left(\frac{H^3}{\epsilon}\right)$$

for finding an ϵ -optimal NE.

Detailed comparisons between the proposed method and prior arts are provided in Table 1 and Table 2. To the best of our knowledge, this work presents the first method that is simultaneously single-loop, symmetric, and achieves fast finite-time last-iterate convergence in terms of the duality gap in both infinite-horizon and finite-horizon settings. From a technical perspective, the infinite-horizon discounted setting is in particular challenging, where ours is the first single-loop algorithm that guarantees an iteration complexity of $\tilde{\mathcal{O}}(1/\epsilon)$ for last-iterate convergence in terms of the duality gap, with clear and improved dependencies on other problem parameters in the meantime. In contrast, several existing works introduce additional problem-dependent constants (Daskalakis et al., 2020; Wei et al., 2021b; Zeng et al., 2022) in the iteration complexity, which can scale rather pessimistically—sometimes even exponentially—with problem dimensions (Li et al., 2021).

Our technical developments require novel ingredients that deviate from prior tools such as error propagation analysis for Bellman operators (Perolat et al., 2015; Patek and Bertsekas, 1999) from a dynamic programming perspective, as well as the gradient dominance condition (Daskalakis et al., 2020; Zeng et al., 2022) from a policy optimization perspective. Importantly, at the core of our analysis lies a carefully-designed one-step error contraction bound for policy learning, together with a set of recursive error bounds for value learning, all of which tailored to the non-Euclidean OMWU update rules that have not been well studied in the setting of Markov games.

Solution type	Reference	Iteration complexity	Single loop	Symmetric	Last-iterate convergence
ϵ -NE	PI-based Methods Zhao et al. (2022) Alacaoglu et al. (2022)	$\tilde{\mathcal{O}}\left(\frac{\ 1/\rho\ _\infty}{(1-\gamma)^3\epsilon}\right)^*$	✗	✗	✓
	VI-based Methods Cen et al. (2021b) Chen et al. (2021a)	$\tilde{\mathcal{O}}\left(\frac{1}{(1-\gamma)^3\epsilon}\right)$	✗	✓	✓
	Daskalakis et al. (2020)	Polynomial*	✓	✗	✗
	Zeng et al. (2022)	$\tilde{\mathcal{O}}\left(\frac{ S ^2\ 1/\rho\ _\infty^5}{(1-\gamma)^{14}c^4\epsilon^3}\right)^*$	✓	✗	✓
	Wei et al. (2021b)	$\tilde{\mathcal{O}}\left(\frac{ S ^3}{(1-\gamma)^9\epsilon^2}\right)$	✓	✓	✗
		$\tilde{\mathcal{O}}\left(\frac{ S ^5(A + B)^{1/2}}{(1-\gamma)^{16}c^4\epsilon^2}\right)$	✓	✓	✓
This Work	$\tilde{\mathcal{O}}\left(\frac{ S }{(1-\gamma)^5\epsilon}\right)$	✓	✓	✓	
ϵ -QRE	VI-based Methods Cen et al. (2021b)	$\tilde{\mathcal{O}}\left(\frac{1}{(1-\gamma)^3}\log^2\frac{1}{\epsilon}\right)$	✗	✓	✓
	Zeng et al. (2022)	$\tilde{\mathcal{O}}\left(\frac{ S ^2\ 1/\rho\ _\infty^5}{(1-\gamma)^{11}c^4\tau^3}\log\frac{1}{\epsilon}\right)^*$	✓	✗	✓
	This Work	$\tilde{\mathcal{O}}\left(\frac{ S }{(1-\gamma)^4\tau}\log\frac{1}{\epsilon}\right)$	✓	✓	✓

Table 1: Comparison of policy optimization methods for finding an ϵ -optimal NE or QRE of two-player zero-sum discounted Markov games in terms of the duality gap. Note that * implies one-sided convergence, i.e., only one of the agents is guaranteed to achieve finite-time convergence to the equilibrium. Here, $c > 0$ refers to some problem-dependent constant. For simplicity and a fair comparison, we replace various notions of concentrability coefficient and distribution mismatch coefficient with a crude upper bound $\|1/\rho\|_\infty$, where ρ is the initial state distribution.

Solution type	Reference	Iteration complexity	Single loop	Symmetric	Last-iterate convergence
ϵ -NE	Zhang et al. (2022) OFTRL	$\tilde{\mathcal{O}}\left(\frac{H^{28/5}}{\epsilon^{6/5}}\right)$	✓	✓	✗
	Zhang et al. (2022) modified OFTRL	$\tilde{\mathcal{O}}\left(\frac{H^4}{\epsilon}\right)$	✓	✓	✗
	Yang and Ma (2022) OFTRL	$\tilde{\mathcal{O}}\left(\frac{H^5}{\epsilon}\right)$	✓	✓	✗
	This Work	$\tilde{\mathcal{O}}\left(\frac{H^3}{\epsilon}\right)$	✓	✓	✓
ϵ -QRE	This Work	$\tilde{\mathcal{O}}\left(\frac{H^2}{\tau}\log\frac{1}{\epsilon}\right)$	✓	✓	✓

Table 2: Comparison of policy optimization methods for finding an ϵ -optimal NE or QRE of two-player zero-sum episodic Markov games in terms of the duality gap.

1.3 RELATED WORKS

Learning in two-player zero-sum matrix games. [Freund and Schapire \(1999\)](#) showed that the average iterate of Multiplicative Weight Update (MWU) method converges to NE at a rate of $\mathcal{O}(1/\sqrt{T})$, which in principle holds for many other no-regret algorithms as well. [Daskalakis et al. \(2011\)](#) deployed the excessive gap technique of Nesterov and improved the convergence rate to $\mathcal{O}(1/T)$, which is achieved later by [Rakhlin and Sridharan, 2013](#)) with a simple modification of MWU method, named Optimistic Mirror Descent (OMD) or more commonly, OMWU. Moving beyond average-iterate convergence, [Bailey and Piliouras \(2018\)](#) demonstrated that MWU updates, despite converging in an ergodic manner, diverge from the equilibrium. [Daskalakis and Panageas \(2018\)](#); [Wei et al. \(2021a\)](#) explored the last-iterate convergence guarantee of OMWU, as-

suming uniqueness of NE. [Cen et al. \(2021b\)](#) established linear last-iterate convergence of entropy-regularized OMWU without uniqueness assumption. [Sokota et al. \(2022\)](#) showed that optimistic update is not necessary for achieving linear last-iterate convergence in the presence of regularization, albeit with a more strict restriction on the step size.

Learning in two-player zero-sum Markov games. In addition to the aforementioned works on policy optimization methods (policy-based methods) for two-player zero-sum Markov games (cf. Table 1 and Table 2), a growing body of works have developed model-based methods ([Liu et al., 2021](#); [Zhang et al., 2020](#); [Li et al., 2022](#)) and value-based methods ([Bai and Jin, 2020](#); [Bai et al., 2020](#); [Chen et al., 2021b](#); [Jin et al., 2021](#); [Sayin et al., 2021](#); [Xie et al., 2020](#)), with a primary focus on learning NE in a sample-efficient manner. Our work, together with prior literatures on policy optimization, focuses instead on learning NE in a computation-efficient manner assuming full-information.

Entropy regularization in RL and games. Entropy regularization is a popular algorithmic idea in RL ([Williams and Peng, 1991](#)) that promotes exploration of the policy. A recent line of works ([Mei et al., 2020](#); [Cen et al., 2021a](#); [Lan, 2022](#); [Zhan et al., 2021](#)) demonstrated that incorporating entropy regularization provably accelerates policy optimization in single-agent MDPs by enabling fast linear convergence. While the positive role of entropy regularization is also verified in various game-theoretic settings, e.g., two-player zero-sum matrix games ([Cen et al., 2021b](#)), zero-sum polymatrix games ([Leonardos et al., 2021](#)), and potential games ([Cen et al., 2022](#)), it remains highly unexplored the interplay between entropy regularization and policy optimization in Markov games with only a few exceptions ([Zeng et al., 2022](#)).

1.4 NOTATIONS

We denote the probability simplex over a set \mathcal{A} by $\Delta(\mathcal{A})$. We use bracket with subscript to index the entries of a vector or matrix, e.g., $[x]_a$ for a -th element of a vector x , or simply $x(a)$ when it is clear from the context. Given two distributions $x, y \in \Delta(\mathcal{A})$, the Kullback-Leibler (KL) divergence from y to x is denoted by $\text{KL}(x \parallel y) = \sum_{a \in \mathcal{A}} x(a)(\log x(a) - \log y(a))$. Finally, we denote by $\|A\|_\infty$ the maximum entrywise absolute value of a matrix A , i.e., $\|A\|_\infty = \max_{i,j} |A_{i,j}|$.

2 ALGORITHM AND THEORY: THE INFINITE-HORIZON SETTING

2.1 PROBLEM FORMULATION

Two-player zero-sum discounted Markov game. A two-player zero-sum discounted Markov game is defined by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{B}, P, r, \gamma)$, with finite state space \mathcal{S} , finite action spaces of the two players \mathcal{A} and \mathcal{B} , reward function $r : \mathcal{S} \times \mathcal{A} \times \mathcal{B} \rightarrow [0, 1]$, transition probability kernel $P : \mathcal{S} \times \mathcal{A} \times \mathcal{B} \rightarrow \Delta(\mathcal{S})$ and discount factor $0 \leq \gamma < 1$. The action selection rule of the max player (resp. the min player) is represented by $\mu : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ (resp. $\nu : \mathcal{S} \rightarrow \Delta(\mathcal{B})$), where the probability of selecting action $a \in \mathcal{A}$ (resp. $b \in \mathcal{B}$) in state $s \in \mathcal{S}$ is specified by $\mu(a|s)$ (resp. $\nu(b|s)$). The probability of transitioning from state s to a new state s' upon selecting the action pair $(a, b) \in \mathcal{A}, \mathcal{B}$ is given by $P(s'|s, a, b)$.

Value function and Q-function. For a given policy pair μ, ν , the state value of $s \in \mathcal{S}$ is evaluated by the expected discounted sum of rewards with initial state $s_0 = s$:

$$\forall s \in \mathcal{S} : \quad V^{\mu, \nu}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t) \mid s_0 = s \right], \quad (1)$$

the quantity the max player seeks to maximize while the min player seeks to minimize. Here, the trajectory $(s_0, a_0, b_0, s_1, \dots)$ is generated according to $a_t \sim \mu(\cdot|s_t)$, $b_t \sim \nu(\cdot|s_t)$ and $s_{t+1} \sim P(\cdot|s_t, a_t, b_t)$. Similarly, the Q -function $Q^{\mu, \nu}(s, a, b)$ evaluates the expected discounted cumulative reward with initial state s and initial action pair (a, b) :

$$\forall (s, a, b) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B} : \quad Q^{\mu, \nu}(s, a, b) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t, b_t) \mid s_0 = s, a_0 = a, b_0 = b \right]. \quad (2)$$

For notation simplicity, we denote by $Q^{\mu,\nu}(s) \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{B}|}$ the matrix $[Q^{\mu,\nu}(s, a, b)]_{(a,b) \in \mathcal{A} \times \mathcal{B}}$, so that

$$\forall s \in \mathcal{S}: \quad V^{\mu,\nu}(s) = \mu(s)^\top Q^{\mu,\nu}(s) \nu(s).$$

Shapley (1953) proved the existence of a policy pair (μ^*, ν^*) that solves the min-max problem $\max_{\mu} \min_{\nu} V^{\mu,\nu}(s)$ for all $s \in \mathcal{S}$ simultaneously, and that the mini-max value is unique. A set of such optimal policy pair (μ^*, ν^*) is called the Nash equilibrium (NE) to the Markov game.

Entropy regularized two-player zero-sum Markov game. Entropy regularization is shown to provably accelerate convergence in single-agent RL (Geist et al., 2019; Mei et al., 2020; Cen et al., 2021a) and facilitate the analysis in two-player zero-sum matrix games (Cen et al., 2021b) as well as Markov games (Cen et al., 2021b; Zeng et al., 2022). The entropy-regularized value function $V_\tau^{\mu,\nu}(s)$ is defined as

$$\forall s \in \mathcal{S}: \quad V_\tau^{\mu,\nu}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \left(r(s_t, a_t, b_t) - \tau \log \mu(a_t | s_t) + \tau \log \nu(b_t | s_t) \right) \middle| s_0 = s \right], \quad (3)$$

where $\tau \geq 0$ is the regularization parameter. Similarly, the regularized Q -function $Q_\tau^{\mu,\nu}$ is given by

$$\forall (s, a, b) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B}: \quad Q_\tau^{\mu,\nu}(s) = r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} [V_\tau^{\mu,\nu}(s')]. \quad (4)$$

It is known that (Cen et al., 2021b) there exists a unique pair of policy (μ_τ^*, ν_τ^*) that solves the min-max entropy-regularized problem $\max_{\mu} \min_{\nu} V_\tau^{\mu,\nu}(s)$, or equivalently

$$\max_{\mu} \min_{\nu} \mu(s)^\top Q_\tau^{\mu,\nu}(s) \nu(s) + \tau \mathcal{H}(\mu(s)) - \tau \mathcal{H}(\nu(s)) \quad (5)$$

for all $s \in \mathcal{S}$, and we call (μ_τ^*, ν_τ^*) the quantal response equilibrium (QRE) (McKelvey and Palfrey, 1995) to the entropy-regularized Markov game. We denote the associated regularized value function and Q -function by $V_\tau^*(s) = V_\tau^{\mu_\tau^*, \nu_\tau^*}(s)$ and $Q_\tau^*(s, a, b) = Q_\tau^{\mu_\tau^*, \nu_\tau^*}(s, a, b)$.

Goal. We seek to find an ϵ -optimal QRE or ϵ -QRE (resp. ϵ -optimal NE or ϵ -NE) $\zeta = (\mu, \nu)$ which satisfies

$$\max_{s \in \mathcal{S}, \mu', \nu'} \left(V_\tau^{\mu', \nu'}(s) - V_\tau^{\mu, \nu'}(s) \right) \leq \epsilon \quad (6)$$

(resp. $\max_{s \in \mathcal{S}, \mu', \nu'} \left(V_\tau^{\mu', \nu'}(s) - V_\tau^{\mu, \nu'}(s) \right) \leq \epsilon$) in a computationally efficient manner. In truth, the solution concept of ϵ -QRE provides an approximation of ϵ -NE with appropriate choice of the regularization parameter τ . Basic calculations tell us that

$$\begin{aligned} V_\tau^{\mu', \nu'}(s) - V_\tau^{\mu, \nu'}(s) &= (V_\tau^{\mu', \nu'}(s) - V_\tau^{\mu, \nu'}(s)) + (V_\tau^{\mu', \nu'}(s) - V_\tau^{\mu', \nu'}(s)) - (V_\tau^{\mu, \nu'}(s) - V_\tau^{\mu', \nu'}(s)) \\ &\leq V_\tau^{\mu', \nu'}(s) - V_\tau^{\mu, \nu'}(s) + \frac{\tau(\log |\mathcal{A}| + \log |\mathcal{B}|)}{1 - \gamma}, \end{aligned}$$

which guarantees that an $\epsilon/2$ -QRE is an ϵ -NE as long as $\tau \leq \frac{(1-\gamma)\epsilon}{2(\log |\mathcal{A}| + \log |\mathcal{B}|)}$. For technical convenience, we assume $\tau \leq \frac{1}{\max\{1, \log |\mathcal{A}| + \log |\mathcal{B}|\}}$ throughout the paper.

Additional notation. For notation convenience, we denote by ζ the concatenation of a policy pair μ and ν , i.e., $\zeta = (\mu, \nu)$. The QRE to the regularized problem is denoted by $\zeta_\tau^* = (\mu_\tau^*, \nu_\tau^*)$. We use shorthand notation $\mu(s)$ and $\nu(s)$ to denote $\mu(\cdot | s)$ and $\nu(\cdot | s)$. In addition, we write $\text{KL}(\mu(s) \parallel \mu'(s))$ and $\text{KL}(\nu(s) \parallel \nu'(s))$ as $\text{KL}_s(\mu \parallel \mu')$ and $\text{KL}_s(\nu \parallel \nu')$, and let

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

By a slight abuse of notation, $\text{KL}_\rho(\zeta \parallel \zeta')$ denotes $\mathbb{E}_{s \sim \rho} [\text{KL}_s(\zeta \parallel \zeta')]$ for $\rho \in \Delta(\mathcal{S})$.

2.2 SINGLE-LOOP ALGORITHM DESIGN

In this section, we propose a single-loop policy optimization algorithm for finding the QRE of the entropy-regularized Markov game, which is generalized from the entropy-regularized OMWU method (Cen et al., 2021b) for solving entropy-regularized matrix games, with a careful orchestrating of the policy update and the value update.

Algorithm 1: Entropy-regularized OMWU for Discounted Two-player Zero-sum Markov Game

- 1 **Input:** Regularization parameter $\tau > 0$, learning rate for policy update $\eta > 0$, learning rate for value update $\{\alpha_t\}_{t=1}^{\infty}$.
 2 **Initialization:** Set $\mu^{(0)}, \bar{\mu}^{(0)}, \nu^{(0)}$ and $\bar{\nu}^{(0)}$ as uniform policies; and set

$$Q^{(0)} = 0, \quad V^{(0)} = \tau(\log |\mathcal{A}| - \log |\mathcal{B}|).$$

3 **for** $t = 0, 1, \dots$ **do**

4 **for all** $s \in \mathcal{S}$ **do in parallel**

5 When $t \geq 1$, update policy pair $\zeta^{(t)}(s)$ as:

$$\begin{cases} \mu^{(t)}(a|s) \propto \mu^{(t-1)}(a|s)^{1-\eta\tau} \exp(\eta[Q^{(t)}(s)\bar{\nu}^{(t)}(s)]_a) \\ \nu^{(t)}(b|s) \propto \nu^{(t-1)}(b|s)^{1-\eta\tau} \exp(-\eta[Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s)]_b) \end{cases} \quad (9a)$$

6 Update policy pair $\bar{\zeta}^{(t+1)}(s)$ as:

$$\begin{cases} \bar{\mu}^{(t+1)}(a|s) \propto \mu^{(t)}(a|s)^{1-\eta\tau} \exp(\eta[Q^{(t)}(s)\bar{\nu}^{(t)}(s)]_a) \\ \bar{\nu}^{(t+1)}(b|s) \propto \nu^{(t)}(b|s)^{1-\eta\tau} \exp(-\eta[Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s)]_b) \end{cases} \quad (9b)$$

7 Update $Q^{(t+1)}(s)$ and $V^{(t+1)}(s)$ as

$$\begin{cases} Q^{(t+1)}(s, a, b) = r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a, b)} [V^{(t)}(s')] \\ V^{(t+1)}(s) = (1 - \alpha_{t+1})V^{(t)}(s) \\ \quad + \alpha_{t+1} [\bar{\mu}^{(t+1)}(s)^\top Q^{(t+1)}(s)\bar{\nu}^{(t+1)}(s) + \tau \mathcal{H}(\bar{\mu}^{(t+1)}(s)) - \tau \mathcal{H}(\bar{\nu}^{(t+1)}(s))] \end{cases} \quad (10)$$

Review: entropy-regularized OMWU for two-player zero-sum matrix games. We briefly review the algorithm design of entropy-regularized OMWU method for two-player zero-sum matrix game (Cen et al., 2021b). The problem of interest can be described as

$$\max_{\mu \in \Delta(\mathcal{A})} \min_{\nu \in \Delta(\mathcal{B})} \mu^\top A \nu + \tau \mathcal{H}(\mu) - \tau \mathcal{H}(\nu), \quad (7)$$

where $A \in \mathbb{R}^{|\mathcal{A}| \times |\mathcal{B}|}$ is the payoff matrix of the game. The update rule of entropy-regularized OMWU with learning rate $\eta > 0$ is defined as follows: $\forall a \in \mathcal{A}, b \in \mathcal{B}$,

$$\begin{cases} \mu^{(t)}(a) \propto \mu^{(t-1)}(a)^{1-\eta\tau} \exp(\eta[A\bar{\nu}^{(t)}]_a) \\ \nu^{(t)}(b) \propto \nu^{(t-1)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \bar{\mu}^{(t)}]_b) \end{cases}, \quad (8a)$$

$$\begin{cases} \bar{\mu}^{(t+1)}(a) \propto \mu^{(t)}(a)^{1-\eta\tau} \exp(\eta[A\bar{\nu}^{(t)}]_a) \\ \bar{\nu}^{(t+1)}(b) \propto \nu^{(t)}(b)^{1-\eta\tau} \exp(-\eta[A^\top \bar{\mu}^{(t)}]_b) \end{cases}. \quad (8b)$$

We remark that the update rule can be alternatively motivated from the perspective of natural policy gradient (Kakade, 2002; Cen et al., 2021a) or mirror descent (Lan, 2022; Zhan et al., 2021) with optimistic updates. In particular, 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. Cen et al. (2021b) established that the last iterate of entropy-regularized OMWU converges to the QRE of the matrix game (7) at a linear rate $(1 - \eta\tau)^t$, as long as the step size η is no larger than $\min \left\{ \frac{1}{2\|A\|_\infty + 2\tau}, \frac{1}{4\|A\|_\infty} \right\}$.

Single-loop algorithm for two-player zero-sum Markov games. In view of the similarity in the problem formulations of (5) and (7), it is tempting to apply the aforementioned method to the Markov game in a state-wise manner, where the Q -function assumes the role of the payoff matrix. It is worth noting, however, that Q -function depends on the policy pair $\zeta = (\mu, \nu)$ and is hence changing concurrently with the update of the policy pair. We take inspiration from Bai et al. (2020); Wei et al. (2021b) and equip the entropy-regularized OMWU method with the following update rule that iteratively approximates the value function in an actor-critic fashion:

$$Q^{(t+1)}(s, a, b) = r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot|s, a, b)} [V^{(t)}(s')],$$

where $V^{(t+1)}$ is updated as a convex combination of the previous $V^{(t)}$ and the regularized game value induced by $Q^{(t+1)}$ as well as the policy pair $\bar{\zeta}^{(t+1)} = (\bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)})$:

$$V^{(t+1)}(s) = (1 - \alpha_{t+1})V^{(t)}(s) + \alpha_{t+1} \left[\bar{\mu}^{(t+1)}(s)^\top Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s) + \tau \mathcal{H}(\bar{\mu}^{(t+1)}(s)) - \tau \mathcal{H}(\bar{\nu}^{(t+1)}(s)) \right]. \quad (11)$$

The update of V becomes more conservative with a smaller learning rate α_t , hence stabilizing the update of policies. However, setting α_t too small slows down the convergence of V to V_τ^* . A key novelty—suggested by our analysis—is the choice of the constant learning rates $\alpha := \alpha_t = \eta\tau$ which updates at a slower timescale than the policy due to $\tau < 1$. This is in sharp contrast to the vanishing sequence $\alpha_t = \frac{2/(1-\gamma)+1}{2/(1-\gamma)+t}$ adopted in [Wei et al. \(2021b\)](#), which is essential in their analysis but inevitably leads to a much slower convergence. We summarize the detailed procedure in Algorithm 1. Last but not least, it is worth noting that the proposed method access the reward via “first-order information”, i.e., either agent can only update its policy with the marginalized value function $Q(s)\nu(s)$ or $Q(s)^\top\mu(s)$. Update rules of this kind are instrumental in breaking the curse of multi-agents in the sample complexity when adopting sample-based estimates in (10), as we only need to estimate the marginalized Q-function rather than its full form ([Li et al., 2022](#); [Chen et al., 2021a](#)).

2.3 THEORETICAL GUARANTEES

Below we present our main results concerning the last-iterate convergence of Algorithm 1 for solving entropy-regularized two-player zero-sum Markov games in the infinite-horizon discounted setting. The proof is postponed to Appendix A.

Theorem 1. *Setting $0 < \eta \leq \frac{(1-\gamma)^3}{32000|\mathcal{S}|}$ and $\alpha_t = \eta\tau$, it holds for all $t \geq 0$ that*

$$\max \left\{ \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t)}), \frac{1}{2|\mathcal{S}|} \sum_{s \in \mathcal{S}} \text{KL}_s(\zeta_\tau^* \parallel \bar{\zeta}^{(t)}), \frac{3\eta}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty \right\} \leq \frac{3000}{(1-\gamma)^2\tau} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t; \quad (12a)$$

and

$$\max_{s \in \mathcal{S}, \mu, \nu} \left(V_\tau^{\mu, \bar{\nu}^{(t)}}(s) - V_\tau^{\bar{\mu}^{(t)}, \nu}(s) \right) \leq \frac{6000|\mathcal{S}|}{(1-\gamma)^3\tau} \max \left\{ \frac{8}{(1-\gamma)^2\tau}, \frac{1}{\eta} \right\} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t. \quad (12b)$$

Theorem 1 demonstrates that as long as the learning rate η is small enough, the last iterate of Algorithm 1 converges at a linear rate for the entropy-regularized Markov game. Compared with prior literatures investigating on policy optimization, our analysis focuses on the last-iterate convergence of non-Euclidean updates in the presence of entropy regularization, which appears to be the first of its kind. Several remarks are in order, with detailed comparisons in Table 1.

- **Linear convergence to the QRE.** Theorem 1 demonstrates that the last iterate of Algorithm 1 takes at most $\tilde{\mathcal{O}}\left(\frac{1}{(1-\gamma)\eta\tau} \log \frac{1}{\epsilon}\right)$ iterations to yield an ϵ -optimal policy in terms of the KL divergence to the QRE $\max_{s \in \mathcal{S}} \text{KL}_s(\zeta_\tau^* \parallel \bar{\zeta}^{(t)}) \leq \epsilon$, the entrywise error of the regularized Q-function $\|Q^{(t)} - Q_\tau^*\|_\infty \leq \epsilon$, as well as the duality gap $\max_{s \in \mathcal{S}, \mu, \nu} (V_\tau^{\mu, \bar{\nu}^{(t)}}(s) - V_\tau^{\bar{\mu}^{(t)}, \nu}(s)) \leq \epsilon$ at once. Minimizing the bound over the learning rate η , the proposed method is guaranteed to find an ϵ -QRE within $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}|}{(1-\gamma)^4\tau} \log \frac{1}{\epsilon}\right)$ iterations, which significantly improves upon the one-side convergence rate of [Zeng et al. \(2022\)](#).
- **Last-iterate convergence to ϵ -optimal NE.** By setting $\tau = \frac{(1-\gamma)\epsilon}{2(\log |\mathcal{A}| + \log |\mathcal{B}|)}$, this immediately leads to provable last-iterate convergence to an ϵ -NE, with an iteration complexity of $\tilde{\mathcal{O}}\left(\frac{|\mathcal{S}|}{(1-\gamma)^5\epsilon}\right)$, which again outperforms the convergence rate in [Wei et al. \(2021b\)](#).

Remark 1. The learning rate η is constrained to be inverse proportional to $|\mathcal{S}|$, which is for the worst case and can be potentially loosened for problems with a small concentration coefficient. We refer interested readers to Appendix A for details.

3 ALGORITHM AND THEORY: THE EPISODIC SETTING

Episodic two-player zero-sum Markov game. An episodic two-player zero-sum Markov game is defined by a tuple $\{\mathcal{S}, \mathcal{A}, \mathcal{B}, H, \{P_h\}_{h=1}^H, \{r_h\}_{h=1}^H\}$, with \mathcal{S} being a finite state space, \mathcal{A} and \mathcal{B} denoting finite action spaces of the two players, and $H > 0$ the horizon length. Every step $h \in [H]$ admits a transition probability kernel $P_h : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ and reward function $r_h : \mathcal{S} \times \mathcal{A} \times \mathcal{B} \rightarrow [0, 1]$. Furthermore, $\mu = \{\mu_h\}_{h=1}^H$ and $\nu = \{\nu_h\}_{h=1}^H$ denote the policies of the two players, where the probability of the max player choosing $a \in \mathcal{A}$ (resp. the min player choosing $b \in \mathcal{B}$) at time h is specified by $\mu_h(a|s)$ (resp. $\nu_h(b|s)$).

Entropy regularized value functions. The value function and Q-function characterize the expected cumulative reward starting from step h by following the policy pair μ, ν . For conciseness, we only present the definition of entropy-regularized value functions below and remark that their un-regularized counterparts $V_h^{\mu, \nu}$ and $Q_h^{\mu, \nu}$ can be obtained by setting $\tau = 0$. We have

$$V_{h, \tau}^{\mu, \nu}(s) = \mathbb{E} \left[\sum_{h'=h}^H [r_{h'}(s_{h'}, a_{h'}, b_{h'}) - \tau \log \mu_{h'}(a_{h'}|s_{h'}) + \tau \log \nu_{h'}(b_{h'}|s_{h'})] \mid s_h = s \right];$$

$$Q_{h, \tau}^{\mu, \nu}(s, a, b) = r_h(s, a, b) + \mathbb{E}_{s' \sim P_h(\cdot|s, a, b)} [V_{h+1, \tau}^{\mu, \nu}(s')].$$

The solution concept of NE and QRE are defined in a similar manner by focusing on the episodic versions of value functions. We again denote the unique QRE by $\zeta_\tau^* = (\mu_\tau^*, \nu_\tau^*)$.

Proposed method and convergence guarantee It is straightforward to adapt Algorithm 1 to the episodic setting with minimal modifications, with detailed procedure showcased in Algorithm 2 (cf. Appendix B). The analysis, which substantially deviates from the discounted setting, exploits the structure of finite-horizon MDP and time-inhomogeneous policies, enabling a much larger range of learning rates as showed in the following theorem.

Theorem 2. Setting $0 < \eta \leq \frac{1}{8H}$ and $\alpha_t = \eta\tau$, it holds for all $h \in [H]$ and $t \geq T_h := (H-h)T_{\text{start}}$ with $T_{\text{start}} = \lceil \frac{1}{\eta\tau} \log H \rceil$ that

$$\|Q_{h, \tau}^* - Q_h^{(t)}\|_\infty \leq (1 - \eta\tau)^{t-T_h} t^{H-h}; \quad (13a)$$

$$\max_{s \in \mathcal{S}, \mu, \nu} (V_{h, \tau}^{\mu, \bar{\nu}^{(t)}}(s) - V_{h, \tau}^{\bar{\mu}^{(t)}, \nu}(s)) \leq 4(1 - \eta\tau)^{t-T_h} \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{8H}{\tau} + 6\eta t^{H-h+1} \right). \quad (13b)$$

Theorem 2 implies that the last iterate of Algorithm 2 takes no more than $\tilde{\mathcal{O}}(HT_{\text{start}} + \frac{H}{\eta\tau} \log \frac{1}{\epsilon}) = \tilde{\mathcal{O}}(\frac{H}{\eta\tau} \log \frac{1}{\epsilon})$ iterations for finding an ϵ -QRE. Minimizing the bound over the learning rate η , Algorithm 2 is guaranteed to find an ϵ -QRE in $\tilde{\mathcal{O}}\left(\frac{H^2}{\tau} \log \frac{1}{\epsilon}\right)$ iterations, which translates into an iteration complexity of $\tilde{\mathcal{O}}\left(\frac{H^3}{\epsilon}\right)$ for finding an ϵ -NE in terms of the duality gap, i.e., $\max_{s \in \mathcal{S}, h \in [H], \mu, \nu} (V_h^{\mu, \bar{\nu}^{(t)}}(s) - V_h^{\bar{\mu}^{(t)}, \nu}(s)) \leq \epsilon$, by setting $\tau = \mathcal{O}\left(\frac{\epsilon}{H(\log |\mathcal{A}| + \log |\mathcal{B}|)}\right)$.

4 DISCUSSION

This work develops policy optimization methods for zero-sum Markov games that feature single-loop and symmetric updates with provable last-iterate convergence guarantees. Our approach yields better iteration complexities in both infinite-horizon and finite-horizon settings, by adopting entropy regularization and non-Euclidean policy update. Important future directions include investigating whether larger learning rates are possible without knowing problem-dependent information a priori, extending the framework to allow function approximation, and designing sample-efficient implementations of the proposed method.

ACKNOWLEDGMENTS

The authors would like to thank Gen Li and Zeyuan Allen-Zhu for valuable discussions. Part of this work was completed while S. Cen was an intern at Meta AI Research. 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, CCF-2106778 and CNS-2148212. S. Cen is also gratefully supported by Wei Shen and Xuehong Zhang Presidential Fellowship, and Nicholas Minnici Dean’s Graduate Fellowship in Electrical and Computer Engineering at Carnegie Mellon University.

REFERENCES

- A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. Optimality and approximation with policy gradient methods in Markov decision processes. In *Conference on Learning Theory*, pages 64–66. PMLR, 2020.
- A. Alacaoglu, L. Viano, N. He, and V. Cevher. A natural actor-critic framework for zero-sum Markov games. In *International Conference on Machine Learning*, pages 307–366. PMLR, 2022.
- Y. Bai and C. Jin. Provable self-play algorithms for competitive reinforcement learning. In *International Conference on Machine Learning*, pages 551–560. PMLR, 2020.
- Y. Bai, C. Jin, and T. Yu. Near-optimal reinforcement learning with self-play. *Advances in neural information processing systems*, 33:2159–2170, 2020.
- J. P. Bailey and G. Piliouras. Multiplicative weights update in zero-sum games. In *Proceedings of the 2018 ACM Conference on Economics and Computation*, pages 321–338, 2018.
- J. Bhandari and D. Russo. Global optimality guarantees for policy gradient methods. *arXiv preprint arXiv:1906.01786*, 2019.
- J. Bhandari and D. Russo. A note on the linear convergence of policy gradient methods. *arXiv preprint arXiv:2007.11120*, 2020.
- S. Cen, C. Cheng, Y. Chen, Y. Wei, and Y. Chi. Fast global convergence of natural policy gradient methods with entropy regularization. *Operations Research*, 2021a.
- S. Cen, Y. Wei, and Y. Chi. Fast policy extragradient methods for competitive games with entropy regularization. *Advances in Neural Information Processing Systems*, 34:27952–27964, 2021b.
- S. Cen, F. Chen, and Y. Chi. Independent natural policy gradient methods for potential games: Finite-time global convergence with entropy regularization. In *2022 IEEE 61th Conference on Decision and Control (CDC)*. IEEE, 2022.
- Z. Chen, S. Ma, and Y. Zhou. Sample efficient stochastic policy extragradient algorithm for zero-sum markov game. In *International Conference on Learning Representations*, 2021a.
- Z. Chen, D. Zhou, and Q. Gu. Almost optimal algorithms for two-player Markov games with linear function approximation. *arXiv preprint arXiv:2102.07404*, 2021b.
- C. Daskalakis and I. Panageas. Last-iterate convergence: Zero-sum games and constrained min-max optimization. *arXiv preprint arXiv:1807.04252*, 2018.
- C. Daskalakis, A. Deckelbaum, and A. Kim. 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, 2011.
- C. Daskalakis, D. J. Foster, and N. Golowich. Independent policy gradient methods for competitive reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 33, pages 5527–5540, 2020.
- Y. Freund and R. E. Schapire. Adaptive game playing using multiplicative weights. *Games and Economic Behavior*, 29(1-2):79–103, 1999.

- M. Geist, B. Scherrer, and O. Pietquin. A theory of regularized Markov decision processes. In *International Conference on Machine Learning*, pages 2160–2169, 2019.
- C. Jin, Q. Liu, Y. Wang, and T. Yu. V-learning—a simple, efficient, decentralized algorithm for multiagent rl. *arXiv preprint arXiv:2110.14555*, 2021.
- S. M. Kakade. A natural policy gradient. In *Advances in neural information processing systems*, pages 1531–1538, 2002.
- S. Khodadadian, P. R. Jhunjhunwala, S. M. Varma, and S. T. Maguluri. On the linear convergence of natural policy gradient algorithm. In *2021 60th IEEE Conference on Decision and Control (CDC)*, pages 3794–3799. IEEE, 2021.
- V. R. Konda and J. N. Tsitsiklis. Actor-critic algorithms. In *Advances in neural information processing systems*, pages 1008–1014. Citeseer, 2000.
- G. Lan. Policy mirror descent for reinforcement learning: Linear convergence, new sampling complexity, and generalized problem classes. *Mathematical programming*, pages 1–48, 2022.
- S. Leonardos, G. Piliouras, and K. Spendlove. Exploration-exploitation in multi-agent competition: convergence with bounded rationality. *Advances in Neural Information Processing Systems*, 34: 26318–26331, 2021.
- G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen. Softmax policy gradient methods can take exponential time to converge. In *Conference on Learning Theory*, pages 3107–3110. PMLR, 2021.
- G. Li, Y. Chi, Y. Wei, and Y. Chen. Minimax-optimal multi-agent RL in zero-sum Markov games with a generative model. *arXiv preprint arXiv:2208.10458*, 2022.
- Q. Liu, T. Yu, Y. Bai, and C. Jin. A sharp analysis of model-based reinforcement learning with self-play. In *International Conference on Machine Learning*, pages 7001–7010. PMLR, 2021.
- R. D. McKelvey and T. R. Palfrey. Quantal response equilibria for normal form games. *Games and economic behavior*, 10(1):6–38, 1995.
- J. Mei, C. Xiao, C. Szepesvari, and D. Schuurmans. On the global convergence rates of softmax policy gradient methods. In *International Conference on Machine Learning*, pages 6820–6829. PMLR, 2020.
- S. D. Patek and D. P. Bertsekas. Stochastic shortest path games. *SIAM Journal on Control and Optimization*, 37(3):804–824, 1999.
- J. Perolat, B. Scherrer, B. Piot, and O. Pietquin. Approximate dynamic programming for two-player zero-sum Markov games. In *International Conference on Machine Learning*, pages 1321–1329. PMLR, 2015.
- J. Peters and S. Schaal. Natural actor-critic. *Neurocomputing*, 71(7-9):1180–1190, 2008.
- A. Rakhlin and K. Sridharan. Optimization, learning, and games with predictable sequences. *arXiv preprint arXiv:1311.1869*, 2013.
- M. Sayin, K. Zhang, D. Leslie, T. Basar, and A. Ozdaglar. Decentralized Q-learning in zero-sum Markov games. *Advances in Neural Information Processing Systems*, 34:18320–18334, 2021.
- J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897, 2015.
- J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- L. S. Shapley. Stochastic games. *Proceedings of the National Academy of Sciences*, 39(10):1095–1100, 1953.

- D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, et al. Mastering the game of Go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- S. Sokota, R. D’Orazio, J. Z. Kolter, N. Loizou, M. Lanctot, I. Mitliagkas, N. Brown, and C. Kroer. A unified approach to reinforcement learning, quantal response equilibria, and two-player zero-sum games. *arXiv preprint arXiv:2206.05825*, 2022.
- R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in neural information processing systems*, pages 1057–1063, 2000.
- J. Van Der Wal. Discounted markov games: Generalized policy iteration method. *Journal of Optimization Theory and Applications*, 25(1):125–138, 1978.
- C.-Y. Wei, C.-W. Lee, M. Zhang, and H. Luo. Linear last-iterate convergence in constrained saddle-point optimization. In *International Conference on Learning Representations (ICLR)*, 2021a.
- C.-Y. Wei, C.-W. Lee, M. Zhang, and H. Luo. Last-iterate convergence of decentralized optimistic gradient descent/ascent in infinite-horizon competitive markov games. In *Conference on learning theory*, pages 4259–4299. PMLR, 2021b.
- R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256, 1992.
- R. J. Williams and J. Peng. Function optimization using connectionist reinforcement learning algorithms. *Connection Science*, 3(3):241–268, 1991.
- L. Xiao. On the convergence rates of policy gradient methods. *arXiv preprint arXiv:2201.07443*, 2022.
- Q. Xie, Y. Chen, Z. Wang, and Z. Yang. Learning zero-sum simultaneous-move markov games using function approximation and correlated equilibrium. In *Conference on learning theory*, pages 3674–3682. PMLR, 2020.
- Y. Yang and C. Ma. $O(T^{-1})$ convergence of optimistic-follow-the-regularized-leader in two-player zero-sum markov games. *arXiv preprint arXiv:2209.12430*, 2022.
- S. Zeng, T. T. Doan, and J. Romberg. Regularized gradient descent ascent for two-player zero-sum Markov games. *Advances in Neural Information Processing Systems*, 35, 2022.
- W. Zhan, S. Cen, B. Huang, Y. Chen, J. D. Lee, and Y. Chi. Policy mirror descent for regularized reinforcement learning: A generalized framework with linear convergence. *arXiv preprint arXiv:2105.11066*, 2021.
- K. Zhang, S. Kakade, T. Basar, and L. Yang. Model-based multi-agent RL in zero-sum Markov games with near-optimal sample complexity. *Advances in Neural Information Processing Systems*, 33, 2020.
- R. Zhang, Q. Liu, H. Wang, C. Xiong, N. Li, and Y. Bai. Policy optimization for Markov games: Unified framework and faster convergence. *Advances in Neural Information Processing Systems*, 35, 2022.
- Y. Zhao, Y. Tian, J. Lee, and S. Du. Provably efficient policy optimization for two-player zero-sum markov games. In *International Conference on Artificial Intelligence and Statistics*, pages 2736–2761. PMLR, 2022.

Appendix

Table of Contents

A	Analysis for the infinite-horizon setting	13
B	Analysis for the episodic setting	17
C	Proof of key lemmas for the discounted setting	20
C.1	Proof of Lemma 1	20
C.2	Proof of Lemma 2	22
C.3	Proof of Lemma 3	25
C.4	Proof of Lemma 4	27
C.5	Proof of Lemma 5	29
C.6	Proof of Lemma 6	30
D	Proof of key lemmas for the episodic setting	30
D.1	Proof of Lemma 9	30
D.2	Proof of Lemma 10	32
E	Proof of auxiliary lemmas	34
E.1	Proof of Lemma 11	34
E.2	Proof of Lemma 12	35
E.3	Proof of Lemma 13	37
E.4	Proof of Lemma 14	37
E.5	Proof of Lemma 16	38
E.6	Proof of Lemma 17	39
F	Further discussion regarding approximate algorithms	40
G	Further discussion regarding Wei et al. (2021b)	41

A ANALYSIS FOR THE INFINITE-HORIZON SETTING

Definition 1. Given $\rho \in \Delta(\mathcal{S})$ with $\rho(s) > 0, \forall s \in \mathcal{S}$, concentrability coefficient $c_\rho(t)$ is defined as

$$c_\rho(t) = \sup_{\substack{x^{(l)} \in \mathcal{A}^{\mathcal{S}}, 1 \leq l \leq t, \\ y^{(l)} \in \mathcal{B}^{\mathcal{S}}, 1 \leq l \leq t}} \left\| \frac{\rho P_{x^{(1)}, y^{(1)}} \cdots P_{x^{(t)}, y^{(t)}}}{\rho} \right\|_\infty,$$

where $P_{x^{(l)}, y^{(l)}} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ is the state transition matrix induced by a pair of deterministic policy $x^{(l)}, y^{(l)}$:

$$[P_{x^{(l)}, y^{(l)}}]_{s, s'} = P(s' | s, x^{(l)}(s), y^{(l)}(s)).$$

Let \mathcal{C}_ρ be the maximum value of $c_\rho(t)$ over $t \geq 0$:

$$\mathcal{C}_\rho = \sup_{t \geq 0} c_\rho(t).$$

In addition, let $\Gamma(\rho)$ be the set of all possible distribution over \mathcal{S} induced by initial distribution ρ and deterministic policy sequences, i.e.,

$$\Gamma(\rho) = \bigcup_{t=0}^{\infty} \{ \rho P_{x^{(1)}, y^{(1)}} \cdots P_{x^{(t)}, y^{(t)}} : x^{(l)} \in \mathcal{A}^{\mathcal{S}}, y^{(l)} \in \mathcal{B}^{\mathcal{S}}, \forall l \in [t] \}$$

We make note that Theorem 1 is the direct corollary of following theorems, by setting ρ to the uniform distribution over \mathcal{S} , where \mathcal{C}_ρ admits a trivial upper bounded $|\mathcal{S}|$.

Theorem 3. *With $0 < \eta \leq \frac{(1-\gamma)^3}{32000\mathcal{C}_\rho}$, and $\alpha_i = \eta\tau$, we have*

$$\max \left\{ \text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(t)}), \frac{1}{2} \text{KL}_\rho(\zeta_\tau^* \parallel \bar{\zeta}^{(t)}), 3\eta \mathbb{E}_{s \sim \rho} \left[\|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty \right] \right\} \leq \frac{3000}{(1-\gamma)^2\tau} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t.$$

Definition 2. *We define regularized minimax mismatch coefficient by*

$$\mathcal{C}_{\rho,\tau}^\dagger = \max \left\{ \max_\mu \left\| \frac{d_\rho^{\mu, \nu_\tau^\dagger(\mu)}}{\rho} \right\|_\infty, \max_\nu \left\| \frac{d_\rho^{\mu_\tau^\dagger(\nu), \nu}}{\rho} \right\|_\infty \right\}.$$

Here, $\nu_\tau^\dagger(\mu)$ denotes the optimal policy of the min player when the max player adopts policy μ :

$$\nu_\tau^\dagger(\mu) = \arg \min_\nu V_\tau^{\mu, \nu}(s),$$

and $\mu_\tau^\dagger(\nu)$ is defined in a symmetric way. The discounted state visitation distribution $d_\rho^{\mu, \nu}$ is defined as

$$d_\rho^{\mu, \nu}(s) = (1-\gamma) \mathbb{E}_{s_0 \sim \rho} \left[\sum_{t=0}^{\infty} \gamma^t P(s_t = s | s_0) \right].$$

Note that this definition parallels that of the (unregularized) minimax mismatch coefficient in (Daskalakis et al., 2020).

Theorem 4. *With $0 < \eta \leq \frac{(1-\gamma)^3}{32000\mathcal{C}_\rho}$, and $\alpha_i = \eta\tau$, we have*

$$\max_{s \in \mathcal{S}, \mu, \nu} \left(V_\tau^{\mu, \bar{\nu}^{(t)}}(s) - V_\tau^{\bar{\mu}^{(t)}, \nu}(s) \right) \leq \frac{6000 \|1/\rho\|_\infty}{(1-\gamma)^3\tau} \max \left\{ \frac{8}{(1-\gamma)^2\tau}, \frac{1}{\eta} \right\} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t,$$

and

$$\max_{\mu, \nu} \left(V_\tau^{\mu, \bar{\nu}^{(t)}}(\rho) - V_\tau^{\bar{\mu}^{(t)}, \nu}(\rho) \right) \leq \frac{6000\mathcal{C}_{\rho,\tau}^\dagger}{(1-\gamma)^3\tau} \max \left\{ \frac{8}{(1-\gamma)^2\tau}, \frac{1}{\eta} \right\} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t.$$

We start with the following lemma. The proof can be found in Appendix C.1. For notational simplicity, we set $Q^{(-1)} = 0$, $\bar{\zeta}^{(-1)} = \bar{\zeta}^{(0)}$ and $\alpha_0 = 1$. It follows from the update rule (9a) that $\bar{\zeta}^{(1)} = \zeta^{(0)} = \bar{\zeta}^{(0)}$.

Lemma 1. *It holds for any step size $0 < \eta \leq 1/\tau$ and $t \geq 0$ that*

$$\begin{aligned} & \text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(t+1)}) - (1-\eta\tau)\text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(t)}) + \left(1 - \eta\tau - \frac{4\eta}{1-\gamma}\right)\text{KL}_\rho(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) \\ & \quad + \eta\tau\text{KL}_\rho(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) + \left(1 - \frac{2\eta}{1-\gamma}\right)\text{KL}_\rho(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + (1-\eta\tau)\text{KL}_\rho(\bar{\zeta}^{(t)} \parallel \zeta^{(t)}) \\ & \quad - \frac{2\eta}{1-\gamma}\text{KL}_\rho(\bar{\zeta}^{(t)} \parallel \bar{\zeta}^{(t-1)}) \\ & \leq \mathbb{E}_{s \sim \rho} \left[2\eta \|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty + \frac{4\eta^2}{1-\gamma} \|Q^{(t)}(s) - Q^{(t+1)}(s)\|_\infty + \frac{12\eta^2}{1-\gamma} \|Q^{(t-1)}(s) - Q^{(t)}(s)\|_\infty \right]. \end{aligned} \tag{14}$$

It remains to bound the terms on the right hand side of (14). By a slight abuse of notation, we denote

$$\|Q^{(t+1)} - Q_\tau^*\|_{\Gamma(\rho)} = \sup_{\chi \in \Gamma(\rho)} \mathbb{E}_{s \sim \chi} \left[\|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty \right],$$

and

$$\|Q^{(t+1)} - Q^{(t)}\|_{\Gamma(\rho)} = \sup_{\chi \in \Gamma(\rho)} \mathbb{E}_{s \sim \chi} \left[\|Q^{(t+1)}(s) - Q^{(t)}(s)\|_\infty \right].$$

The following two lemmas establish a set of recursive bounds that relate $\{\|Q^{(l+1)}(s) - Q_\tau^*(s)\|_{\Gamma(\rho)}\}_{l=0, \dots, t}$ and $\{\|Q^{(l+1)}(s) - Q^{(l)}(s)\|_{\Gamma(\rho)}\}_{l=0, \dots, t}$ with $\{\text{KL}_\rho(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)})\}_{l=0, \dots, t-1}$:

Lemma 2. With $0 < \eta \leq \min\{(1 - \gamma)/180, (1 - \gamma)^2/48\}$, it holds for all $t \geq 1$ that

$$\|Q^{(t+1)} - Q^{(t)}\|_{\Gamma(\rho)} \leq \frac{1 + \gamma}{2} \sum_{l=1}^t \alpha_{l,t} \|Q^{(l)} - Q^{(l-1)}\|_{\Gamma(\rho)} + \frac{4\mathcal{C}_\rho}{\eta} \cdot \sum_{l=1}^t \alpha_{l,t} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}), \quad (15)$$

Here, $\alpha_{l,t}$ is defined as

$$\alpha_{l,t} = \alpha_l \prod_{i=l+1}^t (1 - \alpha_i).$$

When $t = 0$, we have $\|Q^{(1)}(s) - Q^{(0)}(s)\|_{\Gamma(\rho)} \leq 2$.

Proof. The proof can be found in Appendix C.2. \square

Lemma 3. With $0 < \eta \leq (1 - \gamma)^2/16$, it holds for all $t \geq 1$ that

$$\begin{aligned} & \|Q^{(t+1)} - Q_\tau^*\|_{\Gamma(\rho)} \\ & \leq \frac{1 + \gamma}{2} \cdot \sum_{l=0}^t \alpha_{l,t} \left(\|Q^{(l)} - Q_\tau^*\|_{\Gamma(\rho)} + \frac{2\eta}{1 - \gamma} \|Q^{(l)} - Q^{(l-1)}\|_{\Gamma(\rho)} \right) + 2\alpha_{0,t}. \end{aligned} \quad (16)$$

When $t = 0$, we have $\|Q^{(t+1)} - Q_\tau^*\|_{\Gamma(\rho)} \leq \frac{2\gamma}{1 - \gamma}$.

Proof. The proof can be found in Appendix C.3. \square

The following lemma further demystify the complicated recursive bounds showed in Lemma 2 and 3.

Lemma 4. Let $\lambda_{l,t}$ be defined as

$$\lambda_{l,t} = \alpha_l \prod_{i=l+1}^t \left(1 - \frac{1 - \gamma}{4} \cdot \alpha_i \right).$$

Under the assumption of Lemma 2 and 3, it holds for all $t \geq 0$ that

$$\begin{aligned} & \sum_{l=0}^t \lambda_{l+1,t+1} \left[\eta \|Q_\tau^* - Q^{(l+1)}\|_{\Gamma(\rho)} + \frac{12\eta^2}{(1 - \gamma)^2} \|Q^{(l+1)} - Q^{(l)}\|_{\Gamma(\rho)} \right] \\ & \leq \frac{6250\eta\mathcal{C}_\rho}{(1 - \gamma)^3} \sum_{l=0}^{t-1} \lambda_{l+1,t+1} \text{KL}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) + \frac{550\eta}{(1 - \gamma)^2} \lambda_{0,t+1} \end{aligned}$$

Proof. The proof can be found in Appendix C.4. \square

We are now ready to prove our main results. Averaging (14) with weight λ gives

$$\begin{aligned} & \sum_{l=0}^t \lambda_{l+1,t+1} \left[\text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(l+1)}) - (1 - \eta\tau) \text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(l)}) \right. \\ & \quad + \left(1 - \frac{2\eta}{1 - \gamma} \right) \text{KL}_\rho(\zeta^{(l+1)} \parallel \bar{\zeta}^{(l+1)}) + 3\eta \mathbb{E}_{s \sim \rho} \left[\|Q^{(l+1)}(s) - Q_\tau^*(s)\|_\infty \right] \\ & \quad \left. + \left(1 - \eta\tau - \frac{4\eta}{1 - \gamma} \right) \text{KL}_\rho(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) - \frac{2\eta}{1 - \gamma} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) \right] \\ & \leq \sum_{l=0}^t \lambda_{l+1,t+1} \mathbb{E}_{s \sim \rho} \left[5\eta \|Q^{(l+1)}(s) - Q_\tau^*(s)\|_\infty + \frac{4\eta^2}{1 - \gamma} \|Q^{(l+1)}(s) - Q^{(l)}(s)\|_\infty + \frac{13\eta^2}{1 - \gamma} \|Q^{(l-1)}(s) - Q^{(l)}(s)\|_\infty \right] \\ & \leq 5 \sum_{l=0}^t \lambda_{l+1,t+1} \mathbb{E}_{s \sim \rho} \left[\eta \|Q_\tau^*(s) - Q^{(l+1)}(s)\|_{\Gamma(\rho)} + \frac{12\eta^2}{(1 - \gamma)^2} \|Q^{(l+1)}(s) - Q^{(l)}(s)\|_{\Gamma(\rho)} \right] \end{aligned}$$

$$\leq \frac{31250\eta\mathcal{C}_\rho}{(1-\gamma)^3} \sum_{l=0}^{t-1} \lambda_{l+1,t+1} \text{KL}_\rho(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) + \frac{2750\eta}{(1-\gamma)^2} \lambda_{0,t+1}$$

for all $t \geq 0$. Rearranging terms, we have

$$\begin{aligned} & \alpha_{t+1} \left[\text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(t+1)}) + \left(1 - \frac{2\eta}{1-\gamma}\right) \text{KL}_\rho(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + 3\eta \mathbb{E}_{s \sim \rho} \left[\|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty \right] \right] \\ & + \sum_{l=1}^t (\lambda_{l,t+1} - (1-\eta\tau)\lambda_{l+1,t+1}) \text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(l)}) \\ & + \sum_{l=0}^{t-1} \left[\lambda_{l+1,t+1} \left(1 - \eta\tau - \frac{4\eta}{1-\gamma} - \frac{31250\eta\mathcal{C}_\rho}{(1-\gamma)^3}\right) - \lambda_{l+2,t+1} \frac{2\eta}{1-\gamma} \right] \text{KL}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \\ & \leq \frac{2750\eta}{(1-\gamma)^2} \lambda_{0,t+1} + (1-\eta\tau)\lambda_{1,t+1} \text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(0)}) \leq \left(\frac{2750\eta}{(1-\gamma)^2} + \eta\right) \lambda_{0,t+1}. \end{aligned}$$

With $0 < \eta \leq \frac{(1-\gamma)^3}{32000\mathcal{C}_\rho}$, and $\alpha_i = \eta\tau$, we have $\lambda_{l,t+1} - (1-\eta\tau)\lambda_{l+1,t+1} \geq 0$ (c.f. (40)), and

$$\begin{aligned} & \lambda_{l+1,t+1} \left(1 - \eta\tau - \frac{4\eta}{1-\gamma} - \frac{31250\eta\mathcal{C}_\rho}{(1-\gamma)^3}\right) - \lambda_{l+2,t+1} \frac{2\eta}{1-\gamma} \\ & = \eta\tau \prod_{j=l+3}^{t+1} \left(1 - \frac{1-\gamma}{4} \alpha_j\right) \left[\left(1 - \frac{1-\gamma}{4} \eta\tau\right) \left(1 - \eta\tau - \frac{4\eta}{1-\gamma} - \frac{31250\eta\mathcal{C}_\rho}{(1-\gamma)^3}\right) - \frac{2\eta}{1-\gamma} \right] \geq 0. \end{aligned}$$

It follows that

$$\begin{aligned} & \text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(t+1)}) + \left(1 - \frac{2\eta}{1-\gamma}\right) \text{KL}_\rho(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + 3\eta \mathbb{E}_{s \sim \rho} \left[\|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty \right] \\ & \leq \left(\frac{2750}{(1-\gamma)^2\tau} + \frac{1}{\tau}\right) \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^{t+1} < \frac{3000}{(1-\gamma)^2\tau} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^{t+1}. \end{aligned} \quad (17)$$

This proves the bound of $\text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(t+1)})$ and $3\eta \mathbb{E}_{s \sim \rho} [\|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty]$ in Theorem 3. Note that the bound holds trivially for $\text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(0)})$ and $3\eta \mathbb{E}_{s \sim \rho} [\|Q^{(0)}(s) - Q_\tau^*(s)\|_\infty]$. It remains to bound $\text{KL}(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)})$.

Lemma 5. *With $0 < \eta \leq (1-\gamma)/8$, we have*

$$\begin{aligned} & \frac{1}{2} \text{KL}_s(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) + \eta\tau \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\ & \leq (1-\eta\tau) \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t)}) + \frac{2\eta}{1-\gamma} \text{KL}_s(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) + 2\eta \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty. \end{aligned}$$

Proof. See Appendix C.5. □

Combining the above Lemma with (17) gives

$$\begin{aligned} & \frac{1}{2} \text{KL}_\rho(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) + \eta\tau \text{KL}_\rho(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\ & \leq (1-\eta\tau) \left(\text{KL}_\rho(\zeta_\tau^* \parallel \zeta^{(t)}) + \left(1 - \frac{2\eta}{1-\gamma}\right) \text{KL}_\rho(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) + 3\eta \mathbb{E}_{s \sim \rho} \left[\|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty \right] \right) \\ & \leq \frac{3000}{(1-\gamma)^2\tau} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^{t+1}. \end{aligned}$$

We are now ready to bound the duality gap. Before proceeding, we introduce the following two lemmas:

Lemma 6. *It holds for any policy pair μ, ν that*

$$\max_{\mu', \nu'} \left(V_{\tau}^{\mu', \nu}(\rho) - V_{\tau}^{\mu, \nu'}(\rho) \right) \leq \frac{2\mathcal{C}_{\rho, \tau}^{\dagger}}{1 - \gamma} \mathbb{E} \left[\max_{\mu', \nu'} \left(f_s(Q_{\tau}^*, \mu', \nu) - f_s(Q_{\tau}^*, \mu, \nu') \right) \right] \quad (18)$$

and

$$\max_{s \in \mathcal{S}, \mu', \nu'} \left(V_{\tau}^{\mu', \nu}(s) - V_{\tau}^{\mu, \nu'}(s) \right) \leq \frac{2\|1/\rho\|_{\infty}}{1 - \gamma} \mathbb{E} \left[\max_{\mu', \nu'} \left(f_s(Q_{\tau}^*, \mu', \nu) - f_s(Q_{\tau}^*, \mu, \nu') \right) \right]. \quad (19)$$

Proof. Note that (19) is a slight generalization of (Wei et al., 2021b, Lemma 32). The proof can be found in Appendix C.6. \square

Lemma 7 ((Cen et al., 2021b, Lemma 4)). *It holds for all $s \in \mathcal{S}$ and policy pair μ, ν that*

$$\max_{\mu', \nu'} \left(f_s(Q_{\tau}^*, \mu', \nu) - f_s(Q_{\tau}^*, \mu, \nu') \right) \leq \frac{4}{(1 - \gamma)^{2\tau}} \text{KL}_s(\zeta_{\tau}^* \parallel \zeta) + \tau \text{KL}_s(\zeta \parallel \zeta_{\tau}^*).$$

Putting all pieces together, we arrive at

$$\begin{aligned} \max_{\mu, \nu} \left(V_{\tau}^{\mu, \bar{\nu}^{(t)}}(\rho) - V_{\tau}^{\bar{\mu}^{(t)}, \nu}(\rho) \right) &\leq \frac{2\mathcal{C}_{\rho, \tau}^{\dagger}}{1 - \gamma} \left(\frac{4}{(1 - \gamma)^{2\tau}} \text{KL}_{\rho}(\zeta_{\tau}^* \parallel \bar{\zeta}^{(t+1)}) + \tau \text{KL}_{\rho}(\bar{\zeta}^{(t+1)} \parallel \zeta_{\tau}^*) \right) \\ &\leq \frac{2\mathcal{C}_{\rho, \tau}^{\dagger}}{1 - \gamma} \max \left\{ \frac{8}{(1 - \gamma)^{2\tau}}, \frac{1}{\eta} \right\} \left(\frac{1}{2} \text{KL}_{\rho}(\zeta_{\tau}^* \parallel \bar{\zeta}^{(t+1)}) + \eta \tau \text{KL}_{\rho}(\bar{\zeta}^{(t+1)} \parallel \zeta_{\tau}^*) \right) \\ &\leq \frac{6000\mathcal{C}_{\rho, \tau}^{\dagger}}{(1 - \gamma)^{3\tau}} \max \left\{ \frac{8}{(1 - \gamma)^{2\tau}}, \frac{1}{\eta} \right\} \left(1 - \frac{(1 - \gamma)\eta\tau}{4} \right)^t. \end{aligned}$$

We omit the proof for $\max_{s \in \mathcal{S}, \mu, \nu} \left(V_{\tau}^{\mu, \bar{\nu}^{(t)}}(s) - V_{\tau}^{\bar{\mu}^{(t)}, \nu}(s) \right)$ as it follows virtually the same argument.

B ANALYSIS FOR THE EPISODIC SETTING

Throughout the analysis, we restrict our choice of value update step size to $\alpha_t = \eta\tau$. We start with the following lemma which parallels Lemma 11 in the episodic Markov game setting:

Lemma 8. *With $0 < \eta \leq 1/\tau$, it holds for all $s \in \mathcal{S}$, $h \in [H]$ and $t \geq 0$ that*

$$\max \left\{ \|\bar{\mu}_h^{(t+1)}(s) - \mu_h^{(t+1)}(s)\|_1, \|\bar{\nu}_h^{(t+1)}(s) - \nu_h^{(t+1)}(s)\|_1 \right\} \leq 2\eta H. \quad (22)$$

In addition, we have

$$\max \left\{ \|\log \zeta_h^{(t)}(s)\|_{\infty}, \|\log \bar{\zeta}_h^{(t)}(s)\|_{\infty}, \|\log \zeta_{h, \tau}^*(s)\|_{\infty} \right\} \leq \frac{2H}{\tau} \quad (23)$$

Lemma 9. *With $0 < \eta \leq \frac{1}{8H}$, it holds for all $0 \leq t_1 \leq t_2$, $h \in [H]$ and $s \in \mathcal{S}$ that*

$$\begin{aligned} &\text{KL}_s(\zeta_{h, \tau}^* \parallel \zeta_h^{(t_2)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t_2)} \parallel \bar{\zeta}_h^{(t_2)}) \\ &\leq (1 - \eta\tau)^{t_2 - t_1} \left(\text{KL}_s(\zeta_{h, \tau}^* \parallel \zeta_h^{(t_1)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t_1)} \parallel \bar{\zeta}_h^{(t_1)}) \right) \\ &\quad + 4\eta \sum_{l=t_1}^{t_2} (1 - \eta\tau)^{t_2 - l} \|Q_h^{(l)}(s) - Q_{\tau}^*(s)\|_{\infty}. \end{aligned}$$

Proof. See Appendix D.1. \square

Lemma 10. *With $0 < \eta \leq \frac{1}{8H}$, it holds for all $0 < t_1 \leq t_2$, $2 \leq h \leq H$ and $s \in \mathcal{S}$ that*

$$\begin{aligned} &|Q_{h-1}^{(t_2)}(s, a, b) - Q_{h-1, \tau}^*(s, a, b)| \\ &\leq 2(1 - \eta\tau)^{t_2 - t_1} H + 10\eta\tau \mathbb{E}_{s' \sim P_{h-1}(\cdot | s, a, b)} \left[\sum_{l=t_1-1}^{t_2-1} (1 - \eta\tau)^{t_2-1-l} \|Q_h^{(l)}(s) - Q_{h, \tau}^*(s)\|_{\infty} \right] \\ &\quad + \tau(1 - \eta\tau)^{t_2 - t_1} \mathbb{E}_{s' \sim P_{h-1}(\cdot | s, a, b)} \left[\text{KL}_s(\zeta_{h, \tau}^* \parallel \zeta_h^{(t_1-1)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t_1-1)} \parallel \bar{\zeta}_h^{(t_1-1)}) \right]. \end{aligned}$$

Algorithm 2: Entropy-regularized OMWU for Episodic Two-player Zero-sum Markov Game

- 1 **Input:** Regularization parameter $\tau > 0$, learning rate for policy update $\eta > 0$, learning rate for value update $\{\alpha_t\}_{t=1}^\infty$.
- 2 **Initialization:** Set $\mu^{(0)}$, $\bar{\mu}^{(0)}$, $\nu^{(0)}$ and $\bar{\nu}^{(0)}$ as uniform policies; set

$$Q^{(0)} = 0, \quad V^{(0)} = \tau(\log |\mathcal{A}| - \log |\mathcal{B}|).$$

for $t = 0, 1, \dots$ **do**

3 **for all** $h \in [H]$, $s \in \mathcal{S}$ **do in parallel**

4 When $t \geq 1$, update policy pair $\zeta_h^{(t)}(s)$ as:

$$\begin{cases} \mu_h^{(t)}(a|s) \propto \mu_h^{(t-1)}(a|s)^{1-\eta\tau} \exp(\eta[Q_h^{(t)}(s)\bar{\nu}_h^{(t)}(s)]_a) \\ \nu_h^{(t)}(b|s) \propto \nu_h^{(t-1)}(b|s)^{1-\eta\tau} \exp(-\eta[Q_h^{(t)}(s)^\top \bar{\mu}_h^{(t)}(s)]_b) \end{cases} \quad (20a)$$

5 Update policy pair $\bar{\zeta}_h^{(t+1)}(s)$ as:

$$\begin{cases} \bar{\mu}_h^{(t+1)}(a|s) \propto \mu_h^{(t)}(a|s)^{1-\eta\tau} \exp(\eta[Q_h^{(t)}(s)\bar{\nu}_h^{(t)}(s)]_a) \\ \bar{\nu}_h^{(t+1)}(b|s) \propto \nu_h^{(t)}(b|s)^{1-\eta\tau} \exp(-\eta[Q_h^{(t)}(s)^\top \bar{\mu}_h^{(t)}(s)]_b) \end{cases} \quad (20b)$$

6 Update $Q_h^{(t+1)}(s)$ and $V_h^{(t+1)}(s)$ as

$$\begin{cases} Q_h^{(t+1)}(s, a, b) = r_h(s, a, b) + \mathbb{E}_{s' \sim P_h(\cdot|s, a, b)} [V_{h+1}^{(t)}(s')] \\ V_h^{(t+1)}(s) = (1 - \alpha_{t+1})V_h^{(t)}(s) \\ \quad + \alpha_{t+1} [\bar{\mu}_h^{(t+1)}(s)^\top Q_h^{(t+1)}(s)\bar{\nu}_h^{(t+1)}(s) + \tau\mathcal{H}(\bar{\mu}_h^{(t+1)}(s)) - \tau\mathcal{H}(\bar{\nu}_h^{(t+1)}(s))] \end{cases} \quad (21)$$

Proof. See Appendix D.2. □

We prove Theorem 2 by induction. By definition, we have

$$\|Q_{H,\tau}^* - Q_H^{(0)}\|_\infty = \|Q_{H,\tau}^*\|_\infty \leq 1,$$

and $\|Q_{H,\tau}^* - Q_H^{(t)}\|_\infty = \|r_H - r_H\|_\infty = 0$ for $t > 0$. So (13a) holds trivially for $h = H$. When the statement holds for some h , we can invoke Lemma 10 with $t_1 = T_h + 1$ and $t_2 = t \geq T_{h-1}$, which yields

$$\begin{aligned} & \|Q_{h-1}^{(t)} - Q_{h-1,\tau}^*\| \\ & \leq 2(1 - \eta\tau)^{t-T_{h-1}}H + 10\eta\tau \mathbb{E}_{s' \sim P(\cdot|s, a, b)} \left[\sum_{l=T_h}^{t-1} (1 - \eta\tau)^{t-1-l} \|Q_h^{(l)}(s) - Q_{h,\tau}^*(s)\|_\infty \right] \\ & \quad + \tau(1 - \eta\tau)^{t-T_{h-1}} \mathbb{E}_{s' \sim P(\cdot|s, a, b)} \left[\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(T_h)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(T_h)} \parallel \bar{\zeta}_h^{(T_h)}) \right] \\ & \leq 2(1 - \eta\tau)^{t-T_{h-1}}H + 10\eta\tau \mathbb{E}_{s' \sim P(\cdot|s, a, b)} \left[\sum_{l=T_h}^{t-1} (1 - \eta\tau)^{t-T_{h-1}-l} l^{H-h} \right] \\ & \quad + \tau(1 - \eta\tau)^{t-T_{h-1}} \mathbb{E}_{s' \sim P(\cdot|s, a, b)} \left[\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(T_h)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(T_h)} \parallel \bar{\zeta}_h^{(T_h)}) \right] \\ & \leq (1 - \eta\tau)^{t-T_{h-1}} (1 - \eta\tau)^{T_{\text{start}}-1} \left[10H + 10\eta\tau t^{H-h+1} \right], \end{aligned}$$

where the last step results from

$$\tau \left(\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(T_h)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(T_h)} \parallel \bar{\zeta}_h^{(T_h)}) \right)$$

$$\begin{aligned}
&\leq \tau \left(\left\| \log \mu_{h,\tau}^*(s) - \log \mu_h^{(T_h)}(s) \right\|_\infty + \left\| \log \nu_{h,\tau}^*(s) - \log \nu_h^{(T_h)}(s) \right\|_\infty \right. \\
&\quad \left. + \left\| \log \mu_h^{(T_h)}(s) - \log \bar{\mu}_h^{(T_h)}(s) \right\|_\infty + \left\| \log \nu_h^{(T_h)}(s) - \log \bar{\nu}_h^{(T_h)}(s) \right\|_\infty \right) \\
&\leq 8H.
\end{aligned}$$

Therefore, with $T_{\text{start}} = \lceil \frac{1}{\eta\tau} \log H \rceil$ we can guarantee that

$$\begin{aligned}
\left\| Q_{h-1}^{(t)} - Q_{h-1,\tau}^* \right\| &\leq 10(1-\eta\tau)^{t-T_{h-1}}(1-\eta\tau)^{T_{\text{start}}-1} \left[H + \eta\tau t^{H-h+1} \right] \\
&\leq (1-\eta\tau)^{t-T_{h-1}} t^{H-h+1}.
\end{aligned}$$

This completes the proof for (13a). Regarding (13b), we start by the following lemmas, which are simply Lemma 5 and Lemma 7 applied to the episodic setting:

Lemma 5A. *With $0 < \eta \leq \frac{1}{8H}$, we have*

$$\begin{aligned}
&\frac{1}{2} \text{KL}_s(\zeta_{h,\tau}^* \parallel \bar{\zeta}_h^{(t+1)}) + \eta\tau \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_{h,\tau}^*) \\
&\leq (1-\eta\tau) \text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t)}) + 2\eta H \text{KL}_s(\zeta_h^{(t)} \parallel \bar{\zeta}_h^{(t)}) + 2\eta \left\| Q_h^{(t)}(s) - Q_{h,\tau}^*(s) \right\|_\infty.
\end{aligned}$$

Lemma 7A. *It holds for all $h \in [H]$, $s \in \mathcal{S}$ and policy pair μ, ν that*

$$\max_{\mu', \nu'} (f_s(Q_{h,\tau}^*, \mu'_h, \nu_h) - f_s(Q_\tau^*, \mu_h, \nu'_h)) \leq \frac{4H^2}{\tau} \text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h) + \tau \text{KL}_s(\zeta_h \parallel \zeta_{h,\tau}^*).$$

Combining Lemma 9 with Lemma 5A and Lemma 7A, we conclude that for $0 \leq t_1 \leq t_2 - 1$,

$$\begin{aligned}
&\max_{\mu, \nu} (f_s(Q_{h,\tau}^*, \mu_h, \bar{\nu}_h^{(t_2)}) - f_s(Q_\tau^*, \bar{\mu}_h^{(t_2)}, \nu_h)) \\
&\leq \frac{4H^2}{\tau} \text{KL}_s(\zeta_{h,\tau}^* \parallel \bar{\zeta}_h^{(t_2)}) + \tau \text{KL}_s(\bar{\zeta}_h^{(t_2)} \parallel \zeta_{h,\tau}^*) \\
&\leq \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{1}{2} \text{KL}_s(\zeta_{h,\tau}^* \parallel \bar{\zeta}_h^{(t_2)}) + \eta\tau \text{KL}_s(\bar{\zeta}_h^{(t_2)} \parallel \zeta_{h,\tau}^*) \right) \\
&\leq \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left((1-\eta\tau) \text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t_2-1)}) + 2\eta H \text{KL}_s(\zeta_h^{(t_2-1)} \parallel \bar{\zeta}_h^{(t_2-1)}) + 2\eta \left\| Q_h^{(t_2-1)}(s) - Q_{h,\tau}^*(s) \right\|_\infty \right) \\
&\leq \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left((1-\eta\tau)^{t_2-t_1} \left(\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t_1)}) + (1-4\eta H) \text{KL}_s(\zeta_h^{(t_1)} \parallel \bar{\zeta}_h^{(t_1)}) \right) \right. \\
&\quad \left. + 6\eta \sum_{l=t_1}^{t_2} (1-\eta\tau)^{t_2-l} \left\| Q_h^{(l)}(s) - Q_{h,\tau}^*(s) \right\|_\infty \right).
\end{aligned}$$

It is straightforward to verify that the above inequality holds for $0 \leq t_1 \leq t_2$, by omitting the third step. Substitution of (13a) into the above inequality yields

$$\begin{aligned}
&\max_{\mu, \nu} (f_s(Q_{h,\tau}^*, \mu_h, \bar{\nu}_h^{(t)}) - f_s(Q_\tau^*, \bar{\mu}_h^{(t)}, \nu_h)) \\
&\leq \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left((1-\eta\tau)^{t-T_h} \left(\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(T_h)}) + (1-4\eta H) \text{KL}_s(\zeta_h^{(T_h)} \parallel \bar{\zeta}_h^{(T_h)}) \right) \right. \\
&\quad \left. + 6\eta \sum_{l=T_h}^t (1-\eta\tau)^{t-l} (1-\eta\tau)^{l-T_h} t^{H-h} \right) \\
&\leq (1-\eta\tau)^{t-T_h} \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{8H}{\tau} + 6\eta t^{H-h+1} \right). \tag{24}
\end{aligned}$$

We prove the following results instead, where (13b) is a direct conclusion of (25) by summing the two inequalities.

$$\begin{cases} \max_{s \in \mathcal{S}, \mu} \left(V_{h,\tau}^{\mu, \bar{\nu}^{(t)}}(s) - V_{h,\tau}^*(s) \right) \leq 2(1-\eta\tau)^{t-T_h} \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{8H}{\tau} + 6\eta t^{H-h+1} \right) \\ \max_{s \in \mathcal{S}, \mu} \left(V_{h,\tau}^*(s) - V_{h,\tau}^{\bar{\mu}^{(t)}, \nu}(s) \right) \leq 2(1-\eta\tau)^{t-T_h} \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{8H}{\tau} + 6\eta t^{H-h+1} \right) \end{cases} \tag{25}$$

We prove by induction. Note that when $h = H$, we have $V_{H,\tau}^{\mu,\nu}(s) = f_s(r_H, \mu_H, \nu_H) = f_s(Q_{H,\tau}^*, \mu_H, \nu_H)$ and the claim holds by invoking (24). When the claim holds for some $2 \leq h \leq H$, we have

$$\begin{aligned}
& V_{h-1,\tau}^{\mu,\bar{\nu}^{(t)}}(s) - V_{h-1,\tau}^*(s) \\
&= \mu_{h-1}(s)^\top Q_{h-1,\tau}^{\mu,\bar{\nu}^{(t)}}(s) \bar{\nu}_{h-1}^{(t)}(s) + \tau \mathcal{H}(\mu_{h-1}(s)) - \tau \mathcal{H}(\bar{\nu}_{h-1}^{(t)}(s)) \\
&\quad - \mu_{h-1}^*(s)^\top Q_{h-1,\tau}^*(s) \nu_{h-1}^*(s) + \tau \mathcal{H}(\mu_{h-1}^*(s)) - \tau \mathcal{H}(\nu_{h-1}^*(s)) \\
&= f_s(Q_{h-1,\tau}^*, \mu_{h-1}, \bar{\nu}_{h-1}^{(t)}) - f_s(Q_{h-1,\tau}^*, \mu_{h-1}^*, \nu_{h-1}^*) + \mu_{h-1}(s)^\top (Q_{h-1,\tau}^{\mu,\bar{\nu}^{(t)}}(s) - Q_{h-1,\tau}^*(s)) \bar{\nu}_{h-1}^{(t)}(s) \\
&\leq f_s(Q_{h-1,\tau}^*, \mu_{h-1}, \bar{\nu}_{h-1}^{(t)}) - f_s(Q_{h-1,\tau}^*, \bar{\mu}_{h-1}^{(t)}, \nu_{h-1}^*) + \max_{s' \in \mathcal{S}} [V_{h,\tau}^{\mu,\bar{\nu}^{(t)}}(s') - V_{h,\tau}^*(s')] \\
&\leq \max_{\mu_{h-1}, \nu_{h-1}'} \left(f_s(Q_{h-1,\tau}^*, \mu_{h-1}', \bar{\nu}_{h-1}^{(t)}) - f_s(Q_{h-1,\tau}^*, \bar{\mu}_{h-1}^{(t)}, \nu_{h-1}') \right) + \max_{s' \in \mathcal{S}} [V_{h,\tau}^{\mu,\bar{\nu}^{(t)}}(s') - V_{h,\tau}^*(s')] \\
&\leq (1 - \eta\tau)^{t-T_{h-1}} \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{8H}{\tau} + 6\eta t^{H-h+2} \right) \\
&\quad + 2(1 - \eta\tau)^{t-T_h} \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{8H}{\tau} + 6\eta t^{H-h+1} \right) \\
&\leq 2(1 - \eta\tau)^{t-T_{h-1}} \max \left\{ \frac{8H^2}{\tau}, \frac{1}{\eta} \right\} \left(\frac{8H}{\tau} + 6\eta t^{H-h+2} \right).
\end{aligned}$$

Taking maximum over μ verifies the claim for $h - 1$, thereby finishing the proof. The bound for $\max_{s \in \mathcal{S}, \mu} (V_{h,\tau}^*(s) - V_{h,\tau}^{\bar{\mu}^{(t)}, \nu}(s))$ can be established by following a similar argument and is therefore omitted.

C PROOF OF KEY LEMMAS FOR THE DISCOUNTED SETTING

C.1 PROOF OF LEMMA 1

Before proceeding, we shall introduce the following lemma that quantifies the distance between two consecutive updates, whose proof can be found in Appendix E.1.

Lemma 11. *For $0 < \eta \leq 1/\tau$, it holds for all $s \in \mathcal{S}$ and $t \geq 0$ that*

$$\max \left\{ \|\bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s)\|_1, \|\bar{\nu}^{(t+1)}(s) - \nu^{(t+1)}(s)\|_1 \right\} \leq \frac{2\eta}{1-\gamma}$$

and that

$$\max \left\{ \|\bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s)\|_1, \|\bar{\nu}^{(t+1)}(s) - \bar{\nu}^{(t)}(s)\|_1 \right\} \leq \frac{6\eta}{1-\gamma}.$$

For notational simplicity, we use $x \stackrel{1}{=} y$ to denote equivalence up to a global shift for two vectors x, y :

$$x = y + c \cdot \mathbf{1}$$

for some constant $c \in \mathbb{R}$. Taking logarithm on the both sides of the update rule (9a), we get

$$\begin{cases} \log \mu^{(t+1)}(s) - (1 - \eta\tau) \log \mu^{(t)}(s) & \stackrel{1}{=} \eta Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s) \\ \log \nu^{(t+1)}(s) - (1 - \eta\tau) \log \nu^{(t)}(s) & \stackrel{1}{=} -\eta Q^{(t+1)}(s)^\top \bar{\mu}^{(t+1)}(s) \end{cases}. \quad (26)$$

On the other hand, it holds for the optimal policies (μ_τ^*, ν_τ^*) that

$$\begin{cases} \eta\tau \log \mu_\tau^*(s) & \stackrel{1}{=} \eta Q_\tau^*(s) \nu_\tau^*(s) \\ \eta\tau \log \nu_\tau^*(s) & \stackrel{1}{=} -\eta Q_\tau^*(s)^\top \mu_\tau^*(s) \end{cases}. \quad (27)$$

Subtracting (27) from (26) and taking inner product with $\bar{\zeta}^{(t+1)}(s) - \zeta_\tau^*(s)$ gives

$$\begin{aligned}
& \langle \log \zeta^{(t+1)}(s) - (1 - \eta\tau) \log \zeta^{(t)}(s) - \eta\tau \log \zeta_\tau^*(s), \bar{\zeta}^{(t+1)}(s) - \zeta_\tau^*(s) \rangle \\
&= \eta \langle \bar{\mu}^{(t+1)}(s) - \mu_\tau^*(s), Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s) - Q_\tau^*(s) \nu_\tau^*(s) \rangle \\
&\quad - \eta \langle \bar{\nu}^{(t+1)}(s) - \nu_\tau^*(s), Q^{(t+1)}(s)^\top \bar{\mu}^{(t+1)}(s) - Q_\tau^*(s)^\top \mu_\tau^*(s) \rangle \\
&= \eta \langle \bar{\mu}^{(t+1)}(s) - \mu_\tau^*(s), (Q^{(t+1)}(s) - Q_\tau^*(s)) \bar{\nu}^{(t+1)}(s) \rangle \\
&\quad - \eta \langle \bar{\nu}^{(t+1)}(s) - \nu_\tau^*(s), (Q^{(t+1)}(s) - Q_\tau^*(s))^\top \bar{\mu}^{(t+1)}(s) \rangle \\
&= -\eta \langle \mu_\tau^*(s), (Q^{(t+1)}(s) - Q_\tau^*(s)) \bar{\nu}^{(t+1)}(s) \rangle + \eta \langle \nu_\tau^*(s), (Q^{(t+1)}(s) - Q_\tau^*(s))^\top \bar{\mu}^{(t+1)}(s) \rangle \\
&\leq 2\eta \|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty.
\end{aligned} \tag{28}$$

We rewrite the LHS as

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

Rearranging terms, we have

$$\begin{aligned}
& \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t+1)}) - (1 - \eta\tau) \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t)}) + (1 - \eta\tau) \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) \\
&\quad + \eta\tau \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) + \text{KL}_s(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + (1 - \eta\tau) \text{KL}_s(\bar{\zeta}^{(t)} \parallel \zeta^{(t)}) \\
&\quad - \langle \log \bar{\zeta}^{(t+1)}(s) - \log \zeta^{(t+1)}(s), \bar{\zeta}^{(t+1)}(s) - \zeta^{(t+1)}(s) \rangle \\
&\quad - (1 - \eta\tau) \langle \log \zeta^{(t)}(s) - \log \bar{\zeta}^{(t)}(s), \bar{\zeta}^{(t+1)}(s) - \bar{\zeta}^{(t)}(s) \rangle \\
&\leq 2\eta \|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty.
\end{aligned}$$

It remains to bound $\langle \log \bar{\zeta}^{(t+1)}(s) - \log \zeta^{(t+1)}(s), \bar{\zeta}^{(t+1)}(s) - \zeta^{(t+1)}(s) \rangle$ and $\langle \log \zeta^{(t)}(s) - \log \bar{\zeta}^{(t)}(s), \bar{\zeta}^{(t+1)}(s) - \bar{\zeta}^{(t)}(s) \rangle$. Note that

$$\begin{aligned}
& \langle \log \bar{\mu}^{(t+1)}(s) - \log \mu^{(t+1)}(s), \bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s) \rangle \\
&= \eta \langle Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s), \bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s) \rangle \\
&\leq \eta \|Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s)\|_1 \|\bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s)\|_1.
\end{aligned} \tag{29}$$

We bound $\|Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s)\|_1$ as

$$\begin{aligned}
& \|Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s)\|_1 \\
&\leq \|Q^{(t+1)}(s) (\bar{\nu}^{(t)}(s) - \bar{\nu}^{(t+1)}(s))\|_1 + \|(Q^{(t)}(s) - Q^{(t+1)}(s)) \bar{\nu}^{(t)}(s)\|_1 \\
&\leq \frac{2}{1 - \gamma} \|\bar{\nu}^{(t)}(s) - \bar{\nu}^{(t+1)}(s)\|_1 + \|Q^{(t)}(s) - Q^{(t+1)}(s)\|_\infty.
\end{aligned}$$

Plugging the above inequality into (29) and invoking Young's inequality yields

$$\begin{aligned}
& \langle \log \bar{\mu}^{(t+1)}(s) - \log \mu^{(t+1)}(s), \bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s) \rangle \\
&\leq \frac{\eta}{1 - \gamma} \left(\|\bar{\nu}^{(t+1)}(s) - \bar{\nu}^{(t)}(s)\|_1^2 + \|\bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s)\|_1^2 \right)
\end{aligned}$$

$$\begin{aligned}
& + \eta \|Q^{(t)}(s) - Q^{(t+1)}(s)\|_\infty \|\bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s)\|_1 \\
& \leq \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\nu}^{(t+1)} \parallel \bar{\nu}^{(t)}) + \frac{2\eta}{1-\gamma} \text{KL}_s(\mu^{(t+1)} \parallel \bar{\mu}^{(t+1)}) + \frac{2\eta^2}{1-\gamma} \|Q^{(t)}(s) - Q^{(t+1)}(s)\|_\infty,
\end{aligned}$$

where the last step results from Pinsker's inequality and Lemma 11. Similarly, we have

$$\begin{aligned}
& \langle \log \bar{\nu}^{(t+1)}(s) - \log \nu^{(t+1)}(s), \bar{\nu}^{(t+1)}(s) - \nu^{(t+1)}(s) \rangle \\
& \leq \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\mu}^{(t+1)} \parallel \bar{\mu}^{(t)}) + \frac{2\eta}{1-\gamma} \text{KL}_s(\nu^{(t+1)} \parallel \bar{\nu}^{(t+1)}) + \frac{2\eta^2}{1-\gamma} \|Q^{(t)}(s) - Q^{(t+1)}(s)\|_\infty.
\end{aligned}$$

Combining the above two inequalities gives

$$\begin{aligned}
& \langle \log \bar{\zeta}^{(t+1)}(s) - \log \zeta^{(t+1)}(s), \bar{\zeta}^{(t+1)}(s) - \zeta^{(t+1)}(s) \rangle \\
& \leq \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) + \frac{2\eta}{1-\gamma} \text{KL}_s(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + \frac{4\eta^2}{1-\gamma} \|Q^{(t)}(s) - Q^{(t+1)}(s)\|_\infty.
\end{aligned}$$

By a similar argument, when $t \geq 1$:

$$\begin{aligned}
& \langle \log \zeta^{(t)}(s) - \log \bar{\zeta}^{(t)}(s), \bar{\zeta}^{(t+1)}(s) - \bar{\zeta}^{(t)}(s) \rangle \\
& = \eta \langle Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q^{(t-1)}(s) \bar{\nu}^{(t-1)}(s), \bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s) \rangle \\
& \quad - \eta \langle Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) - Q^{(t-1)}(s)^\top \bar{\mu}^{(t-1)}(s), \bar{\nu}^{(t+1)}(s) - \bar{\nu}^{(t)}(s) \rangle \\
& \leq \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t)} \parallel \bar{\zeta}^{(t-1)}) + \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) \\
& \quad + \eta (\|\bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s)\|_1 + \|\bar{\nu}^{(t+1)}(s) - \bar{\nu}^{(t)}(s)\|_1) \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty \\
& \leq \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t)} \parallel \bar{\zeta}^{(t-1)}) + \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) \\
& \quad + \frac{12\eta^2}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty.
\end{aligned}$$

Note that the above inequality trivially holds for $t = 0$, since $\log \zeta^{(0)}(s) = \log \bar{\zeta}^{(0)}(s)$.

Putting pieces together, we conclude for that

$$\begin{aligned}
& \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t+1)}) - (1 - \eta\tau) \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t)}) + \left(1 - \eta\tau - \frac{4\eta}{1-\gamma}\right) \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) \\
& \quad + \eta\tau \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) + \left(1 - \frac{2\eta}{1-\gamma}\right) \text{KL}_s(\zeta^{(t+1)} \parallel \bar{\zeta}^{(t+1)}) + (1 - \eta\tau) \text{KL}_s(\bar{\zeta}^{(t)} \parallel \zeta^{(t)}) \\
& \quad - \frac{2\eta}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t)} \parallel \bar{\zeta}^{(t-1)}) \\
& \leq 2\eta \|Q^{(t+1)}(s) - Q_\tau^*(s)\|_\infty + \frac{4\eta^2}{1-\gamma} \|Q^{(t)}(s) - Q^{(t+1)}(s)\|_\infty + \frac{12\eta^2}{1-\gamma} \|Q^{(t-1)}(s) - Q^{(t)}(s)\|_\infty.
\end{aligned}$$

Averaging s over the distribution ρ completes the proof.

C.2 PROOF OF LEMMA 2

Proof. By definition of Q , it holds for $t \geq 1$ that

$$|Q^{(t+1)}(s, a, b) - Q^{(t)}(s, a, b)| \leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} \left[|V^{(t)}(s') - V^{(t-1)}(s')| \right]. \quad (30)$$

We denote by $f_s(Q, \mu, \nu)$ the one-step entropy-regularized game value at state s , i.e.,

$$f_s(Q, \mu, \nu) = \mu(s)^\top Q(s) \nu(s) + \tau \mathcal{H}(\mu(s)) - \tau \mathcal{H}(\nu(s)).$$

We further simplify the notation by introducing

$$f_s^{(t)} = f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)})$$

By recursively applying the update rule $V^{(t)}(s) = (1 - \alpha_t)V^{(t-1)}(s) + \alpha_t f_s^{(t)}$, we get

$$V^{(t)}(s) = \alpha_{0,t}V^{(0)} + \sum_{l=1}^t \alpha_{l,t}f_s(Q^{(l)}, \bar{\mu}^{(l)}, \bar{\nu}^{(l)}) = \sum_{l=0}^t \alpha_{l,t}f_s^{(l)}.$$

Since $\alpha_0 = 1$, it follows that

$$\sum_{l=0}^t \alpha_{l,t} = \alpha_0 = 1$$

So we have

$$\begin{aligned} |V^{(t)}(s) - V^{(t-1)}(s)| &= \alpha_t |f_s^{(t)} - V^{(t-1)}(s)| \\ &= \alpha_t \sum_{l=0}^{t-1} \alpha_{l,t-1} |f_s^{(t)} - f_s^{(l)}| \\ &\leq \alpha_t \sum_{l=0}^{t-1} \alpha_{l,t-1} \sum_{j=l}^{t-1} |f_s^{(j+1)} - f_s^{(j)}| \end{aligned} \quad (31)$$

The next lemma enables us to upper bound $|f_s^{(t+1)} - f_s^{(t)}|$ with $\|Q^{(t+1)}(s) - Q^{(t)}(s)\|_\infty$ and $\text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)})$ (as well as their $(t-1)$ -th iteration counter parts). The proof is postponed to Appendix E.2.

Lemma 12. For any $t \geq 0$, $\eta \leq (1 - \gamma)/180$, we have

$$\begin{aligned} |f_s^{(t+1)} - f_s^{(t)}| &\leq \|Q^{(t+1)}(s) - Q^{(t)}(s)\|_\infty + \left(\frac{3}{\eta} + \frac{4}{1-\gamma}\right) \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) \\ &\quad + \frac{12\eta}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty + \frac{2}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t)} \parallel \bar{\zeta}^{(t-1)}). \end{aligned}$$

Plugging the above lemma into (31),

$$\begin{aligned} &|V^{(t)}(s) - V^{(t-1)}(s)| \\ &\leq \alpha_t \sum_{l=0}^{t-1} \alpha_{l,t-1} \sum_{j=l}^{t-1} \left[\|Q^{(j+1)}(s) - Q^{(j)}(s)\|_\infty + \left(\frac{3}{\eta} + \frac{4}{1-\gamma}\right) \text{KL}_s(\bar{\zeta}^{(j+1)} \parallel \bar{\zeta}^{(j)}) \right] \\ &\quad + \alpha_t \sum_{l=0}^{t-1} \alpha_{l,t-1} \sum_{j=l}^{t-1} \left[\frac{12\eta}{1-\gamma} \|Q^{(j)}(s) - Q^{(j-1)}(s)\|_\infty + \frac{2}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(j)} \parallel \bar{\zeta}^{(j-1)}) \right] \\ &\leq \alpha_t \sum_{l=0}^{t-1} \alpha_{l,t-1} \sum_{j=l}^{t-1} \left[\left(1 + \frac{12\eta}{1-\gamma}\right) \|Q^{(j+1)}(s) - Q^{(j)}(s)\|_\infty + \left(\frac{3}{\eta} + \frac{6}{1-\gamma}\right) \text{KL}_s(\bar{\zeta}^{(j+1)} \parallel \bar{\zeta}^{(j)}) \right] \\ &\quad + \alpha_t \sum_{l=0}^{t-1} \alpha_{l,t-1} \left[\frac{12\eta}{1-\gamma} \|Q^{(l)}(s) - Q^{(l-1)}(s)\|_\infty + \frac{2}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) \right] \\ &\leq \sum_{j=0}^{t-1} \alpha_{j+1} \sum_{l=0}^j \alpha_{l,t-1} \left[\left(1 + \frac{12\eta}{1-\gamma}\right) \|Q^{(j+1)}(s) - Q^{(j)}(s)\|_\infty + \left(\frac{3}{\eta} + \frac{6}{1-\gamma}\right) \text{KL}_s(\bar{\zeta}^{(j+1)} \parallel \bar{\zeta}^{(j)}) \right] \\ &\quad + \alpha_t \sum_{l=0}^{t-2} \alpha_{l+1,t-1} \left[\frac{12\eta}{1-\gamma} \|Q^{(l+1)}(s) - Q^{(l)}(s)\|_\infty + \frac{2}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right], \end{aligned}$$

where the last step is due to $\alpha_t \leq \alpha_j$ for all $j \leq t$. To continue, by definition of α we have $\alpha_t \alpha_{l+1,t-1} \leq \alpha_{l+1,t-1}(1 - \alpha_t) = \alpha_{l+1,t}$ for $0 \leq l < t$, and that

$$\alpha_{j+1} \sum_{l=0}^j \alpha_{l,t-1} = \alpha_{j+1} \sum_{l=0}^j \left(\prod_{i=l+1}^{t-1} (1 - \alpha_i) - \prod_{i=l}^{t-1} (1 - \alpha_i) \right)$$

$$\begin{aligned}
&= \alpha_{j+1} \prod_{i=j+1}^{t-1} (1 - \alpha_i) \\
&\leq \alpha_{j+1} \prod_{i=j+2}^t (1 - \alpha_i) = \alpha_{j+1,t}.
\end{aligned}$$

Plugging into the inequality above gives

$$\begin{aligned}
&|V^{(t)}(s) - V^{(t-1)}(s)| \\
&\leq \sum_{j=0}^{t-1} \alpha_{j+1,t} \left[\left(1 + \frac{12\eta}{1-\gamma}\right) \|Q^{(j+1)}(s) - Q^{(j)}(s)\|_{\infty} + \left(\frac{3}{\eta} + \frac{6}{1-\gamma}\right) \text{KL}_s(\bar{\zeta}^{(j+1)} \parallel \bar{\zeta}^{(j)}) \right] \\
&\quad + \sum_{l=0}^{t-2} \alpha_{l+1,t} \left[\frac{12\eta}{1-\gamma} \|Q^{(l+1)}(s) - Q^{(l)}(s)\|_{\infty} + \frac{2}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right] \\
&\leq \sum_{l=0}^{t-1} \alpha_{l+1,t} \left[\left(1 + \frac{24\eta}{1-\gamma}\right) \|Q^{(l+1)}(s) - Q^{(l)}(s)\|_{\infty} + \frac{4}{\eta} \text{KL}_s(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right].
\end{aligned}$$

Plugging the above inequality into (30) leads to

$$\begin{aligned}
&|Q^{(t+1)}(s, a, b) - Q^{(t)}(s, a, b)| \\
&\leq \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} \left\{ \sum_{l=0}^{t-1} \alpha_{l+1,t} \left[\left(1 + \frac{24\eta}{1-\gamma}\right) \|Q^{(l+1)}(s') - Q^{(l)}(s')\|_{\infty} + \frac{4}{\eta} \text{KL}_{s'}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right] \right\}.
\end{aligned}$$

When $\eta \leq \frac{(1-\gamma)^2}{48\gamma}$, we have $\gamma(1 + \frac{24\eta}{1-\gamma}) \leq \frac{1+\gamma}{2}$ and hence that

$$\begin{aligned}
&|Q^{(t+1)}(s, a, b) - Q^{(t)}(s, a, b)| \\
&\leq \mathbb{E}_{s' \sim P(\cdot | s, a, b)} \left\{ \frac{1+\gamma}{2} \sum_{l=0}^{t-1} \alpha_{l+1,t} \left[\|Q^{(l+1)}(s') - Q^{(l)}(s')\|_{\infty} + \frac{4}{\eta} \text{KL}_{s'}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right] \right\}.
\end{aligned}$$

Let $x^{(t+1)} \in \mathcal{A}^S$ and $y^{(t+1)} \in \mathcal{B}^S$ be defined as

$$(x^{(t+1)}(s), y^{(t+1)}(s)) = \arg \max_{(a,b) \in \mathcal{A} \times \mathcal{B}} |Q^{(t+1)}(s, a, b) - Q^{(t)}(s, a, b)|.$$

It follows that $\forall \chi \in \Gamma(\rho)$, we have $\chi P_{x^{(t+1)}, y^{(t+1)}} \in \Gamma(\rho)$ and hence

$$\begin{aligned}
&\mathbb{E}_{s \sim \chi} \left[\|Q^{(t+1)}(s) - Q^{(t)}(s)\|_{\infty} \right] \\
&= \mathbb{E}_{\substack{s \sim \chi, \\ a = x^{(t+1)}(s), \\ b = y^{(t+1)}(s)}} \left[|Q^{(t+1)}(s, a, b) - Q^{(t)}(s, a, b)| \right] \\
&\leq \mathbb{E}_{s' \sim \chi P_{x^{(t+1)}, y^{(t+1)}}} \left[\frac{1+\gamma}{2} \sum_{l=0}^{t-1} \alpha_{l+1,t} \left[\|Q^{(l+1)}(s') - Q^{(l)}(s')\|_{\infty} + \frac{4}{\eta} \text{KL}_{s'}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right] \right] \\
&\leq \frac{1+\gamma}{2} \sum_{l=0}^{t-1} \alpha_{l+1,t} \left[\|Q^{(l+1)}(s') - Q^{(l)}(s')\|_{\Gamma(\rho)} + \frac{4}{\eta} \cdot \left\| \frac{\chi P_{x^{(t+1)}, y^{(t+1)}}}{\rho} \right\|_{\infty} \text{KL}_{\rho}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right] \\
&\leq \frac{1+\gamma}{2} \sum_{l=0}^{t-1} \alpha_{l+1,t} \left[\|Q^{(l+1)}(s') - Q^{(l)}(s')\|_{\Gamma(\rho)} + \frac{4\mathcal{C}_{\rho}}{\eta} \text{KL}_{\rho}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) \right]. \tag{32}
\end{aligned}$$

Taking supremum over $\chi \in \Gamma(\rho)$ completes the proof.

When $t = 0$, we have $\|Q^{(0)}(s) - Q^{(1)}(s)\|_{\Gamma(\rho)} = \|Q^{(1)}(s)\|_{\Gamma(\rho)} \leq 2$. \square

C.3 PROOF OF LEMMA 3

Note that it suffices to show for $t \geq 0$, $s \in \mathcal{S}$, $(a, b) \in \mathcal{A} \times \mathcal{B}$:

$$\begin{aligned} & |Q^{(t+1)}(s, a, b) - Q_\tau^*(s, a, b)| \\ & \leq \frac{1+\gamma}{2} \cdot \mathbb{E}_{s' \sim P(s, a, b)} \left[\sum_{l=0}^t \alpha_{l,t} \left[\|Q^{(l)}(s') - Q_\tau^*(s')\|_\infty + \frac{2\eta}{1-\gamma} \|Q^{(l)}(s') - Q^{(l-1)}(s')\|_\infty \right] \right] + 2\alpha_{0,t}. \end{aligned} \quad (33)$$

The remaining step follows a similar argument as (32) and is therefore omitted.

For $t \geq 0$, we have

$$\begin{aligned} Q^{(t+1)}(s, a, b) - Q_\tau^*(s, a, b) &= \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} \left[V^{(t)}(s') - V_\tau^*(s') \right] \\ &= \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} \left[\sum_{l=0}^t \alpha_{l,t} (f_{s'}^{(l)} - f_{s'}^*) \right]. \end{aligned} \quad (34)$$

We start by decomposing $f_s^{(t)} - f_s^*$ as

$$\begin{aligned} f_s^{(t)} - f_s^* &= f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q_\tau^*, \mu_\tau^*, \nu_\tau^*) \\ &= \left(f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \nu_\tau^*) \right) + f_s(Q^{(t)}, \bar{\mu}^{(t)}, \nu_\tau^*) - f_s(Q_\tau^*, \mu_\tau^*, \nu_\tau^*) \\ &\leq \left(f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \nu_\tau^*) \right) + f_s(Q_\tau^*, \bar{\mu}^{(t)}, \nu_\tau^*) - f_s(Q_\tau^*, \mu_\tau^*, \nu_\tau^*) \\ &\quad + \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty \\ &\leq f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \nu_\tau^*) + \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty. \end{aligned}$$

We bound the first two terms with the following lemma:

Lemma 13. *It holds for all $t \geq 0$, $s \in \mathcal{S}$ and $\nu(s) \in \Delta(\mathcal{B})$ that*

$$\begin{aligned} & f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \nu) \\ &= \langle \bar{\nu}^{(t)}(s) - \nu(s), Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) \rangle - \tau \mathcal{H}(\bar{\nu}^{(t)}(s)) + \tau \mathcal{H}(\nu^*(s)) \\ &\leq \frac{2\eta}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty + \frac{2}{1-\gamma} \left(\text{KL}_s(\bar{\mu}^{(t)} \parallel \mu^{(t-1)}) + \text{KL}_s(\mu^{(t-1)} \parallel \bar{\mu}^{(t-1)}) \right) \\ &\quad - \frac{1}{\eta} \left(1 - \frac{4\eta}{1-\gamma} \right) \text{KL}_s(\nu^{(t)} \parallel \bar{\nu}^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\bar{\nu}^{(t)} \parallel \nu^{(t-1)}) \\ &\quad + \frac{1-\eta\tau}{\eta} \text{KL}_s(\nu \parallel \nu^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\nu \parallel \nu^{(t)}). \end{aligned}$$

Proof. See Appendix E.3. □

Applying Lemma 13 with $\nu(s) = \nu_\tau^*(s)$ gives

$$\begin{aligned} f_s^{(t)} - f_s^* &\leq \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty + \frac{2\eta}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty \\ &\quad + \frac{1-\eta\tau}{\eta} \text{KL}_s(\nu_\tau^* \parallel \nu^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\nu_\tau^* \parallel \nu^{(t)}) \\ &\quad - \frac{1}{\eta} \left(1 - \frac{4\eta}{1-\gamma} \right) \text{KL}_s(\nu^{(t)} \parallel \bar{\nu}^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\bar{\nu}^{(t)} \parallel \nu^{(t-1)}) \\ &\quad + \frac{2}{1-\gamma} \left(\text{KL}_s(\bar{\mu}^{(t)} \parallel \mu^{(t-1)}) + \text{KL}_s(\mu^{(t-1)} \parallel \bar{\mu}^{(t-1)}) \right) \end{aligned} \quad (35)$$

By a similar argument, we can derive

$$\begin{aligned}
f_s^* - f_s^{(t)} &\leq \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty + \frac{2\eta}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty \\
&\quad + \frac{1-\eta\tau}{\eta} \text{KL}_s(\mu_\tau^* \parallel \mu^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\mu_\tau^* \parallel \mu^{(t)}) \\
&\quad - \frac{1}{\eta} \left(1 - \frac{4\eta}{1-\gamma}\right) \text{KL}_s(\mu^{(t)} \parallel \bar{\mu}^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\bar{\mu}^{(t)} \parallel \mu^{(t-1)}) \\
&\quad + \frac{2}{1-\gamma} \left(\text{KL}_s(\bar{\nu}^{(t)} \parallel \nu^{(t-1)}) + \text{KL}_s(\nu^{(t-1)} \parallel \bar{\nu}^{(t-1)}) \right).
\end{aligned} \tag{36}$$

Computing (35) + $\frac{1-\gamma}{4}$ · (36) gives

$$\begin{aligned}
&\left(1 - \frac{1-\gamma}{4}\right)(f_s^{(t)} - f_s^*) \\
&\leq \left(1 + \frac{1-\gamma}{4}\right) \left[\|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty + \frac{2\eta}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty \right] \\
&\quad + \frac{1-\eta\tau}{\eta} \left[\text{KL}_s(\nu_\tau^* \parallel \nu^{(t-1)}) + \frac{1-\gamma}{4} \text{KL}_s(\mu_\tau^* \parallel \mu^{(t-1)}) \right] - \frac{1}{\eta} \left[\text{KL}_s(\nu_\tau^* \parallel \nu^{(t)}) + \frac{1-\gamma}{4} \text{KL}_s(\mu_\tau^* \parallel \mu^{(t)}) \right] \\
&\quad + \frac{2}{1-\gamma} \left[\text{KL}_s(\mu^{(t-1)} \parallel \bar{\mu}^{(t-1)}) + \frac{1-\gamma}{4} \text{KL}_s(\nu^{(t-1)} \parallel \bar{\nu}^{(t-1)}) \right] \\
&\quad - \frac{1}{\eta} \left(1 - \frac{4\eta}{1-\gamma}\right) \left[\frac{1-\gamma}{4} \text{KL}_s(\mu^{(t)} \parallel \bar{\mu}^{(t)}) + \text{KL}_s(\nu^{(t)} \parallel \bar{\nu}^{(t)}) \right] \\
&\quad + \left(\frac{2}{1-\gamma} - \frac{1-\eta\tau}{\eta} \cdot \frac{1-\gamma}{4} \right) \text{KL}_s(\bar{\mu}^{(t)} \parallel \mu^{(t-1)}) + \left(\frac{2}{1-\gamma} \cdot \frac{1-\gamma}{4} - \frac{1-\eta\tau}{\eta} \right) \text{KL}_s(\bar{\nu}^{(t)} \parallel \nu^{(t-1)}).
\end{aligned} \tag{37}$$

With $0 < \eta \leq (1-\gamma)^2/16$, we have $\left(\frac{2}{1-\gamma} - \frac{1-\eta\tau}{\eta} \cdot \frac{1-\gamma}{4}\right) \leq 0$, $\left(\frac{2}{1-\gamma} \cdot \frac{1-\gamma}{4} - \frac{1-\eta\tau}{\eta}\right) \leq 0$, and

$$\frac{1}{\eta} \left(1 - \frac{4\eta}{1-\gamma}\right) \cdot \frac{1-\gamma}{4} \geq \frac{2}{1-\gamma} \cdot \frac{1}{1-\eta\tau}.$$

To proceed, we introduce a shorthand notation

$$\begin{aligned}
G^{(t)}(s) &= \frac{1}{\eta} \left[\text{KL}_s(\nu_\tau^* \parallel \nu^{(t)}) + \frac{1-\gamma}{4} \text{KL}_s(\mu_\tau^* \parallel \mu^{(t)}) \right] \\
&\quad + \frac{2}{(1-\gamma)(1-\eta\tau)} \left[\text{KL}_s(\mu^{(t)} \parallel \bar{\mu}^{(t)}) + \text{KL}_s(\nu^{(t)} \parallel \bar{\nu}^{(t)}) \right].
\end{aligned}$$

We can then write (37) as

$$\begin{aligned}
\left(1 - \frac{1-\gamma}{4}\right)(f_s^{(t)} - f_s^*) &\leq \left(1 + \frac{1-\gamma}{4}\right) \left[\|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty + \frac{2\eta}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty \right] \\
&\quad + (1-\eta\tau)G^{(t-1)}(s) - G^{(t)}(s).
\end{aligned} \tag{38}$$

Note that when $t = 0$, we have

$$\begin{aligned}
f_s^{(0)} - f_s^* &= \tau \log |\mathcal{A}| - \tau \log |\mathcal{B}| - \mu_\tau^*(s)^\top Q_\tau^*(s) \nu_\tau^*(s) - \tau \mathcal{H}(\mu_\tau^*(s)) + \tau \mathcal{H}(\nu_\tau^*(s)) \\
&= \max_{\mu(s)} \min_{\nu(s)} f_s(Q^{(0)}, \mu, \nu) - \max_{\mu(s)} \min_{\nu(s)} f_s(Q_\tau^*, \mu, \nu) \\
&\leq \|Q^{(0)}(s) - Q_\tau^*(s)\|_\infty.
\end{aligned} \tag{39}$$

Substitution of (38) and (39) into (34) gives

$$\begin{aligned}
&Q^{(t+1)}(s, a, b) - Q_\tau^*(s, a, b) \\
&= \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} \left[\sum_{l=0}^t \alpha_{l,t} (f_{s'}^{(l)} - f_{s'}^*) \right] \\
&\leq \gamma \mathbb{E}_{s' \sim P(s, a, b)} \left[\alpha_{0,t} \|Q^{(0)}(s') - Q_\tau^*(s')\|_\infty \right]
\end{aligned}$$

$$\begin{aligned}
& + \gamma \cdot \frac{1 + (1 - \gamma)/4}{1 - (1 - \gamma)/4} \mathbb{E}_{s' \sim P(s, a, b)} \left[\sum_{l=1}^t \alpha_{l,t} \left[\|Q^{(l)}(s') - Q_\tau^*(s')\|_\infty + \frac{2\eta}{1 - \gamma} \|Q^{(l)}(s') - Q^{(l-1)}(s')\|_\infty \right] \right] \\
& + \frac{\gamma}{1 - (1 - \gamma)/4} \mathbb{E}_{s' \sim P(s, a, b)} \left[(1 - \eta\tau) \sum_{l=1}^t \alpha_{l,t} G^{(l-1)}(s') - \sum_{l=1}^t \alpha_{l,t} G^{(l)}(s') \right].
\end{aligned}$$

Note that

$$\begin{aligned}
& (1 - \eta\tau) \sum_{l=1}^t \alpha_{l,t} G^{(l-1)}(s') - \sum_{l=1}^t \alpha_{l,t} G^{(l)}(s') \\
& \leq \sum_{l=1}^{t-1} ((1 - \eta\tau)\alpha_{l+1,t} - \alpha_{l,t}) G^{(l)}(s') + \alpha_{1,t} G^{(0)}(s') \\
& \leq \alpha_{1,t} G^{(0)}(s') \leq 2\alpha_{0,t} \eta\tau G^{(0)}(s') \leq 2\alpha_{0,t},
\end{aligned}$$

where the second step is due to

$$\begin{aligned}
(1 - \eta\tau)\alpha_{l+1,t} - \alpha_{l,t} & = ((1 - \eta\tau)\alpha_{l+1} - \alpha_l(1 - \alpha_{l+1})) \prod_{j=l+2}^t \alpha_j \\
& \leq ((1 - \eta\tau)\alpha_{l+1} - \alpha_{l+1} + \alpha_l\alpha_{l+1}) \prod_{j=l+2}^t \alpha_j \\
& = \alpha_{l+1}(\alpha_l - \eta\tau) \prod_{j=l+2}^t \alpha_j \leq 0.
\end{aligned} \tag{40}$$

So we conclude that

$$\begin{aligned}
& Q^{(t+1)}(s, a, b) - Q_\tau^*(s, a, b) \\
& \leq \gamma \cdot \frac{1 + (1 - \gamma)/4}{1 - (1 - \gamma)/4} \mathbb{E}_{s' \sim P(s, a, b)} \left[\sum_{l=0}^t \alpha_{l,t} \left[\|Q^{(l)}(s') - Q_\tau^*(s')\|_\infty + \frac{2\eta}{1 - \gamma} \|Q^{(l)}(s') - Q^{(l-1)}(s')\|_\infty \right] \right] \\
& \quad + 2\alpha_{0,t} \\
& \leq \frac{1 + \gamma}{2} \cdot \mathbb{E}_{s' \sim P(s, a, b)} \left[\sum_{l=0}^t \alpha_{l,t} \left[\|Q^{(l)}(s') - Q_\tau^*(s')\|_\infty + \frac{2\eta}{1 - \gamma} \|Q^{(l)}(s') - Q^{(l-1)}(s')\|_\infty \right] \right] \\
& \quad + 2\alpha_{0,t}.
\end{aligned}$$

The other side of (33) can be obtained by computing $\frac{1-\gamma}{4} \cdot$ (35) + (36) and following a similar argument, and is therefore omitted.

For $t = 0$, we have $|Q^{(1)}(s, a, b) - Q_\tau^*(s, a, b)| \leq \gamma \max_{s' \in \mathcal{S}} |f_{s'}^{(0)} - f_{s'}^*| \leq \frac{2\gamma}{1-\gamma}$.

C.4 PROOF OF LEMMA 4

For $t \geq 1$, let

$$u_t = \eta \|Q_\tau^*(s) - Q^{(t)}(s)\|_{\Gamma(\rho)} + \frac{12\eta^2}{(1-\gamma)^2} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_{\Gamma(\rho)}.$$

It follows that

$$u_1 \leq \frac{2\gamma\eta}{1-\gamma} + \frac{24\eta^2}{(1-\gamma)^3} \leq 1.$$

When $t \geq 1$, invoking Lemma 2 and Lemma 3 gives

$$\begin{aligned}
u_{t+1} &\leq \left(1 - \frac{1-\gamma}{2}\right) \sum_{l=1}^t \alpha_{l,t} \left[\eta \|Q^{(l)} - Q_\tau^*\|_{\Gamma(\rho)} + \left(\frac{2\eta^2}{1-\gamma} + \frac{12\eta^2}{(1-\gamma)^2}\right) \|Q^{(l)} - Q^{(l-1)}\|_{\Gamma(\rho)} \right] \\
&\quad + \frac{48\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \alpha_{l,t} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + 2\alpha_{0,t}\eta + \alpha_{0,t}\eta \|Q^{(0)} - Q_\tau^*\|_{\Gamma(\rho)} \\
&\leq \left(1 - \frac{1-\gamma}{3}\right) \sum_{l=1}^t \alpha_{l,t} u_l + \frac{48\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \alpha_{l,t} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{4\eta}{1-\gamma} \alpha_{0,t}.
\end{aligned} \tag{41}$$

Let

$$\beta_{l,t} = \alpha_l \prod_{i=l+1}^t \left(1 - \frac{1-\gamma}{3} \cdot \alpha_i\right).$$

It follows that for $t \geq 0$,

$$\begin{aligned}
&\sum_{l=1}^{t+1} \alpha_{l,t+1} u_l \\
&= (1 - \alpha_{t+1}) \sum_{l=1}^t \alpha_{l,t} u_l + \alpha_{t+1} u_{t+1} \\
&\leq \left(1 - \frac{1-\gamma}{3} \cdot \alpha_{t+1}\right) \sum_{l=1}^t \alpha_{l,t} u_l + \alpha_{t+1} \frac{48\eta\mathcal{C}_\rho}{(1-\gamma)^2} \cdot \sum_{l=1}^t \alpha_{l,t} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{4\eta}{1-\gamma} \alpha_{t+1} \alpha_{0,t} \\
&\leq \prod_{l=2}^{t+1} \left(1 - \frac{1-\gamma}{3} \cdot \alpha_l\right) \alpha_{1,1} u_1 + \frac{48\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{i=1}^t \beta_{i+1,t+1} \sum_{l=1}^i \alpha_{l,i} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{4\eta}{1-\gamma} \sum_{i=1}^t \alpha_{0,i} \beta_{i+1,t+1} \\
&\leq \beta_{1,t+1} u_1 + \frac{48\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \sum_{i=l}^t \alpha_{l,i} \beta_{i+1,t+1} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{4\eta}{1-\gamma} \sum_{i=1}^t \alpha_{0,i} \beta_{i+1,t+1} \\
&\leq \frac{200\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \beta_{l,t+1} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{18\eta}{1-\gamma} \beta_{0,t+1},
\end{aligned} \tag{42}$$

where the last step is due to the following lemma. Similar lemma has appeared in prior works (see i.e., (Wei et al., 2021b, Lemma 36)). Our version features a simpler proof, which is postponed to Appendix E.4.

Lemma 14. *Let two sequences $\{\delta_i\}, \{\xi_i\}$ be defined as*

$$\delta_i = 1 - c_1 \alpha_i, \quad \text{and} \quad \xi_i = 1 - c_2 \alpha_i,$$

where the constants c_1, c_2 satisfy $0 < c_1 < c_2 < \frac{1}{2\alpha_i}$. For $l \leq t$, let $\delta_{l,t} = \alpha_l \prod_{i=l+1}^t \delta_i$ and $\xi_{l,t} = \alpha_l \prod_{i=l+1}^t \xi_i$. We have

$$\sum_{i=l}^t \xi_{l,i} \delta_{i+1,t} \leq \left(1 + \frac{2}{c_2 - c_1}\right) \delta_{l,t}.$$

Substitution of (42) into (41) gives

$$\begin{aligned}
u_{t+1} &\leq \left(1 - \frac{1-\gamma}{3}\right) \sum_{l=1}^t \alpha_{l,t} u_l + \frac{48\eta}{(1-\gamma)^2} \sum_{l=1}^t \alpha_{l,t} \text{KL}(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{4\eta}{1-\gamma} \alpha_{0,t} \\
&\leq \frac{200\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \beta_{l,t} \text{KL}(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{18\eta}{1-\gamma} \beta_{0,t} + \frac{48\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \alpha_{l,t} \text{KL}(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{4\eta}{1-\gamma} \alpha_{0,t} \\
&\leq \frac{250\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \beta_{l,t} \text{KL}(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{22\eta}{1-\gamma} \beta_{0,t}.
\end{aligned}$$

for $t \geq 1$. It is straightforward to verify that the above inequality holds for $t = 0$ as well.

So we conclude that

$$\begin{aligned}
\sum_{l=0}^t \lambda_{l+1,t+1} u_{l+1} &= \sum_{i=0}^t \lambda_{i+1,t+1} u_{i+1} \\
&\leq \sum_{i=0}^t \lambda_{i+1,t+1} \left[\frac{250\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^i \beta_{l,i} \text{KL}(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{22\eta}{1-\gamma} \beta_{0,i} \right] \\
&= \frac{250\eta\mathcal{C}_\rho}{(1-\gamma)^2} \sum_{l=1}^t \sum_{i=l}^t \beta_{l,i} \lambda_{i+1,t+1} \text{KL}(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{22\eta}{1-\gamma} \sum_{i=0}^t \beta_{0,i} \lambda_{i+1,t+1} \\
&\leq \frac{6250\eta\mathcal{C}_\rho}{(1-\gamma)^3} \sum_{l=1}^t \lambda_{l,t+1} \text{KL}(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) + \frac{550\eta}{(1-\gamma)^2} \lambda_{0,t+1} \\
&= \frac{6250\eta\mathcal{C}_\rho}{(1-\gamma)^3} \sum_{l=0}^{t-1} \lambda_{l+1,t+1} \text{KL}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) + \frac{550\eta}{(1-\gamma)^2} \lambda_{0,t+1},
\end{aligned}$$

where the penultimate step invokes Lemma 14.

C.5 PROOF OF LEMMA 5

Taking logarithm on the both sides of the update rule (9b), we get

$$\begin{cases} \log \bar{\mu}^{(t+1)}(s) - (1-\eta\tau) \log \mu^{(t)}(s) & \stackrel{\text{1}}{=} \eta Q^{(t)}(s) \bar{\nu}^{(t)}(s) \\ \log \bar{\nu}^{(t+1)}(s) - (1-\eta\tau) \log \nu^{(t)}(s) & \stackrel{\text{1}}{=} -\eta Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) \end{cases}. \quad (43)$$

Subtracting (27) from (43) and taking inner product with $\bar{\zeta}^{(t+1)}(s) - \zeta_\tau^*(s)$ gives

$$\begin{aligned}
&\langle \log \bar{\zeta}^{(t+1)}(s) - (1-\eta\tau) \log \zeta^{(t)}(s) - \eta\tau \log \zeta_\tau^*(s), \bar{\zeta}^{(t+1)}(s) - \zeta_\tau^*(s) \rangle \\
&= \eta \langle \bar{\mu}^{(t+1)}(s) - \mu_\tau^*(s), Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q_\tau^*(s) \nu_\tau^*(s) \rangle \\
&\quad - \eta \langle \bar{\nu}^{(t+1)}(s) - \nu_\tau^*(s), Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) - Q_\tau^*(s)^\top \mu_\tau^*(s) \rangle \\
&\leq \eta \langle \bar{\mu}^{(t+1)}(s) - \mu_\tau^*(s), Q^{(t)}(s) (\bar{\nu}^{(t)}(s) - \nu_\tau^*(s)) \rangle \\
&\quad - \eta \langle \bar{\nu}^{(t+1)}(s) - \nu_\tau^*(s), Q^{(t)}(s)^\top (\bar{\mu}^{(t)}(s) - \mu_\tau^*(s)) \rangle + 2\eta \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty \\
&\leq \eta \langle \bar{\mu}^{(t+1)}(s) - \mu_\tau^*(s), Q^{(t)}(s) (\bar{\nu}^{(t)}(s) - \bar{\nu}^{(t+1)}(s)) \rangle \\
&\quad - \eta \langle \bar{\nu}^{(t+1)}(s) - \nu_\tau^*(s), Q^{(t)}(s)^\top (\bar{\mu}^{(t)}(s) - \bar{\mu}^{(t+1)}(s)) \rangle + 2\eta \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty \\
&\leq \frac{2\eta}{1-\gamma} \left(2\text{KL}_s(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) + \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) + \text{KL}_s(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) \right) + 2\eta \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty.
\end{aligned}$$

LHS can be written as

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

So we conclude that

$$\left(1 - \frac{4\eta}{1-\gamma}\right) \text{KL}_s(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) - (1-\eta\tau) \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t)})$$

$$\begin{aligned}
& + \left(1 - \eta\tau - \frac{2\eta}{1-\gamma}\right) \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \zeta^{(t)}) + \eta\tau \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\
& \leq \frac{2\eta}{1-\gamma} \text{KL}_s(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) + 2\eta \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty.
\end{aligned}$$

With $0 < \eta \leq \frac{1-\gamma}{8}$, we have

$$\begin{aligned}
& \frac{1}{2} \text{KL}_s(\zeta_\tau^* \parallel \bar{\zeta}^{(t+1)}) + \eta\tau \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \zeta_\tau^*) \\
& \leq (1 - \eta\tau) \text{KL}_s(\zeta_\tau^* \parallel \zeta^{(t)}) + \frac{2\eta}{1-\gamma} \text{KL}_s(\zeta^{(t)} \parallel \bar{\zeta}^{(t)}) + 2\eta \|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty.
\end{aligned}$$

C.6 PROOF OF LEMMA 6

We have

$$\begin{aligned}
V_\tau^{\mu,\nu}(s) - V_\tau^*(s) & = \mu(s)^\top Q_\tau^{\mu,\nu}(s)\nu(s) + \tau\mathcal{H}(\mu(s)) - \tau\mathcal{H}(\nu(s)) \\
& \quad - \mu_\tau^*(s)^\top Q_\tau^*(s)\nu_\tau^*(s) - \tau\mathcal{H}(\mu_\tau^*(s)) + \tau\mathcal{H}(\nu_\tau^*(s)) \\
& = \mu(s)^\top Q_\tau^{\mu,\nu}(s)\nu(s) - \mu(s)^\top Q_\tau^*(s)\nu(s) + f_s(Q_\tau^*, \mu, \nu) - f_s(Q_\tau^*, \mu_\tau^*, \nu_\tau^*) \\
& = \gamma \mathbb{E}_{\substack{a \sim \mu(\cdot|s), \\ b \sim \nu(\cdot|s), \\ s' \sim P(\cdot|s, a, b)}} [V_\tau^{\mu,\nu}(s') - V_\tau^*(s')] + f_s(Q_\tau^*, \mu, \nu) - f_s(Q_\tau^*, \mu_\tau^*, \nu_\tau^*).
\end{aligned}$$

Applying the inequality recursively and averaging s over ρ , we arrive at

$$V_\tau^{\mu,\nu}(\rho) - V_\tau^*(\rho) = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_\rho^{\mu,\nu}} [f_{s'}(Q_\tau^*, \mu, \nu) - f_{s'}(Q_\tau^*, \mu_\tau^*, \nu_\tau^*)], \quad (44)$$

which is the well-known performance difference lemma applied to the setting of Markov games. It follows that

$$\begin{aligned}
V_\tau^{\mu_\tau^\dagger(\nu),\nu}(\rho) - V_\tau^*(\rho) & = \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_\rho^{\mu_\tau^\dagger(\nu),\nu}} [f_{s'}(Q_\tau^*, \mu_\tau^\dagger(\nu), \nu) - f_{s'}(Q_\tau^*, \mu_\tau^*, \nu_\tau^*)] \\
& \leq \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_\rho^{\mu_\tau^\dagger(\nu),\nu}} [f_{s'}(Q_\tau^*, \mu_\tau^\dagger(\nu), \nu) - f_{s'}(Q_\tau^*, \mu, \nu_\tau^*)] \\
& \leq \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_\rho^{\mu_\tau^\dagger(\nu),\nu}} \left[\max_{\mu', \nu'} (f_{s'}(Q_\tau^*, \mu', \nu) - f_{s'}(Q_\tau^*, \mu, \nu')) \right] \\
& \leq \frac{C_{\rho,\tau}^\dagger}{1-\gamma} \mathbb{E}_{s \sim \rho} \left[\max_{\mu', \nu'} (f_s(Q_\tau^*, \mu', \nu) - f_s(Q_\tau^*, \mu, \nu')) \right].
\end{aligned} \quad (45)$$

A similar argument gives $V_\tau^*(\rho) - V_\tau^{\mu,\nu_\tau^\dagger(\mu)}(\rho) \leq \frac{C_{\rho,\tau}^\dagger}{1-\gamma} \mathbb{E}_{s \sim \rho} \left[\max_{\mu', \nu'} (f_s(Q_\tau^*, \mu', \nu) - f_s(Q_\tau^*, \mu, \nu')) \right]$.

Summing the two inequalities proves (18). Alternatively, we continue from (45) and show that

$$\begin{aligned}
V_\tau^{\mu_\tau^\dagger(\nu),\nu}(s) - V_\tau^*(s) & \leq \frac{1}{1-\gamma} \mathbb{E}_{s' \sim d_s^{\mu_\tau^\dagger(\nu),\nu}} \left[\max_{\mu', \nu'} (f_{s'}(Q_\tau^*, \mu', \nu) - f_{s'}(Q_\tau^*, \mu, \nu')) \right] \\
& \leq \frac{\|1/\rho\|_\infty}{1-\gamma} \mathbb{E}_{s \sim \rho} \left[\max_{\mu', \nu'} (f_s(Q_\tau^*, \mu', \nu) - f_s(Q_\tau^*, \mu, \nu')) \right].
\end{aligned}$$

Summing the inequality with the one for $V_\tau^*(s) - V_\tau^{\mu,\nu_\tau^\dagger(\mu)}(s)$ and taking maximum over $s \in \mathcal{S}$ completes the proof for (19).

D PROOF OF KEY LEMMAS FOR THE EPISODIC SETTING

D.1 PROOF OF LEMMA 9

Following the similar argument of arriving (28), we have

$$\langle \log \zeta_h^{(t+1)}(s) - (1 - \eta\tau) \log \zeta_h^{(t)}(s) - \eta\tau \log \zeta_{h,\tau}^*(s), \bar{\zeta}_h^{(t+1)}(s) - \zeta_{h,\tau}^*(s) \rangle$$

$$\leq 2\eta \|Q_h^{(t+1)}(s) - Q_{h,\tau}^*(s)\|_\infty.$$

We rewrite LHS as

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

Rearranging terms gives

$$\begin{aligned} & \text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t+1)}) - (1 - \eta\tau) \text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t)}) + (1 - \eta\tau) \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_h^{(t)}) \\ & \quad + \eta\tau \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_{h,\tau}^*) + \text{KL}_s(\zeta_h^{(t+1)} \parallel \bar{\zeta}_h^{(t+1)}) \\ & \quad - \langle \log \bar{\zeta}_h^{(t+1)}(s) - \log \zeta_h^{(t+1)}(s), \bar{\zeta}_h^{(t+1)}(s) - \zeta_h^{(t+1)}(s) \rangle \\ & \leq 2\eta \|Q_h^{(t+1)}(s) - Q_{h,\tau}^*(s)\|_\infty. \end{aligned} \tag{46}$$

Note that

$$\begin{aligned} & \langle \log \bar{\mu}_h^{(t+1)}(s) - \log \mu_h^{(t+1)}(s), \bar{\mu}_h^{(t+1)}(s) - \mu_h^{(t+1)}(s) \rangle \\ &= \eta \langle Q_h^{(t)}(s) \bar{\nu}_h^{(t)}(s) - Q_h^{(t+1)}(s) \bar{\nu}_h^{(t+1)}(s), \bar{\mu}_h^{(t+1)}(s) - \mu_h^{(t+1)}(s) \rangle \\ & \leq \eta \|Q_h^{(t)}(s) \bar{\nu}_h^{(t)}(s) - Q_h^{(t+1)}(s) \bar{\nu}_h^{(t+1)}(s)\|_1 \|\bar{\mu}_h^{(t+1)}(s) - \mu_h^{(t+1)}(s)\|_1. \end{aligned} \tag{47}$$

We bound $\|Q_h^{(t)}(s) \bar{\nu}_h^{(t)}(s) - Q_h^{(t+1)}(s) \bar{\nu}_h^{(t+1)}(s)\|_1$ as

$$\begin{aligned} & \|Q_h^{(t)}(s) \bar{\nu}_h^{(t)}(s) - Q_h^{(t+1)}(s) \bar{\nu}_h^{(t+1)}(s)\|_1 \\ & \leq \|Q_h^{(t+1)}(s) (\bar{\nu}_h^{(t)}(s) - \bar{\nu}_h^{(t+1)}(s))\|_1 + \|(Q_h^{(t)}(s) - Q_h^{(t+1)}(s)) \bar{\nu}_h^{(t)}(s)\|_1 \\ & \leq 2H \|\bar{\nu}_h^{(t)}(s) - \bar{\nu}_h^{(t+1)}(s)\|_1 + \|Q_h^{(t)}(s) - Q_h^{(t+1)}(s)\|_\infty \\ & \leq 2H \|\bar{\nu}_h^{(t+1)}(s) - \nu_h^{(t)}(s)\|_1 + 2H \|\nu_h^{(t)}(s) - \bar{\nu}_h^{(t)}(s)\|_1 + \|Q_h^{(t)}(s) - Q_h^{(t+1)}(s)\|_\infty. \end{aligned}$$

Plugging the above inequality into (47) and invoking Young's inequality yields

$$\begin{aligned} & \langle \log \bar{\mu}_h^{(t+1)}(s) - \log \mu_h^{(t+1)}(s), \bar{\mu}_h^{(t+1)}(s) - \mu_h^{(t+1)}(s) \rangle \\ & \leq \eta H \left(\|\bar{\nu}_h^{(t+1)}(s) - \nu_h^{(t)}(s)\|_1^2 + \|\nu_h^{(t)}(s) - \bar{\nu}_h^{(t)}(s)\|_1^2 + 2 \|\bar{\mu}_h^{(t+1)}(s) - \mu_h^{(t+1)}(s)\|_1^2 \right) \\ & \quad + \eta \|Q_h^{(t)}(s) - Q_h^{(t+1)}(s)\|_\infty \|\bar{\mu}_h^{(t+1)}(s) - \mu_h^{(t+1)}(s)\|_1 \\ & \leq 2\eta H \text{KL}_s(\bar{\nu}_h^{(t+1)} \parallel \nu_h^{(t)}) + 2\eta H \text{KL}_s(\nu_h^{(t)} \parallel \bar{\nu}_h^{(t)}) + 4\eta H \text{KL}_s(\mu_h^{(t+1)} \parallel \bar{\mu}_h^{(t+1)}) \\ & \quad + 2\eta^2 H \|Q_h^{(t)}(s) - Q_h^{(t+1)}(s)\|_\infty, \end{aligned}$$

where the last step results from Pinsker's inequality and Lemma 8. Similarly, we have

$$\begin{aligned} & \langle \log \bar{\nu}_h^{(t+1)}(s) - \log \nu_h^{(t+1)}(s), \bar{\nu}_h^{(t+1)}(s) - \nu_h^{(t+1)}(s) \rangle \\ & \leq 2\eta H \text{KL}_s(\bar{\mu}_h^{(t+1)} \parallel \mu_h^{(t)}) + 2\eta H \text{KL}_s(\mu_h^{(t)} \parallel \bar{\mu}_h^{(t)}) + 4\eta H \text{KL}_s(\nu_h^{(t+1)} \parallel \bar{\nu}_h^{(t+1)}) \\ & \quad + 2\eta^2 H \|Q_h^{(t)}(s) - Q_h^{(t+1)}(s)\|_\infty. \end{aligned}$$

Summing the above two inequalities gives

$$\langle \log \bar{\zeta}_h^{(t+1)}(s) - \log \zeta_h^{(t+1)}(s), \bar{\zeta}_h^{(t+1)}(s) - \zeta_h^{(t+1)}(s) \rangle$$

$$\begin{aligned}
&\leq 2\eta H \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_h^{(t)}) + 2\eta H \text{KL}_s(\zeta_h^{(t)} \parallel \bar{\zeta}_h^{(t)}) + 4\eta H \text{KL}_s(\zeta_h^{(t+1)} \parallel \bar{\zeta}_h^{(t+1)}) \\
&\quad + 4\eta^2 H \|Q_h^{(t)}(s) - Q_h^{(t+1)}(s)\|_\infty \\
&\leq 2\eta H \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_h^{(t)}) + 2\eta H \text{KL}_s(\zeta_h^{(t)} \parallel \bar{\zeta}_h^{(t)}) + 4\eta H \text{KL}_s(\zeta_h^{(t+1)} \parallel \bar{\zeta}_h^{(t+1)}) \\
&\quad + \frac{\eta}{2} \left(\|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty + \|Q_h^{(t+1)}(s) - Q_{h,\tau}^*(s)\|_\infty \right),
\end{aligned}$$

where the second step invokes triangular inequality and the fact that $\eta \leq \frac{1}{8H}$. Plugging the above inequality into (46) gives

$$\begin{aligned}
&\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t+1)}) - (1 - \eta\tau) \text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t)}) + (1 - \eta(\tau + 2H)) \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_h^{(t)}) \\
&\quad + \eta\tau \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_{h,\tau}^*) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t+1)} \parallel \bar{\zeta}_h^{(t+1)}) - 2\eta H \text{KL}_s(\zeta_h^{(t)} \parallel \bar{\zeta}_h^{(t)}) \\
&\leq \frac{5\eta}{2} \|Q_h^{(t+1)}(s) - Q_{h,\tau}^*(s)\|_\infty + \frac{\eta}{2} \|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty.
\end{aligned}$$

With $\eta \leq \frac{1}{8H}$, we have $(1 - \eta\tau)(1 - 4\eta H) \geq 2\eta H$ and $1 - \eta(\tau + 2H) \geq 0$. It follows that

$$\begin{aligned}
&\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t+1)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t+1)} \parallel \bar{\zeta}_h^{(t+1)}) + \eta\tau \text{KL}_s(\bar{\zeta}_h^{(t+1)} \parallel \zeta_{h,\tau}^*) \\
&\leq (1 - \eta\tau) \text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t)}) + 2\eta H \text{KL}_s(\zeta_h^{(t)} \parallel \bar{\zeta}_h^{(t)}) \\
&\quad + \frac{5\eta}{2} \|Q_h^{(t+1)}(s) - Q_{h,\tau}^*(s)\|_\infty + \frac{\eta}{2} \|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty \\
&\leq (1 - \eta\tau) \left(\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t)} \parallel \bar{\zeta}_h^{(t)}) \right) \\
&\quad + \frac{5\eta}{2} \|Q_h^{(t+1)}(s) - Q_{h,\tau}^*(s)\|_\infty + \frac{\eta}{2} \|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty.
\end{aligned}$$

Therefore, it holds for $0 \leq t_1 < t_2$ that

$$\begin{aligned}
&\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t_2)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t_2)} \parallel \bar{\zeta}_h^{(t_2)}) + \eta\tau \text{KL}_s(\bar{\zeta}_h^{(t_2)} \parallel \zeta_{h,\tau}^*) \\
&\leq (1 - \eta\tau)^{t_2-t_1} \left(\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t_1)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t_1)} \parallel \bar{\zeta}_h^{(t_1)}) \right) \\
&\quad + \sum_{t'=t_1+1}^{t_2} (1 - \eta\tau)^{t_2-t'} \left[\frac{5\eta}{2} \|Q_h^{(t')}(s) - Q_{h,\tau}^*(s)\|_\infty + \frac{\eta}{2} \|Q_h^{(t'-1)}(s) - Q_{h,\tau}^*(s)\|_\infty \right] \\
&\leq (1 - \eta\tau)^{t_2-t_1} \left(\text{KL}_s(\zeta_{h,\tau}^* \parallel \zeta_h^{(t_1)}) + (1 - 4\eta H) \text{KL}_s(\zeta_h^{(t_1)} \parallel \bar{\zeta}_h^{(t_1)}) \right) \\
&\quad + 4\eta \sum_{l=t_1}^{t_2} (1 - \eta\tau)^{t_2-l} \|Q_h^{(l)}(s) - Q_{h,\tau}^*(s)\|_\infty.
\end{aligned}$$

D.2 PROOF OF LEMMA 10

For $t_2 > 0$, we have

$$\begin{aligned}
&Q_{h-1}^{(t_2)}(s, a, b) - Q_{h-1,\tau}^*(s, a, b) \\
&= \mathbb{E}_{s' \sim P_{h-1}(\cdot | s, a, b)} \left[V_h^{(t_2-1)}(s') - V_{h,\tau}^*(s') \right] \\
&= \mathbb{E}_{s' \sim P_{h-1}(\cdot | s, a, b)} \left[(1 - \eta\tau)^{t_2-t_1} (V_h^{(t_1-1)}(s') - V_{h,\tau}^*(s')) \right. \\
&\quad \left. + \eta\tau \sum_{l=t_1}^{t_2-1} (1 - \eta\tau)^{t_2-1-l} (f_{s'}(Q_h^{(l)}, \bar{\mu}_h^{(l)}, \bar{\nu}_h^{(l)}) - f_{s'}(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*)) \right] \\
&\leq (1 - \eta\tau)^{t_2-t_1} 2H + \mathbb{E}_{s' \sim P_{h-1}(\cdot | s, a, b)} \left[\eta\tau \sum_{l=t_1}^{t_2-1} (1 - \eta\tau)^{t_2-1-l} (f_{s'}(Q_h^{(l)}, \bar{\mu}_h^{(l)}, \bar{\nu}_h^{(l)}) - f_{s'}(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*)) \right].
\end{aligned} \tag{48}$$

We start by decomposing $f_s^{(t)} - f_s^*$ as

$$\begin{aligned}
& f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) - f_s(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*) \\
&= \left(f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) - f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \nu_{h,\tau}^*) \right) + f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \nu_{h,\tau}^*) - f_s(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*) \\
&\leq \left(f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) - f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \nu_{h,\tau}^*) \right) + f_s(Q_{h,\tau}^*, \bar{\mu}^{(t)}, \nu_{h,\tau}^*) - f_s(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*) \\
&\quad + \|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty \\
&\leq f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) - f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \nu_{h,\tau}^*) + \|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty.
\end{aligned}$$

Note that Lemma 13 can be applied to the episodic setting by simply replacing $1/(1-\gamma)$ with H , which yields

$$\begin{aligned}
& f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) - f_s(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*) \\
&\leq \|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty + 2\eta H \|Q_h^{(t)}(s) - Q_h^{(t-1)}(s)\|_\infty \\
&\quad + \frac{1-\eta\tau}{\eta} \text{KL}_s(\nu_{h,\tau}^* \parallel \nu_h^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\nu_{h,\tau}^* \parallel \nu_h^{(t)}) \\
&\quad - \frac{1}{\eta} (1-4\eta H) \text{KL}_s(\nu_h^{(t)} \parallel \bar{\nu}_h^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\bar{\nu}_h^{(t)} \parallel \nu_h^{(t-1)}) \\
&\quad + 2H \left(\text{KL}_s(\bar{\mu}_h^{(t)} \parallel \mu_h^{(t-1)}) + \text{KL}_s(\mu_h^{(t-1)} \parallel \bar{\mu}_h^{(t-1)}) \right). \tag{49}
\end{aligned}$$

By a similar argument,

$$\begin{aligned}
& f_s(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*) - f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) \\
&\leq \|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty + 2\eta H \|Q_h^{(t)}(s) - Q_h^{(t-1)}(s)\|_\infty \\
&\quad + \frac{1-\eta\tau}{\eta} \text{KL}_s(\mu_{h,\tau}^* \parallel \mu_h^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\mu_{h,\tau}^* \parallel \mu_h^{(t)}) \\
&\quad - \frac{1}{\eta} (1-4\eta H) \text{KL}_s(\mu_h^{(t)} \parallel \bar{\mu}_h^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\bar{\mu}_h^{(t)} \parallel \mu_h^{(t-1)}) \\
&\quad + 2H \left(\text{KL}_s(\bar{\nu}_h^{(t)} \parallel \nu_h^{(t-1)}) + \text{KL}_s(\nu_h^{(t-1)} \parallel \bar{\nu}_h^{(t-1)}) \right). \tag{50}
\end{aligned}$$

Computing (49) + $\frac{2}{3}$ · (50) gives

$$\begin{aligned}
& \frac{1}{3} [f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) - f_s(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*)] \\
&\leq \frac{5}{3} [\|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty + 2\eta H \|Q_h^{(t)}(s) - Q_h^{(t-1)}(s)\|_\infty] \\
&\quad + \frac{1-\eta\tau}{\eta} \left[\text{KL}_s(\nu_{h,\tau}^* \parallel \nu_h^{(t-1)}) + \frac{2}{3} \text{KL}_s(\mu_{h,\tau}^* \parallel \mu_h^{(t-1)}) \right] - \frac{1}{\eta} \left[\text{KL}_s(\nu_{h,\tau}^* \parallel \nu_h^{(t)}) + \frac{2}{3} \text{KL}_s(\mu_{h,\tau}^* \parallel \mu_h^{(t)}) \right] \\
&\quad + 2H \left[\text{KL}_s(\mu_h^{(t-1)} \parallel \bar{\mu}_h^{(t-1)}) + \frac{2}{3} \text{KL}_s(\nu_h^{(t-1)} \parallel \bar{\nu}_h^{(t-1)}) \right] \\
&\quad - \frac{1}{\eta} (1-4\eta H) \left[\frac{2}{3} \text{KL}_s(\mu_h^{(t)} \parallel \bar{\mu}_h^{(t)}) + \text{KL}_s(\nu_h^{(t)} \parallel \bar{\nu}_h^{(t)}) \right] \\
&\quad + \left(2H - \frac{1-\eta\tau}{\eta} \cdot \frac{2}{3} \right) \text{KL}_s(\bar{\mu}^{(t)} \parallel \mu^{(t-1)}) + \left(2H \cdot \frac{2}{3} - \frac{1-\eta\tau}{\eta} \right) \text{KL}_s(\bar{\nu}^{(t)} \parallel \nu^{(t-1)}). \tag{51}
\end{aligned}$$

With $\eta \leq \frac{1}{8H}$, we have

$$2H - \frac{1-\eta\tau}{\eta} \cdot \frac{2}{3} \leq 0, \quad 2H \cdot \frac{2}{3} - \frac{1-\eta\tau}{\eta} \leq 0, \quad \text{and} \quad \frac{1}{\eta} (1-\eta\tau)(1-4\eta H) \cdot \frac{2}{3} \geq 2H.$$

Let

$$G_h^{(t)}(s) = \text{KL}_s(\nu_{h,\tau}^* \parallel \nu_h^{(t)}) + \frac{2}{3} \text{KL}_s(\mu_{h,\tau}^* \parallel \mu_h^{(t)})$$

$$+ \frac{2}{3}(1 - 4\eta H) \left[\text{KL}_s(\mu_h^{(t)} \parallel \bar{\mu}_h^{(t)}) + \text{KL}_s(\nu_h^{(t)} \parallel \bar{\nu}_h^{(t)}) \right].$$

We can simplify (51) as

$$\begin{aligned} & f_s(Q_h^{(t)}, \bar{\mu}_h^{(t)}, \bar{\nu}_h^{(t)}) - f_s(Q_{h,\tau}^*, \mu_{h,\tau}^*, \nu_{h,\tau}^*) \\ & \leq 5 \left[\|Q_h^{(t)}(s) - Q_{h,\tau}^*(s)\|_\infty + 2\eta H \|Q_h^{(t)}(s) - Q_h^{(t-1)}(s)\|_\infty \right] + \frac{1 - \eta\tau}{\eta} G_h^{(t-1)}(s) - \frac{1}{\eta} G_h^{(t)}(s). \end{aligned}$$

Plugging the above inequality into (48) gives

$$\begin{aligned} & Q_{h-1}^{(t_2)}(s, a, b) - Q_{h-1,\tau}^*(s, a, b) \\ & \leq (1 - \eta\tau)^{t_2-t_1} 2H \\ & \quad + \mathbb{E}_{s' \sim P_{h-1}(\cdot|s,a,b)} \left[5\eta\tau \sum_{l=t_1}^{t_2-1} (1 - \eta\tau)^{t_2-1-l} \left(\|Q_h^{(l)}(s') - Q_{h,\tau}^*(s')\|_\infty + 2\eta H \|Q_h^{(l)}(s') - Q_h^{(l-1)}(s')\|_\infty \right) \right] \\ & \quad + \mathbb{E}_{s' \sim P_{h-1}(\cdot|s,a,b)} \left[\tau(1 - \eta\tau)^{t_2-t_1} G_h^{(t_1-1)}(s') \right] \\ & \leq (1 - \eta\tau)^{t_2-t_1} 2H \\ & \quad + 10\eta\tau \mathbb{E}_{s' \sim P_{h-1}(\cdot|s,a,b)} \left[\sum_{l=t_1-1}^{t_2-1} (1 - \eta\tau)^{t_2-1-l} \|Q_h^{(l)}(s') - Q_{h,\tau}^*(s')\|_\infty \right] \\ & \quad + \tau(1 - \eta\tau)^{t_2-t_1} \mathbb{E}_{s' \sim P_{h-1}(\cdot|s,a,b)} \left[\text{KL}_{s'}(\zeta_{h,\tau}^* \parallel \zeta_h^{(t_1-1)}) + (1 - 4\eta H) \text{KL}_{s'}(\zeta_h^{(t_1-1)} \parallel \bar{\zeta}_h^{(t_1-1)}) \right]. \end{aligned}$$

E PROOF OF AUXILIARY LEMMAS

E.1 PROOF OF LEMMA 11

We first single out a set of bounds for $V^{(t)}$ and $Q^{(t)}$, which can be obtained by a simple induction:

$$\forall (s, a, b) \in \mathcal{S} \times \mathcal{A} \times \mathcal{B}, \quad \begin{cases} -\frac{\tau \log |\mathcal{B}|}{1-\gamma} \leq V^{(t)}(s) \leq \frac{1+\tau \log |\mathcal{A}|}{1-\gamma} \\ -\frac{\gamma\tau \log |\mathcal{B}|}{1-\gamma} \leq Q^{(t)}(s, a, b) \leq \frac{1+\gamma\tau \log |\mathcal{A}|}{1-\gamma} \end{cases}. \quad (52)$$

We invoke the following lemma to bound several key quantities that will be helpful in the analysis.

Lemma 15 ((Mei et al., 2020, Lemma 24)). *Let $\pi, \pi' \in \Delta(\mathcal{A})$ such that $\pi(a) \propto \exp(\theta(a))$, $\pi'(a) \propto \theta'(a)$ for some $\theta, \theta' \in \mathbb{R}^{|\mathcal{A}|}$. It holds that*

$$\|\pi - \pi'\|_1 \leq \|\theta - \theta'\|_\infty.$$

With this lemma in mind, for any $t \geq 0$, it follows that

$$\begin{aligned} \|\bar{\mu}^{(t+1)}(s) - \mu^{(t+1)}(s)\|_1 & \leq \min_{c \in \mathbb{R}} \|\log \bar{\mu}^{(t+1)}(s) - \log \mu^{(t+1)}(s) - c \cdot \mathbf{1}\|_\infty \\ & \leq \eta \|Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s)\|_\infty \\ & \leq \eta \cdot \frac{1 + \gamma\tau(\log |\mathcal{A}| + \log |\mathcal{B}|)}{1 - \gamma} \leq \frac{2\eta}{1 - \gamma}, \end{aligned}$$

and a similar argument reveals that

$$\|\bar{\nu}^{(t+1)}(s) - \nu^{(t+1)}(s)\|_1 \leq \frac{2\eta}{1 - \gamma}.$$

Next we make note of the fact that when $t \geq 1$,

$$\begin{aligned} \bar{\mu}^{(t+1)}(a|s) & \propto \mu^{(t)}(a|s)^{1-\eta\tau} \exp(\eta Q^{(t)}(s) \bar{\nu}^{(t)}(s)) \\ & \propto \bar{\mu}^{(t)}(a|s)^{1-\eta\tau} \exp\left(\eta [Q^{(t)}(s) \bar{\nu}^{(t)}(s) + (1 - \eta\tau)(Q^{(t)}(s) \bar{\nu}^{(t)}(s) - Q^{(t-1)}(s) \bar{\nu}^{(t-1)}(s))]\right) \\ & \propto \bar{\mu}^{(t)}(a|s) \exp(\eta w^{(t)}(a)), \end{aligned} \quad (53)$$

where

$$w^{(t)} = Q^{(t)}(s)\bar{v}^{(t)}(s) + (1 - \eta\tau)(Q^{(t)}(s)\bar{v}^{(t)}(s) - Q^{(t-1)}(s)\bar{v}^{(t-1)}(s)) - \tau \log \bar{\mu}^{(t)}(s)$$

satisfies

$$\begin{aligned} & \|w^{(t)}\|_\infty \\ & \leq \|Q^{(t)}(s)\bar{v}^{(t)}(s)\|_\infty + \|\tau \log \bar{\mu}^{(t)}(s)\|_\infty + (1 - \eta\tau)\|Q^{(t)}(s)\bar{v}^{(t)}(s) - Q^{(t-1)}(s)\bar{v}^{(t-1)}(s)\|_\infty \\ & \leq \frac{2}{1-\gamma} + \frac{2}{1-\gamma} + \frac{2(1-\eta\tau)}{1-\gamma} \leq \frac{6}{1-\gamma}, \end{aligned}$$

where the second step is due to the following bound:

$$\forall t \geq 0, s \in \mathcal{S}, \quad \max \{ \|\log \zeta^{(t)}(s)\|_\infty, \|\log \bar{\zeta}^{(t)}(s)\|_\infty \} \leq \frac{2}{(1-\gamma)\tau}. \quad (54)$$

Recall that when $t = 0$, we have $\bar{\mu}^{(t+1)} = \bar{\mu}^{(0)}$. So we have

$$\forall s \in \mathcal{S}, t \geq 0, \quad \|\bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s)\|_1 \leq \frac{6\eta}{1-\gamma}.$$

It remains to prove the claim (54).

Proof. It is worth noting that $\mu^{(t)}(s)$ can be always written as $\mu^{(t)}(a|s) \propto \exp(w^{(t)}(a)/\tau)$ for some $w^{(t)} \in \mathbb{R}^{|\mathcal{A}|}$ satisfying

$$\forall a \in \mathcal{A}, \quad -\frac{\gamma\tau \log |\mathcal{B}|}{1-\gamma} \leq w^{(t)}(a) \leq \frac{1 + \gamma\tau \log |\mathcal{A}|}{1-\gamma}.$$

To see this, note that the claim trivially holds for $t = 0$ with $w^{(0)} = \mathbf{0}$. When the statement holds for some $t \geq 0$, we have

$$\begin{aligned} \mu^{(t+1)}(a|s) & \propto \mu^{(t)}(a|s)^{1-\eta\tau} \exp(\eta Q^{(t+1)}(s)\bar{v}^{(t+1)}(s)) \\ & \propto \exp(((1-\eta\tau)w^{(t)} + \eta\tau Q^{(t+1)}(s)\bar{v}^{(t+1)}(s))/\tau) \\ & \propto \exp(w^{(t+1)}/\tau), \end{aligned}$$

with $w^{(t+1)} = (1 - \eta\tau)w^{(t)} + \eta\tau Q^{(t+1)}(s)\bar{v}^{(t+1)}(s)$. We conclude that the claim holds for $t + 1$ by recalling (52). It then follows straightforwardly that

$$\frac{\mu^{(t)}(a_1)}{\mu^{(t)}(a_2)} = \exp\left(\frac{w^{(t)}(a_1) - w^{(t)}(a_2)}{\tau}\right) \leq \exp\left(\frac{1 + \gamma\tau(\log |\mathcal{A}| + \log |\mathcal{B}|)}{(1-\gamma)\tau}\right)$$

for any $a_1, a_2 \in \mathcal{A}$. This allows us to show that

$$\min_{a \in \mathcal{A}} \mu^{(t)}(a) \geq \frac{1}{|\mathcal{A}| \exp\left(\frac{1 + \gamma\tau(\log |\mathcal{A}| + \log |\mathcal{B}|)}{(1-\gamma)\tau}\right)} \sum_{a \in \mathcal{A}} \mu^{(t)}(a) = \frac{1}{|\mathcal{A}| \exp\left(\frac{1 + \gamma\tau(\log |\mathcal{A}| + \log |\mathcal{B}|)}{(1-\gamma)\tau}\right)},$$

which gives

$$\begin{aligned} \|\log \mu^{(t)}\|_\infty & \leq \frac{1 + \gamma\tau(\log |\mathcal{A}| + \log |\mathcal{B}|)}{(1-\gamma)\tau} + \log |\mathcal{A}| \leq \frac{1}{(1-\gamma)\tau} + \frac{\log |\mathcal{A}| + \gamma \log |\mathcal{B}|}{1-\gamma} \\ & \leq \frac{2}{(1-\gamma)\tau}. \end{aligned}$$

□

E.2 PROOF OF LEMMA 12

We decompose the term $f_s(Q^{(t+1)}, \bar{\mu}^{(t+1)}, \bar{v}^{(t+1)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{v}^{(t)})$ as follows:

$$\begin{aligned} & f_s(Q^{(t+1)}, \bar{\mu}^{(t+1)}, \bar{v}^{(t+1)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{v}^{(t)}) \\ & = f_s(Q^{(t+1)}, \bar{\mu}^{(t+1)}, \bar{v}^{(t+1)}) - f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{v}^{(t+1)}) + f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{v}^{(t+1)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{v}^{(t)}) \end{aligned}$$

$$\begin{aligned}
&= \bar{\mu}^{(t+1)}(s)^\top \left(Q^{(t+1)}(s) - Q^{(t)}(s) \right) \bar{\nu}^{(t+1)}(s) \\
&\quad + f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) + f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t+1)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) \\
&\quad + \left[f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)}) + f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t+1)}) \right]
\end{aligned}$$

Note that $|\bar{\mu}^{(t+1)}(s)^\top (Q^{(t+1)}(s) - Q^{(t)}(s)) \bar{\nu}^{(t+1)}(s)| \leq \|Q^{(t+1)}(s) - Q^{(t)}(s)\|_\infty$. For the terms in the bracket, we have

$$\begin{aligned}
&\left| \left[f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)}) + f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t+1)}) \right] \right| \\
&= \left| (\bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s))^\top Q^{(t)}(s) (\bar{\nu}^{(t+1)}(s) - \bar{\nu}^{(t)}(s)) \right| \\
&\leq \frac{2}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t+1)} \| \bar{\zeta}^{(t)}).
\end{aligned}$$

It remains to bound the two difference terms $|f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)})|$ and $|f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t+1)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)})|$. To proceed, we show that

$$\begin{aligned}
&f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t)}) \\
&= \langle \bar{\mu}^{(t)}(s) - \bar{\mu}^{(t+1)}(s), Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) \rangle + \tau \mathcal{H}(\bar{\mu}^{(t)}(s)) - \tau \mathcal{H}(\bar{\mu}^{(t+1)}(s)) \\
&= \langle \bar{\mu}^{(t)}(s) - \bar{\mu}^{(t+1)}(s), Q^{(t)}(s)^\top \bar{\nu}^{(t)}(s) + (1 - \eta\tau)(Q^{(t)}(s)\bar{\nu}^{(t)}(s) - Q^{(t-1)}(s)\bar{\nu}^{(t-1)}(s)) \rangle \\
&\quad + \tau \mathcal{H}(\bar{\mu}^{(t)}(s)) - \tau \mathcal{H}(\bar{\mu}^{(t+1)}(s)) \\
&\quad - (1 - \eta\tau) \langle \bar{\mu}^{(t)}(s) - \bar{\mu}^{(t+1)}(s), Q^{(t)}(s)\bar{\nu}^{(t)}(s) - Q^{(t-1)}(s)\bar{\nu}^{(t-1)}(s) \rangle \\
&= -\frac{1}{\eta} \text{KL}_s(\bar{\mu}^{(t)} \| \bar{\mu}^{(t+1)}) - \frac{1 - \eta\tau}{\eta} \text{KL}_s(\bar{\mu}^{(t+1)} \| \bar{\mu}^{(t)}) \\
&\quad - (1 - \eta\tau) \langle \bar{\mu}^{(t)}(s) - \bar{\mu}^{(t+1)}(s), Q^{(t)}(s)\bar{\nu}^{(t)}(s) - Q^{(t-1)}(s)\bar{\nu}^{(t-1)}(s) \rangle \tag{55}
\end{aligned}$$

Here, the third step results from Lemma 17 along with (53). Recall from previous discussion (cf. (53)) that $\bar{\mu}^{(t+1)}(a|s) \propto \bar{\mu}^{(t)}(a|s) \exp(\eta w^{(t)}(s))$ with some $w^{(t)} \in \mathbb{R}^{|\mathcal{B}|}$ satisfying

$$\|w^{(t)}\|_\infty \leq \frac{6}{1-\gamma}.$$

We can ensure that $\|\eta w^{(t)}\|_\infty \leq 1/30$ with $\eta^{-1} \geq \frac{180}{1-\gamma}$, and the next lemma guarantees $\text{KL}_s(\bar{\mu}^{(t)} \| \bar{\mu}^{(t+1)}) \leq 2\text{KL}_s(\bar{\mu}^{(t+1)} \| \bar{\mu}^{(t)})$ in this case.

Lemma 16. *Let $w \in \mathbb{R}^{|\mathcal{A}|}$, $\pi, \pi' \in \Delta(\mathcal{A})$ satisfy, for each $a \in \mathcal{A}$, $\pi'(a) \propto \pi(a) \exp(w(a))$ with $\|w\|_\infty \leq \frac{1}{30}$. It holds that*

$$\text{KL}(\pi \| \pi') \leq 2\text{KL}(\pi' \| \pi).$$

Therefore, we can continue (55) by showing that

$$\begin{aligned}
&|f_s(Q^{(t)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)})| \\
&\leq \frac{1}{\eta} \text{KL}_s(\bar{\mu}^{(t)} \| \bar{\mu}^{(t+1)}) + \frac{1 - \eta\tau}{\eta} \text{KL}_s(\bar{\mu}^{(t+1)} \| \bar{\mu}^{(t)}) \\
&\quad + \|\bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s)\|_1 \|Q^{(t)}(s)\bar{\nu}^{(t)}(s) - Q^{(t-1)}(s)\bar{\nu}^{(t-1)}(s)\|_\infty \\
&\leq \frac{3}{\eta} \text{KL}_s(\bar{\mu}^{(t+1)} \| \bar{\mu}^{(t)}) + \|\bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s)\|_1 \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty \\
&\quad + \|Q^{(t)}(s)\|_\infty \|\bar{\mu}^{(t+1)}(s) - \bar{\mu}^{(t)}(s)\|_1 \|\bar{\nu}^{(t)}(s) - \bar{\nu}^{(t-1)}(s)\|_1 \\
&\leq \left(\frac{3}{\eta} + \frac{2}{1-\gamma} \right) \text{KL}_s(\bar{\mu}^{(t+1)} \| \bar{\mu}^{(t)}) + \frac{2}{1-\gamma} \text{KL}_s(\bar{\mu}^{(t)} \| \bar{\mu}^{(t-1)}) + \frac{6\eta}{1-\gamma} \|Q^{(t)}(s) - Q^{(t-1)}(s)\|_\infty
\end{aligned}$$

One can bound $|f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t+1)})|$ with similar argument. Putting all pieces together, we arrive at

$$|f_s(Q^{(t+1)}, \bar{\mu}^{(t+1)}, \bar{\nu}^{(t+1)}) - f_s(Q^{(t)}, \bar{\mu}^{(t)}, \bar{\nu}^{(t)})|$$

$$\begin{aligned} &\leq \left\| Q^{(t+1)}(s) - Q^{(t)}(s) \right\|_\infty + \left(\frac{3}{\eta} + \frac{4}{1-\gamma} \right) \text{KL}_s(\bar{\zeta}^{(t+1)} \parallel \bar{\zeta}^{(t)}) + \frac{2}{1-\gamma} \text{KL}_s(\bar{\zeta}^{(t)} \parallel \bar{\zeta}^{(t-1)}) \\ &\quad + \frac{12\eta}{1-\gamma} \left\| Q^{(t)}(s) - Q^{(t-1)}(s) \right\|_\infty. \end{aligned}$$

E.3 PROOF OF LEMMA 13

$$\begin{aligned} &\langle \bar{\nu}^{(t)}(s) - \nu_\tau^*(s), Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) \rangle - \tau \mathcal{H}(\bar{\nu}^{(t)}(s)) + \tau \mathcal{H}(\nu_\tau^*(s)) \\ &= \langle \bar{\nu}^{(t)}(s) - \nu^{(t)}(s), Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) - Q^{(t-1)}(s)^\top \bar{\mu}^{(t-1)}(s) \rangle \\ &\quad + \langle \bar{\nu}^{(t)}(s) - \nu^{(t)}(s), Q^{(t-1)}(s)^\top \bar{\mu}^{(t-1)}(s) \rangle - \tau \mathcal{H}(\bar{\nu}^{(t)}(s)) + \tau \mathcal{H}(\nu^{(t)}(s)) \\ &\quad + \langle \nu^{(t)}(s) - \nu_\tau^*(s), Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) \rangle - \tau \mathcal{H}(\nu^{(t)}(s)) + \tau \mathcal{H}(\mu_\tau^*(s)) \\ &= \langle \bar{\nu}^{(t)}(s) - \nu^{(t)}(s), Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) - Q^{(t-1)}(s)^\top \bar{\mu}^{(t-1)}(s) \rangle \\ &\quad + \frac{1-\eta\tau}{\eta} \text{KL}_s(\nu^{(t)} \parallel \nu^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\nu^{(t)} \parallel \bar{\nu}^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\bar{\nu}^{(t)} \parallel \nu^{(t-1)}) \\ &\quad + \frac{1-\eta\tau}{\eta} \text{KL}_s(\nu_\tau^* \parallel \nu^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\nu_\tau^* \parallel \nu^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\nu^{(t)} \parallel \nu^{(t-1)}) \\ &\leq \left\| \bar{\nu}^{(t)}(s) - \nu^{(t)}(s) \right\|_1 \left\| Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) - Q^{(t-1)}(s)^\top \bar{\mu}^{(t-1)}(s) \right\|_\infty \\ &\quad - \frac{1}{\eta} \text{KL}_s(\nu^{(t)} \parallel \bar{\nu}^{(t)}) - \frac{1-\eta\tau}{\eta} \text{KL}_s(\bar{\nu}^{(t)} \parallel \nu^{(t-1)}) + \frac{1-\eta\tau}{\eta} \text{KL}_s(\nu_\tau^* \parallel \nu^{(t-1)}) - \frac{1}{\eta} \text{KL}_s(\nu_\tau^* \parallel \nu^{(t)}). \end{aligned} \tag{56}$$

The second step results from the following three-point lemma:

Lemma 17 (Regularized 3-point lemma). *Let $x \in \Delta(\mathcal{A})$ be defined as*

$$x(a) \propto y(a)^{1-\eta\tau} \exp(-\eta w(a))$$

for some $w \in \mathbb{R}^{|\mathcal{A}|}$ and $y \in \Delta(\mathcal{A})$. It holds for all $z \in \Delta(\mathcal{A})$ that

$$\frac{\eta}{1-\eta\tau} \left[\langle x - z, w \rangle - \tau \mathcal{H}(x) + \tau \mathcal{H}(z) \right] = \text{KL}(z \parallel y) - \frac{1}{1-\eta\tau} \text{KL}(z \parallel x) - \text{KL}(x \parallel y).$$

We bound the first term in (56) as follows:

$$\begin{aligned} &\left\| \bar{\nu}^{(t)}(s) - \nu^{(t)}(s) \right\|_1 \left\| Q^{(t)}(s)^\top \bar{\mu}^{(t)}(s) - Q^{(t-1)}(s)^\top \bar{\mu}^{(t-1)}(s) \right\|_\infty \\ &\leq \left\| \bar{\nu}^{(t)}(s) - \nu^{(t)}(s) \right\|_1 \left(\left\| (Q^{(t)}(s) - Q^{(t-1)}(s))^\top \bar{\mu}^{(t-1)}(s) \right\|_\infty + \left\| Q^{(t)}(s) (\bar{\mu}^{(t)}(s) - \bar{\mu}^{(t-1)}(s)) \right\|_\infty \right) \\ &\leq \left\| \bar{\nu}^{(t)}(s) - \nu^{(t)}(s) \right\|_1 \left\| Q^{(t)}(s) - Q^{(t-1)}(s) \right\|_\infty + \frac{2}{1-\gamma} \left\| \bar{\nu}^{(t)}(s) - \nu^{(t)}(s) \right\|_1 \left\| \bar{\mu}^{(t)}(s) - \bar{\mu}^{(t-1)}(s) \right\|_1 \\ &\leq \frac{2\eta}{1-\gamma} \left\| Q^{(t)}(s) - Q^{(t-1)}(s) \right\|_\infty + \frac{1}{1-\gamma} \left[2 \left\| \bar{\nu}^{(t)}(s) - \nu^{(t)}(s) \right\|_1^2 \right. \\ &\quad \left. + \left\| \bar{\mu}^{(t)}(s) - \mu^{(t-1)}(s) \right\|_1^2 + \left\| \mu^{(t-1)}(s) - \bar{\mu}^{(t-1)}(s) \right\|_1^2 \right] \\ &\leq \frac{2\eta}{1-\gamma} \left\| Q^{(t)}(s) - Q^{(t-1)}(s) \right\|_\infty + \frac{4}{1-\gamma} \text{KL}_s(\nu^{(t)} \parallel \bar{\nu}^{(t)}) \\ &\quad + \frac{2}{1-\gamma} \text{KL}(\bar{\mu}^{(t)}(s) \parallel \mu^{(t-1)}(s)) + \frac{2}{1-\gamma} \text{KL}(\mu^{(t-1)}(s) \parallel \bar{\mu}^{(t-1)}(s)). \end{aligned}$$

Substitution of the above inequality into (56) completes the proof.

E.4 PROOF OF LEMMA 14

We have

$$\delta_{l,t} = \alpha_l \prod_{i=l+1}^t (1 - c_1 \alpha_i)$$

$$\begin{aligned}
&= \alpha_l \prod_{i=l+1}^t (1 - c_2 \alpha_i + (c_2 - c_1) \alpha_i) \\
&= \alpha_l (c_2 - c_1) \alpha_{l+1} \prod_{i=l+2}^t (1 - c_2 \alpha_i + (c_2 - c_1) \alpha_i) + \alpha_l (1 - c_2 \alpha_{l+1}) \prod_{i=l+2}^t (1 - c_2 \alpha_i + (c_2 - c_1) \alpha_i) \\
&= \alpha_l \sum_{i=l+1}^t (c_2 - c_1) \alpha_i \cdot \prod_{j=l+1}^i (1 - c_2 \alpha_j) \cdot \prod_{k=i+1}^t (1 - c_1 \alpha_k) + \alpha_l \prod_{i=l+1}^t (1 - c_2 \alpha_i) \\
&= (c_2 - c_1) \sum_{i=l+1}^t \xi_{l,i} \delta_{i,t} + \xi_{l,t}.
\end{aligned}$$

Rearranging terms,

$$\begin{aligned}
\sum_{i=l}^t \xi_{l,i} \delta_{i+1,t} &= \alpha_l \delta_{l+1,t} + \sum_{i=l+1}^t \xi_{l,i} \delta_{i+1,t} \\
&= \frac{\alpha_{l+1}}{1 - c_1 \alpha_{l+1}} \delta_{l,t} + \sum_{i=l+1}^t \xi_{l,i} \delta_{i,t} \cdot \frac{\alpha_{i+1}}{\alpha_i (1 - c_1 \alpha_{i+1})} \\
&\leq \delta_{l,t} + 2 \sum_{i=l+1}^t \xi_{l,i} \delta_{i,t} = \delta_{l,t} + \frac{2}{c_2 - c_1} (\delta_{l,t} - \xi_{l,t}) \leq \left(1 + \frac{2}{c_2 - c_1}\right) \delta_{l,t},
\end{aligned}$$

where the inequality is due to $\alpha_{l+1} \leq \alpha \leq 1/2$ and $1 - c_1 \alpha_l \geq 1/2$ for all $l \geq 1$.

E.5 PROOF OF LEMMA 16

Proof. For any $x > -1$, it holds that

$$\begin{aligned}
\log(1+x) &\leq x - \frac{x^2}{2} + \frac{x^3}{3} \\
&\leq x - \frac{x^2}{2} + \frac{|x^3|}{3} = x - \left(\frac{1}{2} - \frac{|x|}{3}\right) x^2,
\end{aligned}$$

and that

$$\begin{aligned}
\log(1+x) &\geq x - \frac{x^2}{2} + \frac{x^3}{3(1+x)^3} \\
&\geq x - \frac{x^2}{2} - \frac{|x^3|}{3(1+x)^3} = x - \left(\frac{1}{2} + \frac{|x|}{3(1+x)^3}\right) x^2.
\end{aligned}$$

Therefore, when $x > -\frac{1}{10}$, we have $(1+x)^3 > \frac{2}{3}$ and thus

$$x - \left(\frac{1}{2} + \frac{|x|}{2}\right) x^2 \leq \log(1+x) \leq x - \left(\frac{1}{2} - \frac{|x|}{3}\right) x^2.$$

Let c be a shorthand notation for $\|w\|_\infty$. The following lemma is standard (see, e.g., (Mei et al., 2020, Lemma 23), (Cen et al., 2021a, Lemma 3)), which ensures that $\|\log \pi - \log \pi'\|_\infty \leq 2c$.

Lemma 18. *Let $\pi, \pi' \in \Delta(\mathcal{A})$ satisfy $\pi(a) \propto \exp(\theta(a))$ and $\pi'(a) \propto \exp(\theta'(a))$ for some $\theta, \theta' \in \mathbb{R}^{|\mathcal{A}|}$. It holds that*

$$\|\log \pi - \log \pi'\|_\infty \leq 2\|\theta - \theta'\|_\infty.$$

Since $c < 1/30$, we have

$$\begin{aligned}
\left| \frac{\pi(a)}{\pi'(a)} - 1 \right| &= \left| \exp\left(\log \frac{\pi(a)}{\pi'(a)}\right) - \exp(0) \right| \leq |\log \pi(a) - \log \pi'(a)| \max\left\{1, \frac{\pi(a)}{\pi'(a)}\right\} \\
&\leq 2c \exp(|2c|) \leq 3c, \forall a \in \mathcal{A}.
\end{aligned}$$

Therefore, we can bound $\text{KL}(\pi \parallel \pi')$ as

$$\begin{aligned}
\text{KL}(\pi \parallel \pi') &= \sum_{a \in \mathcal{A}} \pi(a) \log \frac{\pi(a)}{\pi'(a)} \\
&\leq \sum_{a \in \mathcal{A}} \pi(a) \left(\frac{\pi(a)}{\pi'(a)} - 1 - \left(\frac{1}{2} - c \right) \left(\frac{\pi(a)}{\pi'(a)} - 1 \right)^2 \right) \\
&= \chi^2(\pi; \pi') - \left(\frac{1}{2} - c \right) \sum_{a \in \mathcal{A}} \pi(a) \left(\frac{\pi(a)}{\pi'(a)} - 1 \right)^2 \\
&\leq \chi^2(\pi; \pi') - \left(\frac{1}{2} - c \right) (1 - 3c) \sum_{a \in \mathcal{A}} \pi'(a) \left(\frac{\pi(a)}{\pi'(a)} - 1 \right)^2 \\
&= \left(1 - \left(\frac{1}{2} - c \right) (1 - 3c) \right) \chi^2(\pi; \pi').
\end{aligned} \tag{57}$$

On the other hand, we have

$$\begin{aligned}
\text{KL}(\pi' \parallel \pi) &= \sum_{a \in \mathcal{A}} \pi'(a) \log \frac{\pi'(a)}{\pi(a)} \\
&\geq \sum_{a \in \mathcal{A}} \pi'(a) \left(\frac{\pi'(a)}{\pi(a)} - 1 - \frac{1 + 3c}{2} \left(\frac{\pi'(a)}{\pi(a)} - 1 \right)^2 \right) \\
&= \chi^2(\pi'; \pi) - \frac{1 + 3c}{2} \sum_{a \in \mathcal{A}} \pi'(a) \left(\frac{\pi'(a)}{\pi(a)} - 1 \right)^2 \\
&\geq \chi^2(\pi'; \pi) - \frac{(1 + 3c)^2}{2} \sum_{a \in \mathcal{A}} \pi(a) \left(\frac{\pi'(a)}{\pi(a)} - 1 \right)^2 \\
&= \left(1 - \frac{(1 + 3c)^2}{2} \right) \chi^2(\pi'; \pi).
\end{aligned} \tag{58}$$

By definition, we have

$$\begin{aligned}
\chi^2(\pi; \pi') &= \sum_{a \in \mathcal{A}} \pi'(a) \left(\frac{\pi(a)}{\pi'(a)} - 1 \right)^2 = \sum_{a \in \mathcal{A}} \frac{(\pi(a) - \pi'(a))^2}{\pi'(a)} \\
&\leq \|\pi/\pi'\|_\infty \sum_{a \in \mathcal{A}} \frac{(\pi'(a) - \pi(a))^2}{\pi(a)} \\
&\leq (1 + 3c) \chi^2(\pi'; \pi).
\end{aligned} \tag{59}$$

Combining (57), (58) and (59) gives

$$\text{KL}(\pi \parallel \pi') \leq (1 + 3c) \cdot \frac{1 - (1/2 - c)(1 - 3c)}{1 - (1 + 3c)^2/2} \text{KL}(\pi' \parallel \pi).$$

It is straightforward to verify that the factor is less than 2 when $c \leq 1/30$. \square

E.6 PROOF OF LEMMA 17

Proof. We have

$$\begin{aligned}
\text{KL}(z \parallel y) &= -\mathcal{H}(z) + \mathcal{H}(y) - \langle z - y, \log y \rangle \\
&= -\mathcal{H}(z) + \mathcal{H}(x) - \langle z - x, \log y \rangle - \mathcal{H}(x) + \mathcal{H}(y) - \langle x - y, \log y \rangle \\
&= -\mathcal{H}(z) + \mathcal{H}(x) - \langle z - x, \log x \rangle - \mathcal{H}(x) + \mathcal{H}(y) - \langle x - y, \log y \rangle - \langle z - x, \log y - \log x \rangle \\
&= \text{KL}(z \parallel x) + \text{KL}(x \parallel y) - \frac{\eta}{1 - \eta\tau} \langle z - x, w + \tau \log x \rangle.
\end{aligned}$$

Rearranging terms gives

$$\frac{\eta}{1 - \eta\tau} \langle x - z, w \rangle = \text{KL}(z \parallel y) - \text{KL}(z \parallel x) - \text{KL}(x \parallel y) + \frac{\eta\tau}{1 - \eta\tau} \langle z - x, \log x \rangle.$$

Adding $\frac{\eta\tau}{1-\eta\tau}(-\mathcal{H}(x) + \mathcal{H}(z))$ to both sides, we are left with

$$\begin{aligned} \frac{\eta}{1-\eta\tau} \left[\langle x-z, w \rangle - \tau\mathcal{H}(x) + \tau\mathcal{H}(z) \right] &= \text{KL}(z \| y) - \text{KL}(z \| x) - \text{KL}(x \| y) \\ &\quad - \frac{\eta\tau}{1-\eta\tau} (-\mathcal{H}(z) + \mathcal{H}(x) - \langle z-x, \log x \rangle) \\ &= \text{KL}(z \| y) - \frac{1}{1-\eta\tau} \text{KL}(z \| x) - \text{KL}(x \| y). \end{aligned}$$

□

F FURTHER DISCUSSION REGARDING APPROXIMATE ALGORITHMS

In this section we verify the convergence of the proposed method equipped with inexact value updates in the infinite-horizon setting, where (10) in Algorithm 1 is replaced by

$$\begin{cases} Q^{(t+1)}(s, a, b) &= r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} [V^{(t)}(s')] + \alpha_t \delta^{(t)}(s, a, b) \\ V^{(t+1)}(s) &= (1 - \alpha_{t+1})V^{(t)}(s) \\ &\quad + \alpha_{t+1} [\bar{\mu}^{(t+1)}(s)^\top Q^{(t+1)}(s) \bar{\nu}^{(t+1)}(s) + \tau\mathcal{H}(\bar{\mu}^{(t+1)}(s)) - \tau\mathcal{H}(\bar{\nu}^{(t+1)}(s))] \end{cases}.$$

or equivalently

$$\begin{aligned} Q^{(t+1)}(s, a, b) &= (1 - \alpha_t)Q^{(t)}(s, a, b) \\ &\quad + \alpha_t \left[r(s, a, b) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a, b)} [\bar{\mu}^{(t)}(s)^\top Q^{(t)}(s) \bar{\nu}^{(t)}(s) + \tau\mathcal{H}(\bar{\mu}^{(t)}(s)) - \tau\mathcal{H}(\bar{\nu}^{(t)}(s))] + \delta^{(t)}(s, a, b) \right]. \end{aligned}$$

Here, $\delta(s, a, b)^{(t+1)} \in \mathbb{R}$ represents the error due to approximate evaluation. For simplicity we focus on the case where the policy update rules (9a), (9b) remain unchanged. The following theorems reveal that the algorithm converges linearly to the QRE until it reaches an error floor determined by $\|\delta^{(i)}\|_{\Gamma(\rho)}$:

Theorem 5. *With $0 < \eta \leq \frac{(1-\gamma)^3}{32000\mathcal{C}_\rho}$, and $\alpha_i = \eta\tau$, we have*

$$\begin{aligned} &\max \left\{ \text{KL}_\rho(\zeta_\tau^* \| \zeta^{(t)}), \frac{1}{2} \text{KL}_\rho(\zeta_\tau^* \| \bar{\zeta}^{(t)}), 3\eta \mathbb{E}_{s \sim \rho} [\|Q^{(t)}(s) - Q_\tau^*(s)\|_\infty] \right\} \\ &\leq \frac{3000}{(1-\gamma)^{2\tau}} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t + \frac{1500}{(1-\gamma)\tau} \max_{0 \leq i \leq t} \|\delta^{(i)}\|_{\Gamma(\rho)}. \end{aligned}$$

Theorem 6. *With $0 < \eta \leq \frac{(1-\gamma)^3}{32000\mathcal{C}_\rho}$, and $\alpha_i = \eta\tau$, we have*

$$\begin{aligned} \max_{s \in \mathcal{S}, \mu, \nu} \left(V_\tau^{\mu, \bar{\nu}^{(t)}}(s) - V_\tau^{\bar{\mu}^{(t)}, \nu}(s) \right) &\leq \frac{2\|1/\rho\|_\infty}{1-\gamma} \max \left\{ \frac{8}{(1-\gamma)^{2\tau}}, \frac{1}{\eta} \right\} \\ &\quad \cdot \left[\frac{3000}{(1-\gamma)^{2\tau}} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t + \frac{1500}{(1-\gamma)\tau} \max_{0 \leq i \leq t} \|\delta^{(i)}\|_{\Gamma(\rho)} \right], \end{aligned}$$

and

$$\begin{aligned} \max_{\mu, \nu} \left(V_\tau^{\mu, \bar{\nu}^{(t)}}(\rho) - V_\tau^{\bar{\mu}^{(t)}, \nu}(\rho) \right) &\leq \frac{2\mathcal{C}_{\rho, \tau}^\dagger}{1-\gamma} \max \left\{ \frac{8}{(1-\gamma)^{2\tau}}, \frac{1}{\eta} \right\} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t \\ &\quad \cdot \left[\frac{3000}{(1-\gamma)^{2\tau}} \left(1 - \frac{(1-\gamma)\eta\tau}{4}\right)^t + \frac{1500}{(1-\gamma)\tau} \max_{0 \leq i \leq t} \|\delta^{(i)}\|_{\Gamma(\rho)} \right]. \end{aligned}$$

We remark that $\|\delta^{(t)}\|_{\Gamma(\rho)}$ can be bounded either by ϵ_{stat} or $\mathcal{C}_\rho \epsilon_{\text{stat}}$ with evaluation error guarantee $\max_{s \in \mathcal{S}} \|\delta^{(t)}(s)\|_\infty \leq \epsilon_{\text{stat}}$ and $\mathbb{E}_{s \sim \rho} [\|\delta^{(t)}(s)\|_\infty] \leq \epsilon_{\text{stat}}$ respectively.

The remaining part of this section outlines the proof for the above Theorems. For simplicity, we only highlight the key difference from the previous proof due to evaluation error and omit the proof for corresponding lemmas. We first remark that Lemma 1 depends solely on the policy update rules and hence still holds. The error propagation of $\{\delta^{(l)}\}$ is captured by the following lemmas which parallels Lemma 2 and Lemma 16:

Lemma 19. *With $0 < \eta \leq \min\{(1 - \gamma)/180, (1 - \gamma)^2/48\}$, it holds for all $t \geq 1$ that*

$$\begin{aligned} \|Q^{(t+1)} - Q^{(t)}\|_{\Gamma(\rho)} &\leq \frac{1 + \gamma}{2} \sum_{l=1}^t \alpha_{l,t} \|Q^{(l)} - Q^{(l-1)}\|_{\Gamma(\rho)} + \frac{4\mathcal{C}_\rho}{\eta} \cdot \sum_{l=1}^t \alpha_{l,t} \text{KL}_\rho(\bar{\zeta}^{(l)} \parallel \bar{\zeta}^{(l-1)}) \\ &\quad + \alpha_t \|\delta^{(t)}\|_{\Gamma(\rho)} + \alpha_{t-1} \|\delta^{(t-1)}\|_{\Gamma(\rho)}. \end{aligned} \quad (60)$$

When $t = 0$, we have $\|Q^{(1)}(s) - Q^{(0)}(s)\|_{\Gamma(\rho)} \leq 2 + \alpha_0 \|\delta^{(0)}\|_{\Gamma(\rho)}$.

Lemma 20. *With $0 < \eta \leq (1 - \gamma)^2/16$, it holds for all $t \geq 1$ that*

$$\begin{aligned} &\|Q^{(t+1)} - Q_\tau^*\|_{\Gamma(\rho)} \\ &\leq \frac{1 + \gamma}{2} \cdot \sum_{l=0}^t \alpha_{l,t} \left(\|Q^{(l)} - Q_\tau^*\|_{\Gamma(\rho)} + \frac{2\eta}{1 - \gamma} \|Q^{(l)} - Q^{(l-1)}\|_{\Gamma(\rho)} \right) + 2\alpha_{0,t} + \alpha_t \|\delta^{(t)}\|_{\Gamma(\rho)} \end{aligned} \quad (61)$$

When $t = 0$, we have $\|Q^{(1)} - Q_\tau^*\|_{\Gamma(\rho)} \leq \frac{2\gamma}{1 - \gamma} + \alpha_0 \|\delta^{(0)}\|_{\Gamma(\rho)}$.

Following the similar argument in Lemma 4, we can show that

Lemma 21. *Under the assumption of Lemma 19 and 20, it holds for all $t \geq 0$ that*

$$\begin{aligned} &\sum_{l=0}^t \lambda_{l+1,t+1} \left[\eta \|Q_\tau^* - Q^{(l+1)}\|_{\Gamma(\rho)} + \frac{12\eta^2}{(1 - \gamma)^2} \|Q^{(l+1)} - Q^{(l)}\|_{\Gamma(\rho)} \right] \\ &\leq \frac{6250\eta\mathcal{C}_\rho}{(1 - \gamma)^3} \sum_{l=0}^{t-1} \lambda_{l+1,t+1} \text{KL}(\bar{\zeta}^{(l+1)} \parallel \bar{\zeta}^{(l)}) + \frac{550\eta}{(1 - \gamma)^2} \lambda_{0,t+1} + 60\eta \sum_{l=0}^t \lambda_{l+1,t+1} \alpha_l \|\delta^{(l)}\|_{\Gamma(\rho)}. \end{aligned}$$

With $\alpha_l = \eta\tau$ for $l \geq 1$, we have

$$\begin{aligned} \sum_{l=0}^t \lambda_{l+1,t+1} \alpha_l \|\delta^{(l)}\|_{\Gamma(\rho)} &\leq \lambda_{1,t+1} \|\delta^{(0)}\|_{\Gamma(\rho)} + \max_{1 \leq i \leq t} \|\delta^{(i)}\|_{\Gamma(\rho)} \sum_{l=1}^t \lambda_{l+1,t+1} \alpha_l \\ &\leq \lambda_{1,t+1} \|\delta^{(0)}\|_{\Gamma(\rho)} + \frac{4}{1 - \gamma} \max_{1 \leq i \leq t} \|\delta^{(i)}\|_{\Gamma(\rho)} \leq \frac{5}{1 - \gamma} \max_{0 \leq i \leq t} \|\delta^{(i)}\|_{\Gamma(\rho)}. \end{aligned}$$

It is then straightforward to put together the above lemmas in a similar way to the proof in Appendix A to obtain Theorem 5 and 6.

G FURTHER DISCUSSION REGARDING WEI ET AL. (2021B)

This section demonstrates how the last-iterate convergence result in Wei et al. (2021b, Theorem 2) in terms of the Euclidean distance to the set of NEs can be translated to that of the duality gap. Given any policy pair $\zeta = (\mu, \nu)$ and a NE $\zeta^* = (\mu^*, \nu^*)$, we can invoke performance difference lemma (44) and obtain:

$$\begin{aligned} V^{\mu, \nu}(\rho) - V^*(\rho) &= \frac{1}{1 - \gamma} \mathbb{E}_{s' \sim d_\rho^{\mu, \nu}} [\mu(s')^\top Q^*(s') \nu(s') - \mu^*(s')^\top Q^*(s') \nu^*(s')] \\ &\leq \frac{1}{1 - \gamma} \mathbb{E}_{s' \sim d_\rho^{\mu, \nu}} \left[\max_{\mu'} \mu'(s')^\top Q^*(s') \nu(s') - \mu^*(s')^\top Q^*(s') \nu^*(s') \right] \\ &= \frac{1}{1 - \gamma} \mathbb{E}_{s' \sim d_\rho^{\mu, \nu}} \left[\max_{\mu'} \mu'(s')^\top Q^*(s') \nu(s') - \max_{\mu'} \mu'(s')^\top Q^*(s') \nu^*(s') \right] \\ &\leq \frac{1}{1 - \gamma} \mathbb{E}_{s' \sim d_\rho^{\mu, \nu}} \left[\max_{\mu'} \mu'(s')^\top Q^*(s') (\nu(s') - \nu^*(s')) \right] \\ &\leq \frac{1}{(1 - \gamma)^2} \mathbb{E}_{s' \sim d_\rho^{\mu, \nu}} [\|\nu(s') - \nu^*(s')\|_1]. \end{aligned}$$

Setting μ to the best-response policy of ν , i.e., $\mu = \mu^\dagger(\nu) := \arg \max_{\mu} V^{\mu, \nu}(\rho)$, we get

$$\begin{aligned} \max_{\mu'} V^{\mu', \nu}(\rho) - V^*(\rho) &= V^{\mu^\dagger(\nu), \nu}(\rho) - V^*(\rho) \\ &\leq \frac{1}{(1-\gamma)^2} \mathbb{E}_{s' \sim d_{\rho}^{\mu^\dagger(\nu), \nu}} [\|\nu(s') - \nu^*(s')\|_1] \\ &\leq \frac{\|d_{\rho}^{\mu^\dagger(\nu), \nu}\|_{\infty}}{(1-\gamma)^2} \sum_{s \in \mathcal{S}} \|\nu(s) - \nu^*(s)\|_1. \end{aligned}$$

Similarly, we have

$$V^*(\rho) - \min_{\nu'} V^{\mu, \nu'}(\rho) \leq \frac{\|d_{\rho}^{\mu, \nu^\dagger(\mu)}\|_{\infty}}{(1-\gamma)^2} \sum_{s \in \mathcal{S}} \|\mu(s') - \mu^*(s')\|_1.$$

Taken together, the duality gap can be bounded by the policy's ℓ_1 distance to NE (μ^*, ν^*) as

$$\begin{aligned} \max_{\mu', \nu'} [V^{\mu', \nu}(\rho) - V^{\mu, \nu'}(\rho)] &\leq \frac{1}{(1-\gamma)^2} \sum_{s \in \mathcal{S}} (\|\nu(s') - \nu^*(s')\|_1 + \|\mu(s') - \mu^*(s')\|_1) \\ &\leq \frac{|\mathcal{S}|^{1/2} (|\mathcal{A}| + |\mathcal{B}|)^{1/2}}{(1-\gamma)^2} \left[\sum_{s \in \mathcal{S}} (\|\nu(s') - \nu^*(s')\|_2^2 + \|\mu(s') - \mu^*(s')\|_2^2) \right]^{1/2}, \end{aligned}$$

where the second step results from Cauchy-Schwarz inequality. Finally, recall from [Wei et al. \(2021b, Theorem 2\)](#) that it takes at most

$$\mathcal{O}\left(\frac{|\mathcal{S}|^2}{\eta^4 c^4 (1-\gamma)^4 \epsilon^2}\right)$$

iterations to ensure

$$\frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} (\|\nu(s') - \nu^*(s')\|_2^2 + \|\mu(s') - \mu^*(s')\|_2^2) \leq \epsilon^2,$$

with $\eta^2 = \mathcal{O}((1-\gamma)^5 |\mathcal{S}|^{-1})$. Putting pieces together and minimizing the bound over η , this leads to an iteration complexity of

$$\mathcal{O}\left(\frac{|\mathcal{S}|^5 (|\mathcal{A}| + |\mathcal{B}|)^{1/2}}{(1-\gamma)^{16} c^4 \epsilon^2}\right)$$

to achieve ϵ -NE in a last-iterate fashion.