

LEARNING IMPERFECT INFORMATION EXTENSIVE-FORM GAMES WITH LAST-ITERATE CONVERGENCE UNDER BANDIT FEEDBACK

Anonymous authors

Paper under double-blind review

ABSTRACT

We study learning the approximate Nash equilibrium (NE) policy profile in two-player zero-sum imperfect information extensive-form games (IIEFGs) with last-iterate convergence. The algorithms in previous works studying this problem either require full-information feedback or only have asymptotic convergence rates. In contrast, we study IIEFGs in the formulation of partially observable Markov games (POMGs) with the perfect-recall assumption and bandit feedback, where the knowledge of the game is not known a priori and only the rewards of the experienced information set and action pairs are revealed to the learners in each episode. Our algorithm utilizes a negentropy regularizer weighted by a virtual transition over information set-action space. By carefully designing the virtual transition together with the leverage of the entropy regularization technique, we prove that our algorithm converges to the NE of IIEFGs with a provable finite-time convergence rate of $\tilde{O}(k^{-1/8})$ with high probability under bandit feedback, thus answering the second question of [Fiegel et al. \(2023\)](#) affirmatively.

1 INTRODUCTION

In imperfect information games (IIGs), players operate with limited visibility into the game’s true state, necessitating strategic decision-making based on incomplete information. Notably, the concept of imperfect-information extensive-form games (IIEFGs), as introduced by [Kuhn \(1953\)](#), encapsulates both the intricacies of imperfect information and the sequential nature of players’ moves. This framework aptly represents a broad spectrum of real-world scenarios, such as Poker ([Heinrich et al., 2015](#); [Moravčík et al., 2017](#); [Brown & Sandholm, 2018](#)), Bridge ([Tian et al., 2020](#)), Scotland Yard ([Schmid et al., 2021](#)), and Mahjong ([Li et al., 2020](#); [Kurita & Hoki, 2021](#); [Fu et al., 2022](#)). Extensive research has been devoted to identifying the (approximate) Nash equilibrium (NE) ([Nash Jr, 1950](#)) within IIEFGs. Assuming the condition of perfect recall, where players possess the memory of past events and their implications, various methodologies have been employed to tackle these games. These include linear programming approaches ([Koller & Megiddo, 1992](#); [Von Stengel, 1996](#); [Koller et al., 1996](#)), which leverage mathematical optimization under full game knowledge, first-order optimization techniques ([Hoda et al., 2010](#); [Kroer et al., 2015](#); [2018](#); [Munos et al., 2020](#); [Lee et al., 2021](#); [Liu et al., 2022](#)), which iteratively refine strategies via repeated playthroughs of the games, and counterfactual regret minimization algorithms ([Zinkevich et al., 2007](#); [Lanctot et al., 2009](#); [Johnston et al., 2012](#); [Tammelin, 2014](#); [Schmid et al., 2019](#); [Burch et al., 2019](#); [Liu et al., 2022](#)), which adaptively adjust strategies based on counterfactual outcomes.

In practical scenarios, IIEFGs might involve large information set and action spaces, thwarting the application of linear programming approaches for *computing* the NE in IIEFGs. In this realm, the NE in IIEFGs is typically *learned* from random samples gathered through iterative playthroughs of the game, by Monte-Carlo counterfactual regret minimization (CFR) methods ([Lanctot et al., 2009](#); [Farina et al., 2020](#); [Farina & Sandholm, 2021](#)) or online mirror descent (OMD) and follow-the-regularized-leader (FTRL) frameworks ([Farina et al., 2021](#); [Kozuno et al., 2021](#); [Bai et al., 2022](#); [Fiegel et al., 2023](#)). Notably, [Bai et al. \(2022\)](#) devise an OMD-based approach incorporating “balanced exploration policies” to learn an ε -approximate NE with sample complexity of $\tilde{O}(H^3(XA + YB)/\varepsilon^2)$, where H is the horizon length, X, Y are the sizes of the information set space for the max- and min-player, and A and B are the sizes of the action space for the max- and

min-player. This upper bound is information-theoretically optimal with respect to all parameters except H , up to logarithmic factors. Building upon Bai et al. (2022), Fiegel et al. (2023) make further strides, refining the upper bound to $\tilde{O}(H(XA + YB)/\varepsilon^2)$ by harnessing FTRL with “balanced transitions”, achieving (nearly) optimal learning of IIEFGs in all parameters.

Despite the (nearly) optimal leaning of the ε -NE in IIEFGs by Bai et al. (2022); Fiegel et al. (2023), the algorithms in these works require to average all the policies generated during the running of the algorithms, so as to obtain the final policy profile with ε -NE guarantee. This is typically termed as the *average-iterate convergence*. However, in IIEFGs with large information set and action spaces, such an average operation over policy sets usually induces substantial storage and computation overhead. In cases when the policies in the games are approximated by nonlinear function approximation (e.g., neural networks), which has achieved great empirical success in recent years (Moravčík et al., 2017; Brown & Sandholm, 2018), computing the averaged policy even might be not feasible due to the nonlinearity of such function approximations. This motivates the studies of the learning algorithms with the *last-iterate convergence* guarantee of games including IIEFGs (Lin et al., 2020; Wei et al., 2021a;a; Lee et al., 2021; Cai et al., 2022; Abe et al., 2023; Feng et al., 2023; Cen et al., 2023; Liu et al., 2023). Specifically, Lee et al. (2021); Liu et al. (2023) establish algorithms for learning IIEFGs with last-iterate convergence rate of $\tilde{O}(1/k)$. However, the algorithms of Lee et al. (2021); Liu et al. (2023) require full-information feedback when learning IIEFGs, and thus can not be directly applied in practical cases when the knowledge of the games is not known a priori. The above considerations naturally motivate the following question:

Can we achieve last-iterate convergence for learning IIEFGs with bandit feedback?

Indeed, the same question has also been raised by Fiegel et al. (2023). In this work, we answer this question affirmatively. The main contributions of our work are summarized as follows:

- We propose the first algorithm that learns the approximate NE of IIEFGs with provable last-iterate convergence in the bandit feedback setting. In contrast with the vanilla negentropy regularizer (Lee et al., 2021) and the dilated negentropy regularizer (Lee et al., 2021; Liu et al., 2023) used by previous works to achieve the last-iterate convergence for IIEFGs with full-information feedback, our algorithm leverages the negentropy regularizer weighted by a virtual transition over info-set-action space to regularize the game. Via constructing the loss estimator regularized by such virtual transition weighted negentropy, our algorithm avoids directly regularizing the sequence-form representation of policies and results in a desirable contraction of the KL-divergence between probability measures over the information set-action space, instead of only obtaining the KL-divergence between the sequence-form representation of policies (see Section 4.1 and Section 5.1 for details). Besides, our algorithm does not require any communication or coordination between the two players and is model-free, without requiring the knowledge of the underlying state transition probabilities and the reward functions.
- To efficiently bound the stability term in the one-step analysis of OMD with Bregman divergence induced by the virtual transition weighted negentropy regularizer, we design a virtual transition over the information set-action space that maximizes the minimum visitation probability of all the information sets (see Section 4.1 for more elaboration on this). With such a virtual transition, we finally prove that our algorithm obtains the finite-time last-iterate convergence rate for learning IIEFGs in the bandit feedback setting of $\tilde{O}((X + Y)[(XA + YB)^{1/2} + (X + Y)^{1/4}H]k^{-1/8})$ with high probability (for large enough k), where H is the horizon length, X and Y are the size of the information set spaces of the max- and min-player, A and B are the size of the action spaces of the max- and min-player, and k is the number of episodes. The methodology of our algorithm’s analysis is inspired by the last-iterate convergence learning of the matrix games and the (fully observable) Markov games of Cai et al. (2023), but we provide a refined analysis specifically for IIEFGs to further sharpen the dependence on the parameters when deriving the final convergence rate (see Section 5.1 for details).
- When only obtaining the expected convergence rate is desired, our algorithm can generate a policy profile converging to the NE with a rate of $\tilde{O}((X + Y)[(X^2A + Y^2B)^{1/2} + (X + Y)^{1/4}H]k^{-1/6})$ in expectation. For the problem of learning the NE of IIEFGs in the bandit-feedback setting, we provide an $\Omega(\sqrt{XA + YB}k^{-1/2})$ lower bound of the last-iterate convergence rate.

2 RELATED WORKS

2.1 PARTIALLY OBSERVABLE MARKOV GAMES (POMGs)

With perfect information, learning Markov games (MGs) can be traced back to the seminal work of Littman & Szepesvári (1996) and has since garnered extensive research attention (Littman, 2001; Greenwald & Hall, 2003; Hu & Wellman, 2003; Hansen et al., 2013; Sidford et al., 2018; Lagoudakis & Parr, 2002; Pérolat et al., 2015; Fan et al., 2020; Jia et al., 2019; Cui & Yang, 2021; Zhang et al., 2021; Bai & Jin, 2020; Liu et al., 2021; Zhou et al., 2021; Song et al., 2022; Li et al., 2022; Xiong et al., 2022; Wang et al., 2023; Cui et al., 2023). In scenarios where only imperfect information is available yet the complete knowledge of the game (transitions and rewards) is known, existing research can be categorized into three primary streams. The first stream leverages sequence-form representation of policies to recast the problem as a linear program (Koller & Megiddo, 1992; Von Stengel, 1996; Koller et al., 1996). The second stream translates the problem into a minimax optimization problem and explores first-order algorithms, as exemplified in (Hoda et al., 2010; Kroer et al., 2015; 2018; Munos et al., 2020; Lee et al., 2021; Liu et al., 2022). Lastly, the third stream addresses the problem through CFR, minimizing counterfactual regrets locally within each information set (Zinkevich et al., 2007; Lanctot et al., 2009; Johanson et al., 2012; Tammelin, 2014; Schmid et al., 2019; Burch et al., 2019; Liu et al., 2022).

In the realm where the knowledge of the game is either unknown or only partially accessible, the Monte-Carlo CFR algorithm introduced by Lanctot et al. (2009) pioneers the achievement of the first ε -NE result. This framework has been further generalized and extended by Farina et al. (2020); Farina & Sandholm (2021). Additionally, another line of research focuses on integrating OMD and FTRL frameworks with importance-weighted loss estimators (Farina et al., 2021; Kozuno et al., 2021; Bai et al., 2022; Fiegel et al., 2023). Remarkably, Bai et al. (2022) achieve an ε -approximate NE with sample complexity of $\tilde{O}(H^3(XA + YB)/\varepsilon^2)$ by employing a “balanced” dilated KL-divergence as the distance metric. Building upon this concept, Fiegel et al. (2023) utilize “balanced transitions” and attain a (nearly) optimal sample complexity of $\tilde{O}(H(XA + YB)/\varepsilon^2)$, which matches the information-theoretic lower bound up to logarithmic factors. However, we note that all the algorithms in existing works studying POMGs with bandit feedback only have *average-iterate convergence* guarantees, while we aim to establish the *last-iterate convergence* guarantee.

2.2 LAST-ITERATE CONVERGENCE LEARNING IN GAMES

With full-information feedback, learning in games with last-iterate convergence guarantee has been investigated in strongly monotone games (Mokhtari et al., 2020; Jordan et al., 2024), monotone games (Golowich et al., 2020; Cai et al., 2022; Gorbunov et al., 2022; Cai & Zheng, 2023), Markov games (Cen et al., 2021; 2023), and IIEFGs (Lee et al., 2021; Liu et al., 2023; Bernasconi et al., 2024).

Recently, motivated by the fact that it might be restrictive to require full knowledge of the (noisy) gradient as in the full-information feedback setting, a growing body of works has studied learning in games with last-iterate convergence guarantee in the bandit feedback setting including strongly monotone games (Bravo et al., 2018; Lin et al., 2021) (Bravo et al., 2018; Hsieh et al., 2019; Lin et al., 2021; Drusvyatskiy et al., 2022; Huang & Hu, 2023), matrix games (Cai et al., 2023) and Markov games (Wei et al., 2021b; Chen et al., 2022; Cai et al., 2023). However, the algorithm of Wei et al. (2021b) needs coordinated updates and some prior knowledge of the game, and the algorithm of Chen et al. (2022) requires the players to inform the opponent about the entropy of their own policies. Amongst these works, Cai et al. (2023) remove all the coupling requirements, achieving last-iterate convergences of $\tilde{O}(k^{-1/8})$ for matrix games and of $\tilde{O}(k^{-1/9+\varepsilon})$ for any $\varepsilon > 0$ for irreducible Markov games. We note that all existing works study fully-observable Markov games, while we aim to establish uncoupled algorithms for learning IIEFGs in the formulation of partially-observable Markov games, without requiring the knowledge of the games.

3 PRELIMINARIES

For ease of exposition, we consider IIEFGs in the formulation of POMGs and introduce the preliminaries of them in this section, following previous works (Kozuno et al., 2021; Bai et al., 2022).

Partially Observable Markov Games We study episodic, finite-horizon, two-player zero-sum POMGs, denoted by $\text{POMG}(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mathcal{A}, \mathcal{B}, H, \mathbb{P}, r)$, in which

- H is the horizon length;
- $\mathcal{S} = \bigcup_{h \in [H]} \mathcal{S}_h$ is the finite state space with \mathcal{S}_h as the state space at step h . $S = \sum_{h=1}^H S_h$ is the size of \mathcal{S} where $|\mathcal{S}_h| = S_h, \forall h \in [H]$;
- $\mathcal{X} = \bigcup_{h \in [H]} \mathcal{X}_h$ is the finite space of information sets (short for *infosets* in the following) for the max-player, where $\mathcal{X}_h = \{x(s) : s \in \mathcal{S}_h\}$ is the set of the infosets at step h with $x : \mathcal{S} \rightarrow \mathcal{X}$ as the emission function. $X = \sum_{h=1}^H X_h$ is the size of \mathcal{X} with $|\mathcal{X}_h| = X_h$. The finite space of infosets $\mathcal{Y} = \bigcup_{h \in [H]} \mathcal{Y}_h$ for the min-player and its size are defined analogously;
- \mathcal{A} with $|\mathcal{A}| = A$ and \mathcal{B} with $|\mathcal{B}| = B$ are the finite action spaces for the max-player and min-player, respectively;
- $\mathbb{P} = \{p_0(\cdot) \in \Delta_{\mathcal{S}_1}\} \cup \{p_h(\cdot | s_h, a_h, b_h) \in \Delta_{\mathcal{S}_{h+1}}\}_{(s_h, a_h, b_h) \in \mathcal{S}_h \times \mathcal{A} \times \mathcal{B}, h \in [H-1]}$ are the state transition probabilities, where $p_0(\cdot)$ is the probability distribution of initial states, $p_h(s_{h+1} | s_h, a_h, b_h)$ is the probability of transitioning to the next state s_{h+1} conditioned on (s_h, a_h, b_h) at step h , and $\Delta_{\mathcal{S}_h}$ denotes the probability simplex over \mathcal{S}_h ;
- $r = \{r_h(s_h, a_h, b_h) \in [0, 1]\}_{(s_h, a_h, b_h) \in \mathcal{S}_h \times \mathcal{A} \times \mathcal{B}, h \in [H]}$ are the (randomized) reward functions with $\bar{r}_h(s_h, a_h, b_h)$ as mean for each $r_h(s_h, a_h, b_h)$.

Learning Protocol Define the max-player’s stochastic policy as $\mu = \{\mu_h^k\}_{h \in [H]}$, where $\mu_h^k : \mathcal{X}_h \rightarrow \Delta_{\mathcal{A}}$ denotes the policy at step h during episode k . The set of all such policies for the max-player is denoted by Π_{\max} . Analogously, the min-player’s stochastic policy is specified as $\nu = \{\nu_h^k\}_{h \in [H]}$, with $\nu_h^k : \mathcal{Y}_h \rightarrow \Delta_{\mathcal{B}}$ being the policy at step h during episode k , and the set of all min-player policies is denoted by Π_{\min} . The game proceeds in a finite number of episodes. At the commencement of episode k , the max-player selects a stochastic policy $\mu^k \in \Pi_{\max}$, while the min-player chooses $\nu^k \in \Pi_{\min}$. Meanwhile, an initial state s_1^k is sampled from the distribution $p_0(\cdot)$ by the environment. During each step h within an episode, the max-player and min-player observe their respective infosets $x_h^k := x(s_h^k)$ and $y_h^k := y(s_h^k)$, but they do not directly observe the underlying state s_h^k . Given x_h^k , the max-player samples and executes an action $a_h^k \sim \mu_h^k(\cdot | x_h^k)$, while the min-player concurrently takes an action $b_h^k \sim \nu_h^k(\cdot | y_h^k)$. Upon taking these actions, the max-player and min-player receive rewards $r_h^k := r_h(s_h^k, a_h^k, b_h^k)$ and $-r_h^k$, respectively. Subsequently, the game transitions to the next state $s_{h+1}^k \sim p_h(\cdot | s_h^k, a_h^k, b_h^k)$. The k -th episode will terminate after actions a_H^k and b_H^k are taken conditioned on x_H^k and y_H^k .

Perfect Recall and Tree Structure Following prior works (Kozuno et al., 2021; Bai et al., 2022; Fiegel et al., 2023), we assume that the POMGs adhere to the *tree structure* and the *perfect recall* condition, as defined by Kuhn (1953). Explicitly, the tree structure signifies that for any step $h = 2, \dots, H$ and state $s_h \in \mathcal{S}_h$, there exists a *unique* path $(s_1, a_1, b_1, \dots, s_{h-1}, a_{h-1}, b_{h-1})$ culminating in s_h . The perfect recall condition, meanwhile, is fulfilled for both players, implying that for any $h = 2, \dots, H$ and any infoset $x_h \in \mathcal{X}_h$ of the max-player (analogously for the min-player), there exists a *unique* history $(x_1, a_1, \dots, x_{h-1}, a_{h-1})$ leading to x_h . Furthermore, we introduce the notation $C_{h'}(x_h, a_h) \subset \mathcal{X}_{h'}$ to represent the set of descendants of the infoset-action pair (x_h, a_h) at step $h' \geq h$. Also, we define $C_{h'}(x_h) := \bigcup_{a_h \in \mathcal{A}} C_{h'}(x_h, a_h)$ as the union of descendants across all actions at x_h , and for convenience, let $C(x_h, a_h) := C_{h+1}(x_h, a_h)$ signify the immediate descendants at the subsequent step.

Sequence-form Representations For any pair of product policies (μ, ν) , the tree structure and the perfect recall condition facilitate the *sequence-form representation* of the reaching probability for the state-action tuple (s_h, a_h, b_h) :

$$\mathbb{P}^{\mu, \nu}(s_h, a_h, b_h) = p_{1:h}(s_h) \mu_{1:h}(x(s_h), a_h) \nu_{1:h}(y(s_h), b_h), \quad (1)$$

where $p_{1:h}(s_h) = p_0(s_1) \prod_{h'=1}^{h-1} p_{h'}(s_{h'+1} | s_{h'}, a_{h'}, b_{h'})$ denotes the sequence-form transition probability, and $\mu_{1:h}(x_h, a_h) := \prod_{h'=1}^h \mu_{h'}(a_{h'} | x_{h'})$ and $\nu_{1:h}(y_h, b_h) := \prod_{h'=1}^h \nu_{h'}(b_{h'} | y_{h'})$ represent the sequence-form policies of the max- and min player, respectively. Under the sequence-form representation, we adopt a slight abuse of notation for μ and ν by interpreting them as $\mu = \{\mu_{1:h}\}_{h \in [H]}$

and $\nu = \{\nu_{1:h}\}_{h \in [H]}$.¹ Furthermore, it is clear that Π_{\max} constitutes a convex compact subspace of \mathbb{R}^{XA} that adheres to the constraints $\mu_{1:h}(x_h, a_h) \geq 0$ and $\sum_{a_h \in \mathcal{A}} \mu_{1:h}(x_h, a_h) = \mu_{1:h-1}(x_{h-1}, a_{h-1})$, where (x_{h-1}, a_{h-1}) is such that $x_h \in C(x_{h-1}, a_{h-1})$ (with the understanding that $\mu_{1:0}(x_0, a_0) = 1$ as a base case).

Learning Objective In this work, we consider the learning objective of finding an approximate NE of the POMG. Specifically, for any $\varepsilon \geq 0$, an ε -approximate NE is a pair of product policy (μ, ν) satisfying

$$\text{NEGap}(\mu, \nu) := \max_{\mu^\dagger \in \Pi_{\max}} V^{\mu^\dagger, \nu} - \min_{\nu^\dagger \in \Pi_{\min}} V^{\mu, \nu^\dagger} \leq \varepsilon, \quad (2)$$

where $V^{\mu, \nu} = \mathbb{E}_{\mu, \nu} \left[\sum_{h=1}^H r_h(s_h, a_h, b_h) \right]$ the value function of (μ, ν) with the expectation taken over the randomness of the product policy pair (μ, ν) and the environment. It is known that using regret to NE conversion, an approximate NE can be obtained by averaging all the policies $\{\mu\}_{k=1}^K$ of the max-player generated by an algorithm with sublinear regret (similarly for the min-player) to obtain the average policy pair $(\bar{\mu}, \bar{\nu})$ (see, e.g., Theorem 1 of [Kozuno et al. \(2021\)](#)). This is the so-called *average-iterate convergence* of learning NE. In this work, we are interested in finding the ε -NE with the (finite-time) *last-iterate convergence* guarantee; that is, the algorithm is required to generate an approximate NE policy profile (μ^k, ν^k) such that $\text{NEGap}(\mu^k, \nu^k) \leq \varepsilon_k$ in each episode for finite-time k .

Information Available to the Players In this work, we consider learning POMGs in the bandit feedback setting, where in each episode k , the max-player only observes her experienced trajectory $(x_1^k, a_1^k, r_1^k, \dots, x_H^k, a_H^k, r_H^k)$ of infosets, actions, and rewards, but not the underlying states or the opponent’s infosets and actions. Additionally, the max-player does not have knowledge about the policies adopted by the min-player and also can not receive any information from the min-player and vice versa. Besides, there is no shared randomness between both players; that is, the algorithms of both players need to be fully uncoupled from each other.

Additional Notations We slightly abuse the notation to view x_h as the set $\{s \in \mathcal{S}_h : x(s) = x_h\}$, when writing $s \in x_h$. Given sequence-form representations, for any $\mu \in \Pi_{\max}$ and a sequence of functions $f = (f_h)_{h \in [H]}$ with $f_h : \mathcal{X}_h \times \mathcal{A} \rightarrow \mathbb{R}$, we define $\langle \mu, f \rangle := \sum_{h \in [H], (x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}} \mu_{1:h}(x_h, a_h) f_h(x_h, a_h)$. We denote by \mathcal{F}^k the σ -algebra generated by the random variables $\{(s_h^t, a_h^t, b_h^t, r_h^t)\}_{h \in [H], t \in [k]}$. For brevity, we abbreviate the conditional expectation $\mathbb{E}[\cdot | \mathcal{F}^k]$ as $\mathbb{E}^k[\cdot]$. Throughout this paper, the notation $\tilde{\mathcal{O}}(\cdot)$ suppresses all logarithmic factors.

4 ALGORITHM

In this section, we introduce the proposed algorithm, detailed in [Algorithm 1](#).

4.1 FROM SEQUENCE-FORM REPRESENTATIONS TO PROBABILITY MEASURES OVER INFOSET-ACTION SPACE

With sequence-form representations, we first reformulate the IIEFG into the following bilinear game:

$$f(\mu, \nu) = \mu^\top \mathbf{G} \nu, \quad (4)$$

where $\mathbf{G} \in \mathbb{R}^{XA \times YB}$ is the loss matrix with $\mathbf{G}[(x_h, a_h), (y_h, b_h)] = \sum_{s_h \in x_h \cap y_h} p_{1:h}(s_h) (1 - r_h(s_h, a_h, b_h))$. In this manner, the learning objective is equivalent to finding (μ, ν) such that $\text{NEGap}(\mu, \nu) = \sup_{\mu^\dagger \in \Pi_{\max}} f(\mu, \nu^\dagger) - f(\mu^\dagger, \nu) \leq \varepsilon$. At a high level, we apply the entropy regularizing technique to perturb the bilinear form of the game, as defined in [Eq. \(4\)](#), into a strongly convex-strongly concave structure, ensuring convergence to both the NE of the perturbed game (and thus the NE of the original game in [Eq. \(4\)](#)). This approach

¹The sequence-form representation of policies is defined in a top-down manner and is equivalent to the “treeplex” space of policies defined in a bottom-up manner (see, e.g., [Lee et al. \(2021\)](#)).

Algorithm 1 OMD with Virtual Transition Weighted Negentropy Regularization (max-player)

-
- 1: **Input:** $\eta_k = k^{-\alpha_\eta}, \gamma_k = k^{-\alpha_\gamma}, \varepsilon_k = k^{-\alpha_\varepsilon}$.
 - 2: **Initialize:** $\mu_1(a_h|x_h) = \frac{1}{A}, \forall (x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}, \forall h \in [H]$. Set p^x computed by Algorithm 2.
 - 3: **for** $k = 1, \dots$, **do**
 - 4: **for** $h = 1, \dots, H$ **do**
 - 5: Observes x_h^k , executes $a_h^k \sim \mu_h^k(\cdot|x_h^k)$ and receives r_h^k .
 - 6: For all $(x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}$, set entropy regularized loss estimator as
- $$\hat{\ell}_h^k(x_h, a_h) = \frac{\mathbb{I}_h^k\{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} (1 - r_h^k) + \varepsilon_k \cdot p_{1:h}^x(x_h) \log[p_{1:h}^x \cdot \mu_{1:h}^k](x_h, a_h).$$
- 7: **end for**
 - 8: Update policy
-

$$\mu^{k+1} = \arg \min_{\mu \in \Pi_{\max}^{k+1}} \eta_k \langle \mu, \hat{\ell}^k \rangle + D_\psi(\mu, \mu^k), \quad (3)$$

where $\Pi_{\max}^{k+1} = \{\mu \in \Pi_{\max} : \mu(a_h|x_h) \geq \frac{1}{A(k+1)}, \forall (x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}, \forall h \in [H]\}$.

- 9: **end for**
-

Algorithm 2 Computing virtual transition p^x (max-player)

-
- 1: **Input:** Game tree structure of $\mathcal{X} \times \mathcal{A}$.
 - 2: **Initialization:** Sequence-form representation of virtual transition $q \in \mathbb{R}^X$; array of maximized number of descendant infoset $c \in \mathbb{R}^X, d \in \mathbb{R}^{XA}$. For all x_H in \mathcal{X}_H , set $c[x_H] = 1$.
 - 3: **for** $h = H - 1$ to 1 **do**
 - 4: **for** x_h in \mathcal{X}_h **do**
 - 5: **for** a_h in \mathcal{A} **do**
 - 6: Compute $d[x_h, a_h] = \sum_{x_{h+1} \in C(x_h, a_h)} c[x_{h+1}]$.
 - 7: **end for**
 - 8: Compute $c[x_h] = \max_{a \in \mathcal{A}} d[x_h, a]$.
 - 9: **end for**
 - 10: **end for**
 - 11: **for** x_1 in \mathcal{X}_1 **do**
 - 12: Compute $q_{1:1}(x_1) = \frac{c[x_1]}{\sum_{x_1 \in \mathcal{X}_1} c[x_1]}$.
 - 13: **end for**
 - 14: **for** $h = 1$ to $H - 1$ **do**
 - 15: **for** x_h, a_h in $\mathcal{X}_h \times \mathcal{A}$ **do**
 - 16: **for** x_{h+1} in $C(x_h, a_h)$ **do**
 - 17: Compute $q_{1:h+1}(x_{h+1}) = q_{1:h}(x_h) \cdot \frac{c[x_{h+1}]}{\sum_{x_{h+1} \in C(x_h, a_h)} c[x_{h+1}]}$.
 - 18: **end for**
 - 19: **end for**
 - 20: **end for**
 - 21: **return** q .
-

builds upon previous research that has explored last-iterate convergence learning in Markov games with full-information feedback (Cen et al., 2021; Chen et al., 2022; Cen et al., 2023), matrix games and Markov games with bandit feedback (Cai et al., 2023), and IIEFGs with full-information feedback (Liu et al., 2023). Specifically, we consider the following perturbed game as a surrogate:

$$f_k(\mu, \nu) = \mu^\top \mathbf{G}\nu + \varepsilon_k \psi(\mu) - \varepsilon_k \psi(\nu), \quad (5)$$

where ψ is some strongly convex regularizer to be used in OMD and $\varepsilon_k > 0$ serves as the knob to control the strength of the entropy regularization in episode k . Intuitively, due to the strongly convex-strongly concave property of the perturbed game, one is able to find the approximate NE of it with last-iterate convergence using OMD. On the other hand, by gradually decreasing ε_k to be moderately small, the approximate NE of the perturbed game in Eq. (5) will also serve as an approximate NE of the original game in Eq. (4).

The crucial aspect lies in selecting an appropriate regularizer ψ . Initially, the first candidate that might come to mind is the utilization of the vanilla negentropy regularizer $\psi(\mu) = \sum_{h,x_h,a_h} \mu_{1:h}(x_h, a_h) \log \mu_{1:h}(x_h, a_h)$, which has been utilized to achieve the last-iterate convergence for IIEFGs with full-information feedback (Lee et al., 2021) and matrix games, the special case of IIEFGs, with bandit feedback (Cai et al., 2023). However, in IIEFGs with bandit feedback, though using the vanilla negentropy regularizer results in a convergence of the Bregman divergence, it is generally hard to control the NE gap since it directly regularizes the sequence-form representation policies. The other natural approach is considering using the dilated negentropy $\psi(\mu) = \sum_{h,x_h,a_h} \mu_{1:h}(x_h, a_h) \log \left(\frac{\mu_{1:h}(x_h, a_h)}{\mu_{1:h}(x_h)} \right)$ (Kroer et al., 2015; Kozuno et al., 2021). Indeed, the dilated negentropy has also been used to achieve the last-iterate convergence of the IIEFGs with full-information feedback (Lee et al., 2021; Liu et al., 2023; Bernasconi et al., 2024). However, in contrast with the full-information feedback setting, leveraging the entropy regularization technique to obtain the finite-time convergence guarantee in the bandit feedback setting requires the probability of selecting each action a_h given each infoset x_h being lower bounded to prevent the stability term in the analysis of OMD from being prohibitively largely. This essentially requires constraining the optimization of OMD onto a subset of the entire space of the sequence-form representations of policies Π_{\max} . Nevertheless, this will also make the stability term of OMD using the dilated negentropy in conjunction with the regularization technique hard to control, as bounding the stability term of the OMD with dilated negentropy critically relies upon its closed-form update solution (see, e.g., Lemma 7 of Kozuno et al. (2021)), which no longer holds in the case where the policy update of OMD is constrained onto a subset of Π_{\max} .

To cope with the above difficulties, we instead consider using the negentropy regularizer weighted by a kind of *virtual transition* p^x over the infoset-action space $\mathcal{X} \times \mathcal{A}$:

$$\psi_{p^x}(\mu) = \sum_{h,x_h,a_h} p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)),$$

where $p_h^x(\cdot|x_h, a_h) \in \Delta_{\mathcal{C}(x_h, a_h)}$ is a transition probability over $\mathcal{X}_h \times \mathcal{A} \times \mathcal{X}_{h+1}$ and $p_{1:h}^x(x_h) = p_0^x(x_1) \prod_{h'=1}^{h-1} p_{h'+1}^x(x_{h'+1}|x_{h'}, a_{h'})$ is its sequence-form representation. Note that $p_h^x(x_{h+1}|x(s_h), a_h)$ is not necessarily to be the true transition probability $\mathbb{P}^{\mu^k, \nu^k}(x_{h+1}|x(s_h), a_h) = \sum_{s_{h+1} \in \mathcal{X}_{h+1}, b_h \in \mathcal{B}} p(s_{h+1}|s_h, a_h, b_h) \nu^k(b_h|y(s_h))$ experienced by the max-player in episode k . Also, notice that $\psi_{p^x}(\cdot)$ is dependent on the chosen virtual transition p^x and we drop the dependence in the subscript of $\psi_{p^x}(\cdot)$ on p^x when the context is clear for brevity. We remark that similar ideas leveraging negentropy weighted by the transition over infoset-action space have also been exploited by Bai et al. (2022); Fiegel et al. (2023). However, we would like to underscore that the design of our virtual transition p^* over infoset-action space is different from those of Bai et al. (2022); Fiegel et al. (2023) and we aim to establish the last-iterate convergence of IIEFGs while they can only guarantee the average-iterate convergence, necessitating different theoretical analysis. Besides, one can see that the constructed virtual transition p^x is well-defined by the perfect recall condition and $p_{1:h}^x \cdot \mu_{1:h}$ with $[p_{1:h}^x \cdot \mu_{1:h}](x_h, a_h) = p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)$ is a probability measure over the infoset-action space $\mathcal{X}_h \times \mathcal{A}$ at step h . Therefore, we actually regularize the probability measures over $\mathcal{X}_h \times \mathcal{A}$ instead of directly regularizing the sequence-form representation μ , which tackles the difficulties of using the vanilla negentropy and the dilated negentropy as mentioned above. The other nice property of virtual transition weighted negentropy is that $D_\psi(\mu_1, \mu_2) = \text{KL}(p^x \mu_1, p^x \mu_2)$, facilitating bounding the final NE gap as we shall see in Section 5.1.

With regularizer ψ specified, the derivative of $f_k(\mu, \nu)$ w.r.t. $\mu(x_h, a_h)$ is $\frac{\partial f_k(\mu, \nu)}{\partial \mu_{1:h}(x_h, a_h)} = \mathcal{G}\nu[(x_h, a_h)] + \varepsilon_k \cdot p_{1:h}^x(x_h) [\log[p_{1:h}^x \cdot \mu_{1:h}](x_h, a_h) + 1]$. Since $[p_{1:h}^x \cdot \mu_{1:h}] \in \Delta_{\mathcal{X}_h \times \mathcal{A}}$ for any μ , the constant 1 in the above display does not affect the optimization of OMD. On the other hand, in the bandit feedback setting, an (optimistically biased) loss estimate $\frac{\mathbb{I}\{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} (1 - r_h^k)$ of $\mathcal{G}\nu[(x_h, a_h)]$ in episode k is constructed (Kozuno et al., 2021), where $\gamma_k > 0$ is the implicit exploration parameter (Neu, 2015). This specifies the final entropy regularized loss estimator used by Algorithm 1 on Line 6.

With the constructed loss estimator, Algorithm 1 then uses OMD to update policy. Since now the entropy regularized loss estimator is considered, the variance of the loss estimator will be prohibitively large if running OMD on the entire space of the sequence-form representations Π_{\max} , eventually leading to an unbounded stability term of OMD. Hence we constrain the feasible set of the OMD

as a subset Π_{\max}^{k+1} of Π_{\max} , where each $\mu \in \Pi_{\max}^{k+1}$ satisfying $\mu(a_h|x_h)$ is lower bounded for all $(x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}$ and $h \in [H]$ (Line 8).

4.2 VIRTUAL TRANSITION WITH MAXIMIZED MINIMUM VISITATION PROBABILITY

As elaborated in Section 4.1, our Algorithm 1 leverages a virtual transition weighted negentropy to regularize the loss estimator and induce the Bregman divergence used in OMD. It remains to specify an appropriate virtual transition p^x . The upside of employing such virtual transition p^x lies in that it implicitly helps to operate the update of OMD in the space of probability measures over infoset-action pairs instead of the sequence-form representations of policies. However, this also comes at the expense of enlarging the stability term of OMD. Specifically, upon applying the virtual transition to weight the negentropy, the stability term associated with OMD at each information set x_h will be enlarged by (approximately) a multiplicative factor of $1/p^x(x_h)$. This enlargement arises intuitively from the fact that, at each x_h , the Bregman divergence induced by ψ undergoes a downscaling, proportional to $p^x(x_h)$, thereby resulting in a relative increase in the stability term. Therefore, to ensure that the stability term is well-controlled, we design the following p^x which maximizes the minimum visitation probability of all x_h in its sequence-form representation:

$$p^x = \arg \max_{q \in \mathbb{P}^x} \min_{x_h \in \mathcal{X}_h, h \in [H]} q_{1:h}(x_h). \quad (6)$$

In the above display, we denote by \mathbb{P}^x the set of all the valid virtual transitions over infoset-action space. We note that such a virtual transition p^x can be efficiently computed by Algorithm 2 via backward dynamic programming.

Computation Due to the fact the update of OMD is now constrained onto a subset Π_{\max}^k of the entire space Π_{\max} of the sequence-form representation policies, the computation of Eq. (3) generally does not have a closed-form solution. We hereby provide an algorithm, which computes an approximate solution to Eq. (3), detailed in Algorithm 3. In particular, Algorithm 3 utilizes a Frank–Wolfe-type procedure to compute the update in Eq. (3). In particular, there will be T iterations in Algorithm 3, and in each iteration t , the policy will be updated towards the direction that minimizes the gradient of the objective function w.r.t. policy $\mu^{(t-1)}$ in iteration $t - 1$ by dynamic programming in Algorithm 4. We defer the details of Algorithm 3 and Algorithm 4 to Appendix F.

5 ANALYSIS

In this section, we first present the upper bound of the last-iterate convergence rate of our Algorithm 1. Then the lower bound for the problem of learning IIEFGs with bandit feedback and last-iterate convergence guarantee will be provided.

5.1 UPPER BOUND OF LAST-ITERATE CONVERGENCE

Theorem 5.1. *If Algorithm 1 is adopted by both players, for any $k \geq 1$, with probability at least $1 - \tilde{\mathcal{O}}(\delta)$, it holds that*

$$\text{NEG}_{\text{gap}}(\mu^k, \nu^k) = \mathcal{O} \left(\left[(XA + YB)^{\frac{1}{2}} k^{-\frac{1}{8}} + (XA + YB)^{\frac{1}{2}} H k^{-\frac{3}{8}} + (X^2A + Y^2B)^{\frac{1}{2}} k^{-\frac{1}{4}} + (X + Y)^{\frac{1}{4}} H k^{-\frac{1}{8}} \right] \cdot (X + Y) (\log(XAk/\delta) + \log(YBk/\delta)) \log^{\frac{1}{2}}(k) + k^{-\frac{1}{8}} H (\ln(XA) + \ln(YB)) + (XAB + YBH)/k \right).$$

Remark 5.2. *Ignoring the poly-logarithmic terms and when k is large enough (specifically, $k \geq \max\{H^4, (X^2A + Y^2B)^4 / (XA + YB)^4, (XA + YB)^{8/7} / (X + Y)^{10/7}\}$), we have $\text{NEG}_{\text{gap}}(\mu^k, \nu^k) = \tilde{\mathcal{O}}((X + Y)[(XA + YB)^{1/2} + (X + Y)^{1/4} H] k^{-1/8})$. Besides, when only obtaining an expected last-iterate convergence rate is desired, our Algorithm 1 has an improved last-iterate convergence rate of $\tilde{\mathcal{O}}((X + Y)[(X^2A + Y^2B)^{1/2} + (X + Y)^{1/4} H] k^{-1/6})$ in expectation, the details of which are deferred to Appendix C. Though the last-iterate convergence rate of our Algorithm 1 is inferior to the $\tilde{\mathcal{O}}(1/k)$ convergence rate by Lee et al. (2021); Liu et al. (2023), we note that both their algorithms can only work in the full-information setting. Further, we remark that the algorithm of Lee et al. (2021) needs the assumption that the NE of the IIEFG considered is unique, and the algorithm of Liu et al. (2023) requires both players being controlled by a central controller; and thus the algorithm*

of Liu et al. (2023) is not uncoupled. In contrast, our algorithm can work in the bandit feedback setting, is fully uncoupled between the two players, and can still guarantee a regret of order $\tilde{O}(k^{7/8})$ when the opponent of the max-player is an adversary. More importantly, we show in Section 5.2 that the lower bound of the convergence rate for learning IIEFGs with bandit feedback, last-iterate convergence guarantee, and uncoupled algorithms will be of order $\Omega(k^{-1/2})$ (for large enough k).

Proof Sketch of Theorem 5.1 We postpone the complete proof of Theorem 5.1 to Appendix B. Here we provide a proof sketch of it.

We denote by $\xi^{k,*} := (\mu^{k,*}, \nu^{k,*})$ the unique NE in the regularized game f_k in Eq. (5), where there is only a unique NE since f_k is strongly convex in μ and strongly concave in ν . We first show that in each episode k , the product policy $\xi^k := (\mu^k, \nu^k)$ generated by the algorithm will approach $\xi^{k,*}$ close enough by showing that the Bregman divergence $D_\psi(\xi^{k,*}, \xi^k)$ is an (approximate) contraction mapping. In particular, we show that

$$D_\psi(\xi^{k+1,*}, \xi^{k+1}) \lesssim (1 - \eta_k \varepsilon_k) D_\psi(\xi^{k,*}, \xi^k) + \eta_k^2 (X\bar{\tau}_k + Y\bar{r}_k) + \eta_k^2 (X^2A + Y^2B) + \eta_k \rho_k + \eta_k \sigma_k + \eta_k^2 \varepsilon_k^2 H^2 (XA + YB) + \omega_k, \quad (7)$$

where we denote

$$\begin{aligned} \bar{\tau}_k &= \frac{1}{X} \sum_{h, x_h, a_h} \frac{1}{p_{1:h}^x(x_h)} \left(\frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} - 1 \right), \\ \bar{r}_k &= \frac{1}{Y} \sum_{h, y_h, b_h} \frac{1}{p_{1:h}^y(y_h)} \left(\frac{\mathbb{I}_h^k \{y_h, b_h\}}{\nu_{1:h}^k(y_h, b_h) + \gamma_k} - 1 \right), \\ \rho_k &= \sum_{h, x_h, a_h} \mu_{1:h}^k(x_h, a_h) \left[(\mathbf{G}\nu^k)[(x_h, a_h)] - \left(\frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} (1 - r_h^k) \right) \right] \\ &\quad + \sum_{h, y_h, b_h} \nu_{1:h}^k(y_h, b_h) \left[(1 - (\mathbf{G}^\top \mu^k)[(y_h, b_h)]) - \frac{\mathbb{I}_h^k \{y_h, b_h\}}{\nu_{1:h}^k(y_h, b_h) + \gamma_k} r_h^k \right], \\ \sigma_k &= \sum_{h, x_h, a_h} \mu_{1:k}^{k,*}(x_h, a_h) \left[(\mathbf{G}\nu^k)[(x_h, a_h)] - \left(\frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma} (1 - r_h^k) \right) \right] \\ &\quad + \sum_{h, y_h, b_h} \nu_{1:k}^{k,*}(y_h, b_h) \left[\frac{\mathbb{I}_h^k \{y_h, b_h\}}{\nu_{1:h}^k(y_h, b_h) + \gamma_k} r_h^k - (1 - (\mathbf{G}^\top \mu^k)[(y_h, b_h)]) \right], \\ \omega_k &= D_\psi(\mu^{k+1,*}, \mu^{k+1}) - D_\psi(\mu^{k,*}, \mu^{k+1}) + D_\psi(\nu^{k+1,*}, \nu^{k+1}) - D_\psi(\nu^{k,*}, \nu^{k+1}). \end{aligned}$$

Expanding the above recursion, we can bound $D_\psi(\xi^{k+1,*}, \xi^{k+1})$ as

$$\begin{aligned} D_\psi(\xi^{k+1,*}, \xi^{k+1}) &\lesssim \underbrace{\sum_{i=1}^k w_k^i \eta_i \rho_i}_{\text{Term 1}} + \underbrace{\sum_{i=1}^k w_k^i \eta_i \sigma_i}_{\text{Term 2}} + \underbrace{(XA + YB) H^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2}_{\text{Term 3}} \\ &\quad + \underbrace{\sum_{i=1}^k w_k^i \eta_i^2 (X\bar{\tau}_i + Y\bar{r}_i)}_{\text{Term 4}} + \underbrace{\sum_{i=1}^k w_k^i \eta_i^2 (X^2A + Y^2B)}_{\text{Term 5}} + \underbrace{\sum_{i=1}^k w_k^i \omega_i}_{\text{Term 6}}, \quad (8) \end{aligned}$$

where $w_k^i = \prod_{j=i+1}^k (1 - \eta_j \varepsilon_j)$ is the contraction parameter. Then we bound each of the above terms in by Lemma B.4 - Lemma B.9 in Appendix B.2. Note that we follow a similar analysis scheme of Cai et al. (2023) to bound the last-iterate convergence of learning matrix games with bandit feedback. However, we also remark that the straightforward application of their analysis will not address our problem of learning IIEFGs with bandit feedback, since we leverage a different regularizer and a new virtual transition p^x computed by Algorithm 2, which serves as a core ingredient of the analysis in deriving the contraction of Eq. (7) and bounding **Term 6**. Besides, compared with the analysis of Cai et al. (2023), the additional **Term 5** in Eq. (8) comes from the fact that we

486 establish a refined analysis in the case of IIEFGs to further sharpen the dependence on X and A (as
487 well as Y and B) of the final convergence rate.

488 Further, one can see that the NE policy profile $\xi^{k,*}$ of the perturbed game in Eq. (5) is also an
489 approximate NE of the original game in Eq. (4), enabling to bound $\text{NEGap}(\xi^k)$ using $\text{NEGap}(\xi^{k,*})$
490 together with the distance between ξ^k and $\xi^{k,*}$ weighted by the virtual transitions as bellow:
491

$$492 \quad \text{NEGap}(\xi^k) \leq \text{NEGap}(\xi^{k,*}) + X \|p^x(\mu^k - \mu^{k,*})\|_1 + Y \|p^y(\nu^k - \nu^{k,*})\|_1, \quad (9)$$

493 where $\text{NEGap}(\xi^{k,*})$ can be controlled by Lemma B.2. Due to the constructed virtual transition
494 p^x and p^y , the second and the third term in Eq. (9) are actually the ℓ_1 -norm of the difference
495 between the probability measures over infoset-action spaces, which thus turns out to be bounded by
496 $\mathcal{O}(\sqrt{\text{KL}(p^x \mu^{k,*}, p^x \mu^k)})$ and $\mathcal{O}(\sqrt{\text{KL}(p^y \nu^{k,*}, p^y \nu^k)})$ by Pinsker's inequality. Also, thanks to the
497 virtual transition weighted negentropy ψ , one can see that $\text{KL}(p^x \mu^{k,*}, p^x \mu^k) = D_\psi(\mu^{k,*}, \mu^k)$ (and
498 similarly on the min-player side). Therefore, the proof can be concluded by substituting Eq. (8) into
499 Eq. (9) and then using Lemma B.2 and Lemma B.4 - Lemma B.9.
500

501 5.2 LOWER BOUND OF LAST-ITERATE CONVERGENCE

502 **Theorem 5.3.** *For any algorithm Alg that both players adopt to generate policy profile (μ^k, ν^k)
503 and is uncoupled between both players, there exists an IIEFG instance such that the lower bound
504 of the last-iterate convergence of learning this IIEFG in the bandit-feedback setting satisfies
505 $\text{NEGap}(\mu^k, \nu^k) = \Omega(\sqrt{XA + YB}k^{-1/2})$, when $k \geq \max(XA, YB)$.*

506 *Proof Sketch.* The idea of the proof is to leverage the fact that if an uncoupled algorithm can learn
507 the NE of IIEFGs with a last-iterate convergence guarantee of $\tilde{\Theta}(k^{-\alpha})$ ($\alpha \in [0, 1]$) in the bandit
508 feedback setting, then it can be used to learn IIEFGs where the opponent is an adversary with a regret
509 of order $\tilde{\Theta}(k^{1-\alpha})$. Therefore, considering that the hardness of minimizing regret of IIEFGs with an
510 adversarial opponent is equivalent to minimizing regret on a bandit problem with AX arms (Bai
511 et al., 2022; Fiegel et al., 2023), the proof of Theorem 5.3 can be completed by contradiction. \square
512

513 **Remark 5.4.** *Compared with the lower bound of the convergence rate above, the upper bound in
514 Theorem 5.1 is loose by a factor of $\tilde{\mathcal{O}}((X + Y)k^{3/8})$ (for large enough X, Y, A and B). We believe
515 one of the promising approaches to improve the upper bound of the convergence rate might be to
516 consider using the optimistic OMD/FTRL, which utilizes accelerated techniques from the optimiza-
517 tion perspective and is typically used to achieve the $\tilde{\mathcal{O}}(1/k)$ convergence rate for learning IIEFGs
518 with last-iterate convergence in the full-information setting. One of the main difficulties of using
519 optimistic OMD/FTRL in conjunction with the regularization technique to achieve a faster last-
520 iterate convergence rate of learning IIEFGs in the bandit feedback setting is that the loss estimator
521 constructed in the bandit feedback setting (either unbiased or optimistically biased) to serve as a
522 surrogate of the true loss would have undesirably large variance, making the stability of optimistic
523 OMD/FTRL hard to be controlled even in the special case of learning matrix games. We leave the
524 possible improvement of our convergence upper bound as our future study.*

525 6 CONCLUSION

526 In this work, we make the first step to establishing the algorithm that learns an approximate NE of
527 IIEFGs in the bandit feedback setting with finite-time last-iterate convergence. Our algorithm is fully
528 uncoupled between the two players involved in the games and does not require any coordination,
529 communication, or shared randomness between these players. We prove that our algorithm achieves
530 the last-iterate convergence of order $\tilde{\mathcal{O}}((X + Y)[(XA + YB)^{1/2} + (X + Y)^{1/4}H]k^{-1/8})$ with high
531 probability and of order $\tilde{\mathcal{O}}((X + Y)[(X^2A + Y^2B)^{1/2} + (X + Y)^{1/4}H]k^{-1/6})$ in expectation (for large
532 enough k). Also, we provide the lower bound of order $\Omega(\sqrt{XA + YB}k^{-1/2})$ for learning IIEFGs
533 with last-iterate convergence guarantee in the bandit feedback setting. An interesting problem might
534 be closing the gap between the established convergence upper and lower bound, which still remains
535 open in the special case of learning matrix games with the last-iterate convergence guarantee in the
536 bandit feedback setting. We will leave the investigation of this for our future research endeavors.
537
538
539

REFERENCES

- 540
541
542 Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, Kentaro Toyoshima, and Atsushi Iwasaki. Last-iterate
543 convergence with full and noisy feedback in two-player zero-sum games. In *International Con-*
544 *ference on Artificial Intelligence and Statistics, 25-27 April 2023, Palau de Congressos, Valencia,*
545 *Spain*, volume 206 of *Proceedings of Machine Learning Research*, pp. 7999–8028. PMLR, 2023.
- 546 Yu Bai and Chi Jin. Provable self-play algorithms for competitive reinforcement learning. In *Pro-*
547 *ceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July*
548 *2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pp. 551–560.
549 PMLR, 2020.
- 550 Yu Bai, Chi Jin, and Tiancheng Yu. Near-optimal reinforcement learning with self-play. *Advances*
551 *in neural information processing systems*, 33:2159–2170, 2020.
- 552
553 Yu Bai, Chi Jin, Song Mei, and Tiancheng Yu. Near-optimal learning of extensive-form games with
554 imperfect information. In *International Conference on Machine Learning, ICML 2022, 17-23 July*
555 *2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*, pp.
556 1337–1382. PMLR, 2022.
- 557 Martino Bernasconi, Alberto Marchesi, and Francesco Trovò. Learning extensive-form perfect
558 equilibria in two-player zero-sum sequential games. In *International Conference on Artificial*
559 *Intelligence and Statistics, 2-4 May 2024, Palau de Congressos, Valencia, Spain*, volume 238 of
560 *Proceedings of Machine Learning Research*, pp. 2152–2160. PMLR, 2024.
- 561 Mario Bravo, David Leslie, and Panayotis Mertikopoulos. Bandit learning in concave n-person
562 games. *Advances in Neural Information Processing Systems*, 31, 2018.
- 563
564 Noam Brown and Tuomas Sandholm. Superhuman ai for heads-up no-limit poker: Libratus beats
565 top professionals. *Science*, 359(6374):418–424, 2018.
- 566 Neil Burch, Matej Moravcik, and Martin Schmid. Revisiting CFR+ and alternating updates. *J. Artif.*
567 *Intell. Res.*, 64:429–443, 2019.
- 568
569 Yang Cai and Weiqiang Zheng. Doubly optimal no-regret learning in monotone games. In *Interna-*
570 *tional Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*,
571 volume 202 of *Proceedings of Machine Learning Research*, pp. 3507–3524. PMLR, 2023.
- 572 Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Finite-time last-iterate convergence for learn-
573 ing in multi-player games. In *Advances in Neural Information Processing Systems 35: Annual*
574 *Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA,*
575 *USA, November 28 - December 9, 2022, 2022*.
- 576
577 Yang Cai, Haipeng Luo, Chen-Yu Wei, and Weiqiang Zheng. Uncoupled and convergent learning
578 in two-player zero-sum markov games with bandit feedback. In *Advances in Neural Information*
579 *Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023,*
580 *NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023, 2023*.
- 581 Shicong Cen, Yuting Wei, and Yuejie Chi. Fast policy extragradient methods for competitive games
582 with entropy regularization. In *Advances in Neural Information Processing Systems 34: Annual*
583 *Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14,*
584 *2021, virtual*, pp. 27952–27964, 2021.
- 585 Shicong Cen, Yuejie Chi, Simon Shaolei Du, and Lin Xiao. Faster last-iterate convergence of policy
586 optimization in zero-sum markov games. In *The Eleventh International Conference on Learning*
587 *Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- 588
589 Ziyi Chen, Shaocong Ma, and Yi Zhou. Sample efficient stochastic policy extragradient algorithm
590 for zero-sum markov game. In *International Conference on Learning Representation, 2022*.
- 591 Qiwen Cui and Lin F. Yang. Minimax sample complexity for turn-based stochastic game. In *Pro-*
592 *ceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence, UAI 2021,*
593 *Virtual Event, 27-30 July 2021*, volume 161 of *Proceedings of Machine Learning Research*, pp.
1496–1504. AUAI Press, 2021.

- 594 Qiwen Cui, Kaiqing Zhang, and Simon S. Du. Breaking the curse of multiagents in a large state
595 space: RL in markov games with independent linear function approximation. In *The Thirty Sixth*
596 *Annual Conference on Learning Theory, COLT 2023, 12-15 July 2023, Bangalore, India*, volume
597 195 of *Proceedings of Machine Learning Research*, pp. 2651–2652. PMLR, 2023.
- 598
- 599 Dmitriy Drusvyatskiy, Maryam Fazel, and Lillian J Ratliff. Improved rates for derivative free gra-
600 dient play in strongly monotone games. In *2022 IEEE 61st Conference on Decision and Control*
601 *(CDC)*, pp. 3403–3408. IEEE, 2022.
- 602
- 603 Jianqing Fan, Zhaoran Wang, Yuchen Xie, and Zhuoran Yang. A theoretical analysis of deep q-
604 learning. In *Proceedings of the 2nd Annual Conference on Learning for Dynamics and Control,*
605 *LADC 2020, Online Event, Berkeley, CA, USA, 11-12 June 2020*, volume 120 of *Proceedings of*
606 *Machine Learning Research*, pp. 486–489. PMLR, 2020.
- 607
- 608 Gabriele Farina and Tuomas Sandholm. Model-free online learning in unknown sequential decision
609 making problems and games. In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI*
610 *2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021,*
611 *The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual*
612 *Event, February 2-9, 2021*, pp. 5381–5390. AAAI Press, 2021.
- 613
- 614 Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Stochastic regret minimization in
615 extensive-form games. In *Proceedings of the 37th International Conference on Machine Learn-*
616 *ing, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning*
617 *Research*, pp. 3018–3028. PMLR, 2020.
- 618
- 619 Gabriele Farina, Robin Schmucker, and Tuomas Sandholm. Bandit linear optimization for sequen-
620 tial decision making and extensive-form games. In *Thirty-Fifth AAAI Conference on Artificial*
621 *Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intel-*
622 *ligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence,*
623 *EAAI 2021, Virtual Event, February 2-9, 2021*, pp. 5372–5380. AAAI Press, 2021.
- 624
- 625 Yi Feng, Hu Fu, Qun Hu, Ping Li, Ioannis Panageas, Bo Peng, and Xiao Wang. On the last-
626 iterate convergence in time-varying zero-sum games: Extra gradient succeeds where optimism
627 fails. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural*
628 *Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,*
629 *2023*, 2023.
- 630
- 631 Côme Fiegel, Pierre Ménard, Tadashi Kozuno, Rémi Munos, Vianney Perchet, and Michal Valko.
632 Adapting to game trees in zero-sum imperfect information games. In *International Conference*
633 *on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of
634 *Proceedings of Machine Learning Research*, pp. 10093–10135. PMLR, 2023.
- 635
- 636 Haobo Fu, Weiming Liu, Shuang Wu, Yijia Wang, Tao Yang, Kai Li, Junliang Xing, Bin Li, Bo Ma,
637 Qiang Fu, and Wei Yang. Actor-critic policy optimization in a large-scale imperfect-information
638 game. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual*
639 *Event, April 25-29, 2022*. OpenReview.net, 2022.
- 640
- 641 Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. Tight last-iterate convergence rates
642 for no-regret learning in multi-player games. In *Advances in Neural Information Processing*
643 *Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020,*
644 *December 6-12, 2020, virtual*, 2020.
- 645
- 646 Eduard Gorbunov, Adrien B. Taylor, and Gauthier Gidel. Last-iterate convergence of optimistic gra-
647 dient method for monotone variational inequalities. In *Advances in Neural Information Process-*
ing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS
2022, New Orleans, LA, USA, November 28 - December 9, 2022, 2022.
- 648
- 649 Amy Greenwald and Keith Hall. Correlated q-learning. In *Machine Learning, Proceedings of the*
650 *Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*,
651 pp. 242–249. AAAI Press, 2003.

- 648 Thomas Dueholm Hansen, Peter Bro Miltersen, and Uri Zwick. Strategy iteration is strongly poly-
649 nomial for 2-player turn-based stochastic games with a constant discount factor. *J. ACM*, 60(1):
650 1:1–1:16, 2013.
- 651 Johannes Heinrich, Marc Lanctot, and David Silver. Fictitious self-play in extensive-form games.
652 In *Proceedings of the 32nd International Conference on Machine Learning, ICML 2015, Lille,*
653 *France, 6-11 July 2015*, volume 37 of *JMLR Workshop and Conference Proceedings*, pp. 805–
654 813. JMLR.org, 2015.
- 655 Samid Hoda, Andrew Gilpin, Javier Peña, and Tuomas Sandholm. Smoothing techniques for com-
656 puting nash equilibria of sequential games. *Math. Oper. Res.*, 35(2):494–512, 2010.
- 657 Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. On the convergence
658 of single-call stochastic extra-gradient methods. In *Advances in Neural Information Processing*
659 *Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019,*
660 *December 8-14, 2019, Vancouver, BC, Canada*, pp. 6936–6946, 2019.
- 661 Junling Hu and Michael P. Wellman. Nash q-learning for general-sum stochastic games. *J. Mach.*
662 *Learn. Res.*, 4:1039–1069, 2003.
- 663 Yuanhanqing Huang and Jianghai Hu. Zeroth-order learning in continuous games via residual pseu-
664 do-gradient estimates. *arXiv preprint arXiv:2301.02279*, 2023.
- 665 Zeyu Jia, Lin F. Yang, and Mengdi Wang. Feature-based q-learning for two-player stochastic games.
666 *CoRR*, abs/1906.00423, 2019.
- 667 Michael Johanson, Nolan Bard, Marc Lanctot, Richard G. Gibson, and Michael Bowling. Efficient
668 nash equilibrium approximation through monte carlo counterfactual regret minimization. In *In-*
669 *ternational Conference on Autonomous Agents and Multiagent Systems, AAMAS 2012, Valencia,*
670 *Spain, June 4-8, 2012 (3 Volumes)*, pp. 837–846. IFAAMAS, 2012.
- 671 Michael Jordan, Tianyi Lin, and Zhengyuan Zhou. Adaptive, doubly optimal no-regret learning in
672 strongly monotone and exp-concave games with gradient feedback. *Operations Research*, 2024.
- 673 Daphne Koller and Nimrod Megiddo. The complexity of two-person zero-sum games in extensive
674 form. *Games and economic behavior*, 4(4):528–552, 1992.
- 675 Daphne Koller, Nimrod Megiddo, and Bernhard Von Stengel. Efficient computation of equilibria
676 for extensive two-person games. *Games and economic behavior*, 14(2):247–259, 1996.
- 677 Tadashi Kozuno, Pierre Ménard, Rémi Munos, and Michal Valko. Learning in two-player zero-
678 sum partially observable markov games with perfect recall. In *Advances in Neural Information*
679 *Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021,*
680 *NeurIPS 2021, December 6-14, 2021, virtual*, pp. 11987–11998, 2021.
- 681 Christian Kroer, Kevin Waugh, Fatma Kiliç-Karzan, and Tuomas Sandholm. Faster first-order
682 methods for extensive-form game solving. In *Proceedings of the Sixteenth ACM Conference on*
683 *Economics and Computation, EC ’15, Portland, OR, USA, June 15-19, 2015*, pp. 817–834. ACM,
684 2015.
- 685 Christian Kroer, Gabriele Farina, and Tuomas Sandholm. Solving large sequential games with
686 the excessive gap technique. In *Advances in Neural Information Processing Systems 31: Annual*
687 *Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018,*
688 *Montréal, Canada*, pp. 872–882, 2018.
- 689 Harold W Kuhn. Extensive games. *Proceedings of the National Academy of Sciences*, 36(10):
690 570–576, 1950.
- 691 HW Kuhn. Extensive games and the problem of information. *Contributions to the Theory of Games*,
692 (24):193, 1953.
- 693 Moyuru Kurita and Kunihito Hoki. Method for constructing artificial intelligence player with ab-
694 stractions to markov decision processes in multiplayer game of mahjong. *IEEE Trans. Games*, 13
695 (1):99–110, 2021.

- 702 Michail G. Lagoudakis and Ronald Parr. Value function approximation in zero-sum markov games.
703 In *UAI '02, Proceedings of the 18th Conference in Uncertainty in Artificial Intelligence, University of Alberta, Edmonton, Alberta, Canada, August 1-4, 2002*, pp. 283–292. Morgan Kaufmann,
704 2002.
705
- 706 Marc Lanctot, Kevin Waugh, Martin Zinkevich, and Michael H. Bowling. Monte carlo sampling for
707 regret minimization in extensive games. In *Advances in Neural Information Processing Systems*
708 *22: 23rd Annual Conference on Neural Information Processing Systems 2009. Proceedings of a*
709 *meeting held 7-10 December 2009, Vancouver, British Columbia, Canada*, pp. 1078–1086. Curran
710 Associates, Inc., 2009.
711
- 712 Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay,
713 Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, et al.
714 Openspiel: A framework for reinforcement learning in games. *arXiv preprint arXiv:1908.09453*,
715 2019.
716
- 717 Chung-Wei Lee, Christian Kroer, and Haipeng Luo. Last-iterate convergence in extensive-form
718 games. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural*
719 *Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pp. 14293–
720 14305, 2021.
- 721 Gen Li, Yuejie Chi, Yuting Wei, and Yuxin Chen. Minimax-optimal multi-agent RL in markov
722 games with a generative model. In *NeurIPS*, 2022.
723
- 724 Junjie Li, Sotetsu Koyamada, Qiwei Ye, Guoqing Liu, Chao Wang, Ruihan Yang, Li Zhao, Tao
725 Qin, Tie-Yan Liu, and Hsiao-Wuen Hon. Suphx: Mastering mahjong with deep reinforcement
726 learning. *CoRR*, abs/2003.13590, 2020.
- 727 Tianyi Lin, Zhengyuan Zhou, Panayotis Mertikopoulos, and Michael I. Jordan. Finite-time last-
728 iterate convergence for multi-agent learning in games. In *Proceedings of the 37th International*
729 *Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of
730 *Proceedings of Machine Learning Research*, pp. 6161–6171. PMLR, 2020.
731
- 732 Tianyi Lin, Zhengyuan Zhou, Wenjia Ba, and Jiawei Zhang. Doubly optimal no-regret online learn-
733 ing in strongly monotone games with bandit feedback. *arXiv preprint arXiv:2112.02856*, 2021.
- 734 Michael L. Littman. Friend-or-foe q-learning in general-sum games. In *Proceedings of the*
735 *Eighteenth International Conference on Machine Learning (ICML 2001), Williams College,*
736 *Williamstown, MA, USA, June 28 - July 1, 2001*, pp. 322–328. Morgan Kaufmann, 2001.
737
- 738 Michael L. Littman and Csaba Szepesvári. A generalized reinforcement-learning model: Con-
739 vergence and applications. In *Machine Learning, Proceedings of the Thirteenth International*
740 *Conference (ICML '96), Bari, Italy, July 3-6, 1996*, pp. 310–318. Morgan Kaufmann, 1996.
- 741 Mingyang Liu, Asuman E. Ozdaglar, Tiancheng Yu, and Kaiqing Zhang. The power of regular-
742 ization in solving extensive-form games. In *The Eleventh International Conference on Learning*
743 *Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
744
- 745 Qinghua Liu, Tiancheng Yu, Yu Bai, and Chi Jin. A sharp analysis of model-based reinforcement
746 learning with self-play. In *Proceedings of the 38th International Conference on Machine Learn-*
747 *ing, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning*
748 *Research*, pp. 7001–7010. PMLR, 2021.
- 749 Weiming Liu, Huacong Jiang, Bin Li, and Houqiang Li. Equivalence analysis between counter-
750 factual regret minimization and online mirror descent. In *International Conference on Machine*
751 *Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings*
752 *of Machine Learning Research*, pp. 13717–13745. PMLR, 2022.
753
- 754 Aryan Mokhtari, Asuman E. Ozdaglar, and Sarath Pattathil. Convergence rate of $o(1/k)$ for op-
755 timistic gradient and extragradient methods in smooth convex-concave saddle point problems.
SIAM J. Optim., 30(4):3230–3251, 2020.

- 756 Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor
757 Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial
758 intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.
- 759 Rémi Munos, Julien Pérolat, Jean-Baptiste Lespiau, Mark Rowland, Bart De Vylder, Marc Lanctot,
760 Finbarr Timbers, Daniel Hennes, Shayegan Omidshafiei, Audrunas Gruslys, Mohammad Ghesh-
761 laghi Azar, Edward Lockhart, and Karl Tuyls. Fast computation of nash equilibria in imperfect
762 information games. In *Proceedings of the 37th International Conference on Machine Learning,*
763 *ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning*
764 *Research*, pp. 7119–7129. PMLR, 2020.
- 765 John F Nash Jr. Equilibrium points in n-person games. *Proceedings of the national academy of*
766 *sciences*, 36(1):48–49, 1950.
- 767 Gergely Neu. Explore no more: Improved high-probability regret bounds for non-stochastic ban-
768 dits. In *Advances in Neural Information Processing Systems 28: Annual Conference on Neural*
769 *Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada*, pp.
770 3168–3176, 2015.
- 771 Julien Pérolat, Bruno Scherrer, Bilal Piot, and Olivier Pietquin. Approximate dynamic programming
772 for two-player zero-sum markov games. In *Proceedings of the 32nd International Conference on*
773 *Machine Learning, ICML 2015, Lille, France, 6-11 July 2015*, volume 37 of *JMLR Workshop and*
774 *Conference Proceedings*, pp. 1321–1329. JMLR.org, 2015.
- 775 Martin Schmid, Neil Burch, Marc Lanctot, Matej Moravcik, Rudolf Kadlec, and Michael Bowling.
776 Variance reduction in monte carlo counterfactual regret minimization (VR-MCCFR) for extensive
777 form games using baselines. In *The Thirty-Third AAAI Conference on Artificial Intelligence,*
778 *AAAI 2019, The Thirty-First Innovative Applications of Artificial Intelligence Conference, IAAI*
779 *2019, The Ninth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2019,*
780 *Honolulu, Hawaii, USA, January 27 - February 1, 2019*, pp. 2157–2164. AAAI Press, 2019.
- 781 Martin Schmid, Matej Moravcik, Neil Burch, Rudolf Kadlec, Joshua Davidson, Kevin Waugh,
782 Nolan Bard, Finbarr Timbers, Marc Lanctot, Zach Holland, Elnaz Davoodi, Alden Christianson,
783 and Michael Bowling. Player of games. *CoRR*, abs/2112.03178, 2021.
- 784 Aaron Sidford, Mengdi Wang, Xian Wu, Lin Yang, and Yinyu Ye. Near-optimal time and sam-
785 ple complexities for solving markov decision processes with a generative model. In *Advances*
786 *in Neural Information Processing Systems 31: Annual Conference on Neural Information Pro-*
787 *cessing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pp. 5192–5202,
788 2018.
- 789 Ziang Song, Song Mei, and Yu Bai. When can we learn general-sum markov games with a large
790 number of players sample-efficiently? In *The Tenth International Conference on Learning Rep-*
791 *resentations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- 792 Finnegan Southey, Michael P Bowling, Bryce Larson, Carmelo Piccione, Neil Burch, Darse Billings,
793 and Chris Rayner. Bayes’ bluff: Opponent modelling in poker. *arXiv preprint arXiv:1207.1411*,
794 2012.
- 795 Oskari Tammelin. Solving large imperfect information games using CFR+. *CoRR*, abs/1407.5042,
796 2014.
- 797 Yuandong Tian, Qucheng Gong, and Yu Jiang. Joint policy search for multi-agent collaboration
800 with imperfect information. In *Advances in Neural Information Processing Systems 33: Annual*
801 *Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12,*
802 *2020, virtual*, 2020.
- 803 Bernhard Von Stengel. Efficient computation of behavior strategies. *Games and Economic Behavior*,
804 14(2):220–246, 1996.
- 805 Yuanhao Wang, Qinghua Liu, Yu Bai, and Chi Jin. Breaking the curse of multiagency: Provably
806 efficient decentralized multi-agent RL with function approximation. In *The Thirty Sixth Annual*
807 *Conference on Learning Theory, COLT 2023, 12-15 July 2023, Bangalore, India*, volume 195 of
808 *Proceedings of Machine Learning Research*, pp. 2793–2848. PMLR, 2023.
- 809

810 Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Last-iterate convergence of
811 decentralized optimistic gradient descent/ascent in infinite-horizon competitive markov games.
812 In *Conference on Learning Theory, COLT 2021, 15-19 August 2021, Boulder, Colorado, USA*,
813 volume 134 of *Proceedings of Machine Learning Research*, pp. 4259–4299. PMLR, 2021a.

814 Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Last-iterate convergence of
815 decentralized optimistic gradient descent/ascent in infinite-horizon competitive markov games.
816 In *Conference on learning theory*, pp. 4259–4299. PMLR, 2021b.

817

818 Wei Xiong, Han Zhong, Chengshuai Shi, Cong Shen, and Tong Zhang. A self-play posterior sam-
819 pling algorithm for zero-sum markov games. In *International Conference on Machine Learning,*
820 *ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine*
821 *Learning Research*, pp. 24496–24523. PMLR, 2022.

822 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Model-free reinforcement learning: from clipped
823 pseudo-regret to sample complexity. In *Proceedings of the 38th International Conference on*
824 *Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of*
825 *Machine Learning Research*, pp. 12653–12662. PMLR, 2021.

826

827 Dongruo Zhou, Quanquan Gu, and Csaba Szepesvári. Nearly minimax optimal reinforcement learn-
828 ing for linear mixture markov decision processes. In *Conference on Learning Theory, COLT 2021,*
829 *15-19 August 2021, Boulder, Colorado, USA*, volume 134 of *Proceedings of Machine Learning*
830 *Research*, pp. 4532–4576. PMLR, 2021.

831 Martin Zinkevich, Michael Johanson, Michael H. Bowling, and Carmelo Piccione. Regret mini-
832 mization in games with incomplete information. In *Advances in Neural Information Processing*
833 *Systems 20, Proceedings of the Twenty-First Annual Conference on Neural Information Process-*
834 *ing Systems, Vancouver, British Columbia, Canada, December 3-6, 2007*, pp. 1729–1736. Curran
835 Associates, Inc., 2007.

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

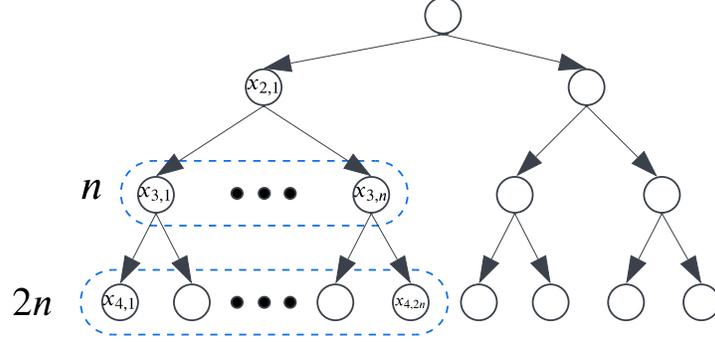
861

862

863

A MORE DISCUSSIONS ON VIRTUAL TRANSITION PROBABILITIES

A.1 ILLUSTRATION ON THE FAILURE OF USING UNIFORM VIRTUAL TRANSITION



$$p_{1:H}(x_{4,1}) = \dots = p_{1:H}(x_{4,2n}) = \frac{1}{4n}$$

Figure 1: An illustrative example where using uniform virtual transition p fails to guarantee $\min_{x_h \in \mathcal{X}_h, h \in [H]} p_{1:h}(x_h) \geq 1/X$.

On the IIEFG instance shown in Figure 1, there is only one action a and $H = 4$. Each infoset x in the game tree of this instance satisfies $|C(x, a)| = 2$ except for infoset $x_{2,1}$, which is such that $|C(x_{2,1}, a)| = n$ with some $n \geq 2$. Now suppose the uniform distribution p is used as a virtual transition over infoset-action spaces. Then for all the descendants $\{x_{4,i}\}_{i=1}^{2n}$ on step $h = 4$ of infoset $x_{2,1}$, one can see that $p_{1:H}(x_{4,i}) = \frac{1}{2} \cdot \frac{1}{n} \cdot \frac{1}{2} = \frac{1}{4n}$, while there are only $X = 9 + 3n$ infosets in total. Thus, it will happen that $p_{1:H}(x_{4,i}) < \frac{1}{X}$ when $n > 9$.

Actually, one can easily construct an IIEFG instance such that $\min_{x_H \in \mathcal{X}_H} p_{1:H}(x_H) \leq \mathcal{O}(\frac{1}{n^m})$ and $X = \mathcal{O}(mn + c)$ with c as a parameter that depends on m but not n for uniform virtual transition p . Therefore, when using uniform distribution p as a virtual transition, $\max_{x_H \in \mathcal{X}_H} 1/p_{1:H}(x_H)$ might be prohibitively large and lead to a convergence rate with much worse dependence on X than the virtual transition constructed in our Algorithm 2.

A.2 BALANCED EFFECTS OF THE PROPOSED VIRTUAL TRANSITION PROBABILITY

Lemma A.1. For any $h \in [H]$ and $x_h \in \mathcal{X}_h$, the constructed virtual transition p^x guarantees that $1/p_{1:h}^x(x_h) \leq X$.

Proof. Clearly, $p_{1:h}^x(\cdot)$ is minimized at $h = H$ for some $x_H \in \mathcal{X}_H$ by the definition of virtual transition. By the construction of $p_{1:h}^x(\cdot)$ in Algorithm 2, one can deduce that $\forall x_H \in \mathcal{X}_H$, it holds that (understanding $\{(x_h, a_h)\}_{h \in [H-1]}$ as the unique trajectory leading to x_H below)

$$\begin{aligned} p_{1:H}^x(x_H) &= q[x_H] \\ &= q[x_{H-1}] \cdot \frac{c[x_H]}{\sum_{x'_H \in C(x_{H-1}, a_{H-1})} c[x'_H]} \\ &= q[x_{H-2}] \cdot \frac{c[x_{H-1}]}{\sum_{x'_{H-1} \in C(x_{H-2}, a_{H-2})} c[x'_{H-1}]} \cdot \frac{c[x_H]}{\sum_{x'_H \in C(x_{H-1}, a_{H-1})} c[x'_H]} \\ &= q[x_{H-2}] \cdot \frac{c[x_{H-1}]}{\sum_{x'_{H-1} \in C(x_{H-2}, a_{H-2})} c[x'_{H-1}]} \cdot \frac{c[x_H]}{d[x_{H-1}, a_{H-1}]} \\ &\stackrel{(i)}{\geq} q[x_{H-2}] \cdot \frac{c[x_{H-1}]}{\sum_{x'_{H-1} \in C(x_{H-2}, a_{H-2})} c[x'_{H-1}]} \cdot \frac{c[x_H]}{c[x_{H-1}]} \end{aligned}$$

$$\begin{aligned}
&= q[x_{H-2}] \cdot \frac{c[x_H]}{\sum_{x'_{H-1} \in C(x_{H-2}, a_{H-2})} c[x'_{H-1}]} \\
&\geq \dots \\
&\geq \frac{c[x_H]}{\sum_{x_1 \in \mathcal{X}_1} c[x_1]} \\
&\geq \frac{c[x_H]}{X_H} \\
&\geq \frac{c[x_H]}{X} \\
&= \frac{1}{X},
\end{aligned}$$

where $c[\cdot]$, $q[\cdot]$, and $d[\cdot, \cdot]$ are defined in our Algorithm 2; and (i) is due to $c[x_{H-1}] = \max_{a \in \mathcal{A}} d[x_{H-1}, a] \geq d[x_{H-1}, a_{H-1}]$. \square

The property shown in this lemma of our constructed virtual transition p^x serves as a key ingredient in the analysis (say, when bounding our **Term 4** and when establishing the final convergence upper bound of the NE gap in the proof of Theorem 5.1) as we shall see.

B PROOF OF HIGH-PROBABILITY LAST-ITERATE CONVERGENCE RATE

Lemma B.1 (One-step analysis of OMD with virtual transition weighted negentropy regularized loss). *Let*

$$\begin{aligned}
\mu' = \operatorname{argmin}_{\tilde{\mu} \in \Omega} & \sum_{h, x_h, a_h} \tilde{\mu}_{1:h}(x_h, a_h) (\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))) \\
& + \frac{1}{\eta} D_\psi(\tilde{\mu}, \mu),
\end{aligned}$$

for some convex set $\Omega \subseteq \Pi_{\max}$, $\ell \in \mathbb{R}_{\geq 0}^{XA}$, and $\varepsilon \in \left[0, \frac{1}{\eta}\right]^{XA}$. Then $\forall u \in \Omega$.

$$\begin{aligned}
&\langle \mu' - \mu, \ell + \varepsilon p \log p \mu \rangle \\
&\leq \sum_{h, (x_h, a_h)} \left[\frac{\eta}{p_{1:h}^x(x_h)} \mu_{1:h}(x_h, a_h) \ell^2(x_h, a_h) + \eta \varepsilon^2(x_h, a_h) \log^2(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)) \right],
\end{aligned}$$

where $(\varepsilon p \log p \mu)[(x_h, a_h)] := \varepsilon p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))$.

Proof. The common one-step analysis of OMD shows that

$$\langle \mu' - \mu, \ell + \varepsilon p \log p \mu \rangle \leq \frac{1}{\eta} (D_\psi(u, \mu) - D_\psi(u, \mu') - D_\psi(\mu', \mu)).$$

Then, to upper bound $\langle \mu - \mu', \ell + \varepsilon p \log p \mu \rangle - \frac{1}{\eta} D_\psi(\mu', \mu)$, notice that

$$\begin{aligned}
&\langle \mu - \mu', \ell + \varepsilon p \log p \mu \rangle - \frac{1}{\eta} D_\psi(\mu', \mu) \\
&\leq \sup_{v \in \mathbb{R}_{\geq 0}^{XA}} \left(\langle \mu - v, \ell + \varepsilon p \log p \mu \rangle - \frac{1}{\eta} D_\psi(v, \mu) \right) \\
&= \langle \mu, \ell + \varepsilon p \log p \mu \rangle - \inf_{v \in \mathbb{R}_{\geq 0}^{XA}} \left(\langle v, \ell + \varepsilon p \log p \mu \rangle + \frac{1}{\eta} D_\psi(v, \mu) \right).
\end{aligned}$$

Further, the first-order optimality condition $\ell + \varepsilon p \log p \mu + \frac{1}{\eta} (\nabla \psi(v) - \nabla \psi(\mu)) = 0$ implies that

$$\log \frac{v_{1:h}(x_h, a_h)}{\mu_{1:h}(x_h, a_h)} = -\frac{\eta}{p_{1:h}^x(x_h)} [\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))].$$

Hence, one can see that

$$v_{1:h}(x_h, a_h) = \mu_{1:h}(x_h, a_h) \exp\left(-\frac{\eta}{p_{1:h}^x(x_h)} [\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))]\right). \quad (10)$$

Therefore, we have

$$\begin{aligned} & \langle \mu - \mu', \ell + \varepsilon p \log(p\mu) \rangle - \frac{1}{\eta} D_\psi(\mu', \mu) \\ &= \sum_{h, (x_h, a_h)} \left[(\mu_{1:h}(x_h, a_h) - v_{1:h}(x_h, a_h)) (\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))) \right. \\ & \quad \left. - \frac{1}{\eta} \left(p_{1:h}^x(x_h) v_{1:h}(x_h, a_h) \log \frac{v_{1:h}(x_h, a_h)}{\mu_{1:h}(x_h, a_h)} - p_{1:h}^x(x_h) (v_{1:h}(x_h, a_h) - \mu_{1:h}(x_h, a_h)) \right) \right] \\ &= \sum_{h, (x_h, a_h)} \left[(\mu_{1:h}(x_h, a_h)) (\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))) \right. \\ & \quad \left. + \frac{p_{1:h}^x(x_h)}{\eta} \left(\exp\left(-\frac{\eta}{p_{1:h}^x(x_h)} [\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))]\right) - 1 \right) \mu_{1:h}(x_h, a_h) \right] \\ &= \sum_{h, (x_h, a_h)} \frac{p_{1:h}^x(x_h)}{\eta} \mu_{1:h}(x_h, a_h) \left[\frac{\eta}{p_{1:h}^x(x_h)} (\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))) \right. \\ & \quad \left. + \left(\exp\left(-\frac{\eta}{p_{1:h}^x(x_h)} [\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))]\right) - 1 \right) \right] \\ &\leq \sum_{h, (x_h, a_h)} \frac{\eta}{p_{1:h}^x(x_h)} \mu_{1:h}(x_h, a_h) \ell^2(x_h, a_h) \\ & \quad + \sum_{h, (x_h, a_h)} \frac{p_{1:h}^x(x_h)}{\eta} \left[\mu_{1:h}(x_h, a_h) (\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))) \right. \\ & \quad \left. + \mu_{1:h}(x_h, a_h) \exp\left(-\frac{\eta}{p_{1:h}^x(x_h)} [\ell(x_h, a_h) + \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))]\right) \right. \\ & \quad \left. - \mu_{1:h}(x_h, a_h) \exp\left(-\frac{\eta}{p_{1:h}^x(x_h)} \ell(x_h, a_h)\right) \right] \\ &= \sum_{h, (x_h, a_h)} \frac{\eta}{p_{1:h}^x(x_h)} \mu_{1:h}(x_h, a_h) \ell^2(x_h, a_h) \\ & \quad + \sum_{h, (x_h, a_h)} \frac{1}{\eta} \left[\mu_{1:h}(x_h, a_h) \eta \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)) \right. \\ & \quad \left. + \exp\left(-\frac{\eta}{p_{1:h}^x(x_h)} \ell(x_h, a_h)\right) \left((p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))^{1-\eta \varepsilon(x_h, a_h)} - p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h) \right) \right] \\ &\leq \sum_{h, (x_h, a_h)} \frac{\eta}{p_{1:h}^x(x_h)} \mu_{1:h}(x_h, a_h) \ell^2(x_h, a_h) \\ & \quad + \sum_{h, (x_h, a_h)} \frac{1}{\eta} \left[\mu_{1:h}(x_h, a_h) \eta \varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)) \right. \\ & \quad \left. - \eta \varepsilon(x_h, a_h) (p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h))^{1-\eta \varepsilon(x_h, a_h)} \ln(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)) \right] \\ &\leq \sum_{h, (x_h, a_h)} \left[\frac{\eta}{p_{1:h}^x(x_h)} \mu_{1:h}(x_h, a_h) \ell^2(x_h, a_h) + \eta \varepsilon^2(x_h, a_h) \log^2(p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)) \right], \end{aligned}$$

where in the second equality we substitute $v_{1:h}(x_h, a_h)$ with Eq. (10), in first inequality comes from the fact that $\frac{\eta}{p_{1:h}^x(x_h)} \ell(x_h, a_h) \leq (\eta \ell(x_h, a_h) / p_{1:h}^x(x_h))^2 - \exp(\eta \ell(x_h, a_h) / p_{1:h}^x(x_h))$ and the

forth equality follows from

$$\begin{aligned} & \exp\left(-\frac{\eta}{p_{1:h}^x(x_h)}[\ell(x_h, a_h) + \varepsilon(x_h, a_h)p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))]\right) \\ &= \exp\left(-\frac{\eta}{p_{1:h}^x(x_h)}\ell(x_h, a_h)\right) \left((p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))^{1-\eta\varepsilon(x_h, a_h)}\right). \end{aligned}$$

The last two inequalities can be derived by following calculations:

$$\begin{aligned} & \exp\left(-\frac{\eta}{p_{1:h}^x(x_h)}\ell(x_h, a_h)\right) \left((p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))^{1-\eta\varepsilon(x_h, a_h)} - p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h)\right) \\ & \leq (p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))^{1-\eta\varepsilon(x_h, a_h)} - p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h) \\ & \leq \eta\varepsilon(x_h, a_h)(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))^{1-\eta\varepsilon(x_h, a_h)} \ln(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h)), \end{aligned}$$

and

$$\begin{aligned} & \mu_{1:h}(x_h, a_h) \eta\varepsilon(x_h, a_h) p_{1:h}^x(x_h) \log(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h)) \\ & \quad - \eta\varepsilon(x_h, a_h)(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))^{1-\eta\varepsilon(x_h, a_h)} \ln(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h)) \\ &= -\eta\varepsilon(x_h, a_h) \log(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h)) \left((p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))^{1-\eta\varepsilon(x_h, a_h)} - p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h)\right) \\ & \leq \eta^2\varepsilon^2(x_h, a_h)(\log^2(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h)))^{1-\eta\varepsilon(x_h, a_h)} \\ & \leq \eta^2\varepsilon^2(x_h, a_h)(\log^2(p_{1:h}^x(x_h)\mu_{1:h}(x_h, a_h))). \end{aligned}$$

□

Lemma B.2. $\forall k \geq 1$, we have

$$\text{NEGap}(\xi^{k,*}) = \mathcal{O}\left(\varepsilon_k H(\ln(XA) + \ln(YB)) + \frac{XAH}{k} + \frac{YBH}{k}\right). \quad (11)$$

Proof. $\forall (\mu', \nu') \in \Pi_{\max} \times \Pi_{\min}$, we have

$$\begin{aligned} & f(\mu^{k,*}, \nu') - f(\mu', \nu^{k,*}) \\ &= f(\mu^{k,*}, \nu') - f(\mu^{k,*}, \nu) + f(\mu^{k,*}, \nu) - f(\mu, \nu^{k,*}) + f(\mu, \nu^{k,*}) - f(\mu', \nu^{k,*}). \end{aligned}$$

First notice that $\forall (\mu, \nu) \in \Pi_{\max}^k \times \Pi_{\min}^k$,

$$\begin{aligned} & f(\mu^{k,*}, \nu) - f(\mu, \nu^{k,*}) \\ &= f(\mu^{k,*}, \nu) - f_k(\mu^{k,*}, \nu) + f_k(\mu^{k,*}, \nu) - f_k(\mu, \nu^{k,*}) + f_k(\mu, \nu^{k,*}) - f(\mu, \nu^{k,*}) \\ &= -(\varepsilon_k \psi(\mu^{k,*}) - \varepsilon_k \psi(\nu)) + (\varepsilon_k \psi(\mu) - \varepsilon_k \psi(\nu^{k,*})) \\ & \leq -\varepsilon_k \psi(\mu^{k,*}) - \varepsilon_k \psi(\nu^{k,*}) \\ & \leq \varepsilon_k H(\ln(XA) + \ln(YB)). \end{aligned}$$

To bound $f(\mu^{k,*}, \nu') - f(\mu^{k,*}, \nu)$, we have

$$\begin{aligned} & f(\mu^{k,*}, \nu') - f(\mu^{k,*}, \nu) \\ & \leq \langle \nabla_{\nu} f(\mu^{k,*}, \nu), \nu' - \nu \rangle \\ & \leq \|\nabla_{\nu} f(\mu^{k,*}, \nu)\|_1 \|\nu' - \nu\|_{\infty} \\ & \leq YB \left(1 - \left(1 - \frac{B-1}{Bk}\right)^H\right) \\ & \leq YB \left(1 - \left(1 - \frac{B}{Bk}\right)^H\right) \\ & \leq YB \left(1 - \left(1 - \frac{1}{k}\right)^H\right) \end{aligned}$$

$$= \mathcal{O}\left(\frac{YBH}{k}\right).$$

Similarly, we have

$$f(\mu, \nu^{k,*}) - f(\mu', \nu^{k,*}) \leq \mathcal{O}\left(\frac{XAH}{k}\right).$$

Putting all the above together completes the proof. \square

B.1 CONVERGENCE RATE OF THE CONTRACTION OF THE BREGMAN DIVERGENCE

Lemma B.3 (Contraction on Bregman divergence).

$$\begin{aligned} D_\psi(\xi^{k+1,*}, \xi^{k+1}) &\leq \underbrace{\sum_{i=1}^k w_k^i \eta_i \rho_i}_{\text{Term 1}} + \underbrace{\sum_{i=1}^k w_k^i \eta_i \sigma_i}_{\text{Term 2}} \\ &+ \underbrace{XA(\log X + H \log(Ak))^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2 + YB(\log Y + H \log(Bk))^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2}_{\text{Term 3}} \\ &+ \underbrace{\sum_{i=1}^k w_k^i \eta_i^2 (X\tau_i + Y\bar{\tau}_i)}_{\text{Term 4}} + \underbrace{\sum_{i=1}^k w_k^i \eta_i^2 (X^2A + Y^2B)}_{\text{Term 5}} + \underbrace{\sum_{i=1}^k w_k^i \omega_i}_{\text{Term 6}}. \end{aligned}$$

Proof. Recall we denote $[p^x \mu](x_h, a_h) := p_{1:h}^x(x_h) \mu_{1:h}(x_h, a_h)$.

$$f_k(\mu^k, \nu^k) - f_k(\mu^{k,*}, \nu^k) = (\mu^k - \mu^{k,*})^\top \mathbf{G} \nu^k + \varepsilon_k (\psi(\mu^k) - \psi(\mu^{k,*})).$$

For the first term in the above display, we have

$$\begin{aligned} &(\mu^k - \mu^{k,*})^\top \mathbf{G} \nu^k \\ &= (\mu^k - \mu^{k,*})^\top (\mathbf{G} \nu^k + g^k - g^k) \\ &= (\mu^k - \mu^{k,*})^\top g^k + (\mu^k)^\top (\mathbf{G} \nu^k - g^k) - (\mu^{k,*})^\top (\mathbf{G} \nu^k - g^k) \\ &= (\mu^k - \mu^{k,*})^\top g^k + \sum_{h, x_h, a_h} \mu_{1:k}^k(x_h, a_h) [(\mathbf{G} \nu^k)[(x_h, a_h)] - g^k[(x_h, a_h)]] \\ &\quad - \sum_{h, x_h, a_h} \mu_{1:h}^{k,*}(x_h, a_h) [(\mathbf{G} \nu^k)[(x_h, a_h)] - g^k[(x_h, a_h)]] \\ &= (\mu^k - \mu^{k,*})^\top g^k \\ &\quad + \sum_{h, x_h, a_h} \mu_{1:h}^k(x_h, a_h) \left[(\mathbf{G} \nu^k)[(x_h, a_h)] - \left(\frac{\mathbb{I}_h^k\{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} (1 - r_h^k) + \varepsilon_k p_{1:h}^x(x_h) \log [p^x \mu^k](x_h, a_h) \right) \right] \\ &\quad - \sum_{h, x_h, a_h} \mu_{1:h}^{k,*}(x_h, a_h) \left[(\mathbf{G} \nu^k)[(x_h, a_h)] - \left(\frac{\mathbb{I}_h^k\{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} (1 - r_h^k) + \varepsilon_k p_{1:h}^x(x_h) \log [p^x \mu^k](x_h, a_h) \right) \right]. \end{aligned}$$

For the second term, we have

$$\begin{aligned} &\psi(\mu^k) - \psi(\mu^{k,*}) \\ &= \sum_{h, x_h, a_h} [p^x \mu^k](x_h, a_h) \log [p^x \mu^k](x_h, a_h) \\ &\quad - \sum_{h, x_h, a_h} [p^x \mu^{k,*}](x_h, a_h) \log [p^x \mu^{k,*}](x_h, a_h) \end{aligned}$$

$$\begin{aligned}
&= \sum_{h, x_h, a_h} ([p^x \mu^k](x_h, a_h) - [p^x \mu^{k,*}](x_h, a_h)) \log [p^x \mu^k](x_h, a_h) \\
&\quad - \sum_{h, x_h, a_h} [p^x \mu^{k,*}](x_h, a_h) (\log [p^x \mu^{k,*}](x_h, a_h) - \log [p^x \mu^k](x_h, a_h)) \\
&= \sum_{h, x_h, a_h} ([p^x \mu^k](x_h, a_h) - [p^x \mu^{k,*}](x_h, a_h)) \log [p^x \mu^k](x_h, a_h) - D_\psi(\mu^{k,*}, \mu^k).
\end{aligned}$$

We then arrive at

$$\begin{aligned}
&f_k(\mu_k, v_k) - f_k(\mu_k^*, v_k) \\
&= (\mu^k - \mu^{k,*})^\top g^k + \underbrace{\sum_{h, x_h, a_h} \mu_{1:k}^k(x_h, a_h) \left[(\mathbf{G} \nu^k)[(x_h, a_h)] - \left(\frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} (1 - r_h^k) \right) \right]}_{=:\underline{\rho}_k} \\
&\quad - \underbrace{\sum_{h, x_h, a_h} \mu_{1:k}^{k,*}(x_h, a_h) \left[(\mathbf{G} \nu^k)[(x_h, a_h)] - \left(\frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} (1 - r_h^k) \right) \right]}_{=:\underline{\sigma}_k} - \varepsilon_k D_\psi(\mu^{k,*}, \mu^k) \\
&\leq \frac{1}{\eta_k} (D_\psi(\mu^{k,*}, \mu^k) - D_\psi(\mu^{k,*}, \mu^{k+1})) - \varepsilon_k D_\psi(\mu^{k,*}, \mu^k) + \underline{\rho}_k + \underline{\sigma}_k \\
&\quad + \sum_{h, x_h, a_h} \eta_k \left(\frac{1}{p_{1:h}^x(x_h)} \mu_{1:h}^k(x_h, a_h) \hat{\ell}_h^k(x_h, a_h)^2 + \varepsilon_k^2(x_h, a_h) \log^2(p_{1:h}^x(x_h) \mu_{1:h}^k(x_h, a_h)) \right) \\
&\leq \frac{(1 - \eta_k \varepsilon_k) D_\psi(\mu^{k,*}, \mu^k) - D_\psi(\mu^{k,*}, \mu^{k+1})}{\eta_k} + \underline{\rho}_k + \underline{\sigma}_k \\
&\quad + \sum_{h, x_h, a_h} \eta_k \left(\frac{1}{p_{1:h}^x(x_h)} \frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} + \varepsilon_k^2 \log^2 \left(\underbrace{p_{1:h}^x(x_h) \mu_{1:h}^k(x_h, a_h)}_{m := \min_{h, (x_h, a_h)} p_{1:h}^x(x_h) \mu_{1:h}^k(x_h, a_h)} \right) \right) \\
&\leq \frac{(1 - \eta_k \varepsilon_k) D_\psi(\mu^{k,*}, \mu^k) - D_\psi(\mu^{k,*}, \mu^{k+1})}{\eta_k} + \underline{\rho}_k + \underline{\sigma}_k \\
&\quad + \eta_k \underbrace{\sum_{h, x_h, a_h} \frac{1}{p_{1:h}^x(x_h)} \left(\frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} - 1 \right)}_{=:\underline{X} \tau_k} + \eta_k X^2 A + \eta_k X A \log^2 m \\
&\leq \frac{(1 - \eta_k \varepsilon_k) D_\psi(\mu^{k,*}, \mu^k) - D_\psi(\mu^{k,*}, \mu^{k+1})}{\eta_k} + \underline{\rho}_k + \underline{\sigma}_k + \eta_k \underline{X} \tau_k + \eta_k X^2 A + \eta_k X A \log^2 m.
\end{aligned}$$

Rearranging shows that

$$\begin{aligned}
&D_\psi(\mu^{k+1,*}, \mu^{k+1}) \\
&\leq (1 - \eta_k \varepsilon_k) D_\psi(\mu^{k,*}, \mu^k) + \eta_k (f_k(\mu^{k,*}, v_k) - f_k(\mu_k, v_k)) \\
&\quad + \eta_k^2 X A \log^2 m + \eta_k^2 X A \tau_k + \eta_k^2 X^2 A + \eta_k \underline{\rho}_k + \eta_k \underline{\sigma}_k + \underbrace{D_\psi(\mu^{k+1,*}, \mu^{k+1}) - D_\psi(\mu^{k,*}, \mu^{k+1})}_{=:\underline{\omega}_k}.
\end{aligned}$$

Analogously, for the min-player, we have

$$\begin{aligned}
&D_\psi(\nu^{k+1,*}, \nu^{k+1}) \\
&\leq (1 - \eta_k \varepsilon_k) D_\psi(\nu^{k,*}, \nu^k) + \eta_k (f_k(\mu^k, \nu^k) - f_k(\mu_k, \nu^{k,*})) \\
&\quad + \eta_k^2 Y B \log^2 m + \eta_k^2 Y \bar{\tau}_k + \eta_k^2 Y^2 B + \eta_k \bar{\rho}_k + \eta_k \bar{\sigma}_k + \bar{w}_k,
\end{aligned}$$

where

$$\bar{\tau}_k := \frac{1}{Y} \sum_{h, y_h, b_h} \frac{1}{p_{1:h}^y(y_h)} \left(\frac{\mathbb{I}_h^k \{y_h, b_h\}}{\nu_{1:h}^k(y_h, b_h) + \gamma_k} - 1 \right)$$

$$\begin{aligned}
1188 \quad \bar{\rho}_k &:= \sum_{h, y_h, b_h} \nu_{1:h}^k(y_h, b_h) \left[\left(1 - (\mathbf{G}^\top \mu^k) [(y_h, b_h)]\right) - \frac{\mathbb{I}_h^k \{y_h, b_h\} r_h^k}{\nu_{1:h}^k(y_h, b_h) + \gamma_{kk}} \right] \\
1189 \quad & \\
1190 \quad \bar{\sigma}_k &:= \sum_{h, y_h, b_h} \nu_{1:h}^{k,*}(y_h, b_h) \left[\frac{\mathbb{I}_h^k \{y_h, b_h\} r_h^k}{\nu_{1:h}^k(y_h, b_h) + \gamma_{kk}} - \left(1 - (\mathbf{G}^\top \mu^k) [(y_h, b_h)]\right) \right] \\
1191 \quad & \\
1192 \quad \bar{\omega}_k &:= D_\psi(\nu^{k+1,*}, \nu^{k+1}) - D_\psi(\nu^{k,*}, \nu^{k+1}) . \\
1193 \quad & \\
1194 \quad & \\
1195 \quad &
\end{aligned}$$

1196 Combining both sides and noticing that $f_k(\mu^{k,*}, \nu^k) - f_k(\mu^k, \nu^{k,*}) \leq 0$, we have

$$\begin{aligned}
1197 \quad & \\
1198 \quad & \\
1199 \quad D_\psi(\xi^{k+1,x}, \xi^{k+1}) \\
1200 \quad &\leq (1 - \eta_k \varepsilon_k) D_\psi(\xi^{k,*}, \xi^k) + \eta_k^2 (X \bar{\tau}_k + Y \bar{\tau}_k) + \eta_k^2 (X^2 A + Y^2 B) + \eta_k \rho_k + \eta_k \sigma_k + \omega_k \\
1201 \quad &+ \eta_k^2 X A \varepsilon_k^2 (\log X + H \log(Ak))^2 + \eta_k^2 Y B \varepsilon_k^2 (\log Y + H \log(Bk))^2 . \\
1202 \quad & \\
1203 \quad &
\end{aligned}$$

1204 Now expanding the recursion in the above display leads to

$$\begin{aligned}
1205 \quad & \\
1206 \quad D_\psi(\xi^{k+1,*}, \xi^{k+1}) &\leq \underbrace{\sum_{i=1}^k w_k^i \eta_i \rho_i}_{\text{Term 1}} + \underbrace{\sum_{i=1}^k w_k^i \eta_i \sigma_i}_{\text{Term 2}} \\
1207 \quad & \\
1208 \quad & \\
1209 \quad &+ \underbrace{X A (\log X + H \log(Ak))^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2 + Y B (\log Y + H \log(Bk))^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2}_{\text{Term 3}} \\
1210 \quad & \\
1211 \quad & \\
1212 \quad & \\
1213 \quad &+ \underbrace{\sum_{i=1}^k w_k^i \eta_i^2 (X \bar{\tau}_i + Y \bar{\tau}_i)}_{\text{Term 4}} + \underbrace{\sum_{i=1}^k w_k^i \eta_i^2 (X^2 A + Y^2 B)}_{\text{Term 5}} + \underbrace{\sum_{i=1}^k w_k^i \omega_i}_{\text{Term 6}} , \\
1214 \quad & \\
1215 \quad & \\
1216 \quad & \\
1217 \quad &
\end{aligned}$$

1218 where $w_k^i = \prod_{j=i+1}^k (1 - \eta_j \varepsilon_j)$. □

1220 B.2 BOUNDING CONTRACTION TERMS

1221 **Lemma B.4** (Bounding **Term 1**).

$$1222 \quad \text{Term 1} \leq (X A + Y B) \ln(k) k^{-\alpha_{\gamma_k} + \alpha_\varepsilon} + k^{-\frac{\alpha_k}{2} + \frac{\alpha_\varepsilon}{2}} \log\left(\frac{k^2}{\delta}\right) .$$

1226 *Proof.* Recall

$$\begin{aligned}
1227 \quad & \\
1228 \quad & \\
1229 \quad \text{Term 1} &= \sum_{i=1}^k w_k^i \eta_i \rho_i = \sum_{i=1}^k w_k^i \eta_i \underline{\rho}_i + \sum_{i=1}^k w_k^i \eta_i \bar{\rho}_i . \\
1230 \quad & \\
1231 \quad &
\end{aligned}$$

1232 To bound $\sum_{i=1}^k w_k^i \eta_i \underline{\rho}_i$, note that

$$\begin{aligned}
1233 \quad & \\
1234 \quad &\sum_{i=1}^k w_k^i \eta_i \underline{\rho}_i \\
1235 \quad & \\
1236 \quad &= \sum_{i=1}^k w_k^i \eta_i \langle \mu^i, \ell^{i,x} - \hat{\ell}^{i,x} \rangle \\
1237 \quad & \\
1238 \quad &= X A \sum_{i=1}^k w_k^i \eta_i \gamma_{ki} + H \sqrt{2 \sum_{i=1}^k (w_k^i \eta_i)^2 \log \frac{k^2}{\delta}} \\
1239 \quad & \\
1240 \quad & \\
1241 \quad &
\end{aligned}$$

$$\begin{aligned}
&\leq XA \sum_{i=1}^k \left[i^{-\alpha_{\gamma k} - \alpha_{\eta}} \prod_{j=i+1}^k (1 - j^{-\alpha_{\eta} - \alpha_{\varepsilon}}) \right] + \sqrt{\log\left(\frac{k^2}{\delta}\right) \sum_{i=1}^k \left[i^{-2\alpha_{\eta}} \left(\prod_{j=i+1}^k (1 - j^{-\alpha_{\eta} - \alpha_{\varepsilon}}) \right)^2 \right]} \\
&\leq XA \sum_{i=1}^k \left[i^{-\alpha_{\gamma} - \alpha_{\eta}} \prod_{j=i+1}^k (1 - j^{-\alpha_{\eta} - \alpha_{\varepsilon}}) \right] + \sqrt{\log\left(\frac{k^2}{\delta}\right) \sum_{i=1}^k \left[i^{-2\alpha_{\eta}} \left(\prod_{j=i+1}^k (1 - j^{-\alpha_{\eta} - \alpha_{\varepsilon}}) \right)^2 \right]} \\
&\leq XA \ln(k) k^{-\alpha_{\gamma} + \alpha_{\varepsilon}} + \sqrt{\log\left(\frac{k^2}{\delta}\right) \ln(k) k^{-\alpha_{\gamma} + \alpha_{\varepsilon}}} \\
&\leq XA \ln(k) k^{-\alpha_{\gamma} + \alpha_{\varepsilon}} + k^{-\frac{\alpha_{\eta}}{2} + \frac{\alpha_{\varepsilon}}{2}} \log\left(\frac{k^2}{\delta}\right),
\end{aligned}$$

where the second equality is given by Lemma B.13 and the third inequality comes from Lemma E.1.

Analogously, we have

$$\sum_{i=1}^k w_k^i \eta_i \bar{\rho}_i \leq YB \ln(k) k^{-\alpha_{\gamma} + \alpha_{\varepsilon}} + k^{-\frac{\alpha_k}{2} + \frac{\alpha_{\varepsilon}}{2}} \log\left(\frac{k^2}{\delta}\right).$$

Hence

$$\mathbf{Term 1} \leq (XA + YB) \ln(k) k^{-\alpha_{\gamma} + \alpha_{\varepsilon}} + k^{-\frac{\alpha_k}{2} + \frac{\alpha_{\varepsilon}}{2}} \log\left(\frac{k^2}{\delta}\right).$$

□

Lemma B.5 (Bounding Term 2).

$$\mathbf{Term 2} \leq k^{-\alpha_{\eta} + \alpha_{\gamma k}} \log \frac{k^2}{\delta}.$$

Proof.

$$\begin{aligned}
\mathbf{Term 2} &= \sum_{i=1}^k w_k^i \eta_i \sigma_i \\
&= \sum_{i=1}^k w_k^i \eta_i \underline{\sigma}_i + \sum_{i=1}^k w_k^i \eta_i \bar{\sigma}_i \\
&\leq \max_{1 \leq i \leq k} \frac{\eta_i w_k^i}{\gamma_{k k}} \log \frac{k^2}{\delta} \quad (\text{with probability } 1 - \frac{k^2}{\delta}) \\
&\leq k^{-\alpha_{\eta} + \alpha_{\gamma k}} \log \frac{k^2}{\delta},
\end{aligned}$$

where the last inequality is due to Lemma E.2.

□

Lemma B.6 (Bounding Term 3).

$$\mathbf{Term 3} \leq \left(XA \left(\log X + H \log(Ak) \right)^2 + YB \left(\log Y + H \log(Bk) \right)^2 \right) k^{-\alpha_{\eta} - \alpha_{\varepsilon}}.$$

Proof.

Term 3

$$\begin{aligned}
&= XA \left(\log X + H \log(Ak) \right)^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2 + YB \left(\log Y + H \log(Bk) \right)^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2 \\
&\leq \left(XA \left(\log X + H \log(Ak) \right)^2 + YB \left(\log Y + H \log(Bk) \right)^2 \right) k^{-2(\alpha_{\eta} + \alpha_{\varepsilon}) + \alpha_{\eta} + \alpha_{\varepsilon}} \\
&= \left(XA \left(\log X + H \log(Ak) \right)^2 + YB \left(\log Y + H \log(Bk) \right)^2 \right) k^{-\alpha_{\eta} - \alpha_{\varepsilon}},
\end{aligned}$$

where the inequality follows from Lemma E.1.

□

1296 **Lemma B.7** (Bounding **Term 4**).

$$1297 \quad \mathbf{Term\ 4} \leq k^{\alpha_{\gamma k} - 2\alpha_{\eta}} (X + Y) \log\left(\frac{1}{\delta}\right).$$

1299 *Proof.*

1300 **Term 4**

$$1301 \quad = \sum_{i=1}^k w_k^i \eta_i^2 (X \bar{\tau}_i + Y \bar{r}_i)$$

$$1302 \quad = \sum_{i=1}^k w_k^i \eta_i^2 \left(X \cdot \frac{1}{X} \sum_{h, x_h, a_h} \frac{1}{p_{1:h}^x(x_h)} \left(\frac{\mathbb{I}_h^k \{x_h, a_h\}}{\mu_{1:h}^k(x_h, a_h) + \gamma_k} - 1 \right) \right.$$

$$1303 \quad \left. + Y \cdot \frac{1}{Y} \sum_{h, y_h, b_h} \frac{1}{p_{1:h}^y(y_h)} \left(\frac{\mathbb{I}_h^k \{y_h, b_h\}}{\nu_{1:h}^k(y_h, b_h) + \gamma_k} - 1 \right) \right)$$

$$1304 \quad \leq \max_{1 \leq i \leq k} \frac{w_k^i \eta_i^2 (X + Y)}{\gamma_k} \log\left(\frac{1}{\delta}\right)$$

$$1305 \quad \leq k^{\alpha_{\gamma} - 2\alpha_{\eta}} (X + Y) \log\left(\frac{1}{\delta}\right),$$

1306 where the first inequality follows from that $\frac{1}{X} \frac{1}{p_{1:h}^x(x_h)} \leq 1$ for all (x_h, a_h) guaranteed by Lemma
1307 **A.1** together with the use of Lemma **B.15**. \square

1308 **Lemma B.8** (Bounding **Term 5**).

$$1309 \quad \mathbf{Term\ 5} = (X^2 A + Y^2 B) k^{-\alpha_{\eta} + \alpha_{\epsilon}}.$$

1310 *Proof.*

$$1311 \quad \mathbf{Term\ 5} = \sum_{i=1}^k w_k^i \eta_i^2 (X^2 A + Y^2 B)$$

$$1312 \quad \leq (X^2 A + Y^2 B) k^{-2\alpha_{\eta} + \alpha_{\eta} + \alpha_{\epsilon}}$$

$$1313 \quad = (X^2 A + Y^2 B) k^{-\alpha_{\eta} + \alpha_{\epsilon}},$$

1314 where the inequality is given by Lemma **E.1**. \square

1315 **Lemma B.9** (Bounding **Term 6**).

$$1316 \quad \mathbf{Term\ 6} \leq (X + Y)^{\frac{1}{2}} (H \log(Ak) + H \log(Bk)) \cdot \log(k) k^{-\min\{1, \frac{3}{2} - \frac{\alpha_{\epsilon}}{2}\} + \alpha_{\eta} + \alpha_{\epsilon}}.$$

1317 *Proof.* To begin with, note that $\min_{(x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}, h \in [H]} \mu_{1:h}^k(x_h, a_h) \geq \frac{1}{(Ak)^H}$ due to the definition
1318 of Π_{\max} in Algorithm **1**. Similarly, $\min_{(y_h, b_h) \in \mathcal{Y}_h \times \mathcal{B}, h \in [H]} \nu_{1:h}^k(y_h, b_h) \geq \frac{1}{(Bk)^H}$ holds for the
1319 min-player. Further,

1320 **Term 6**

$$1321 \quad = \sum_{i=1}^k w_k^i \omega_i$$

$$1322 \quad \leq (X + Y)^{\frac{1}{2}} \log\left(\frac{1}{(Ak)^H}\right) \log\left(\frac{1}{(Bk)^H}\right) \sum_{i=1}^k w_k^i i^{-\min\{1, \frac{3}{2} - \frac{\alpha_{\epsilon}}{2}\}}$$

$$1323 \quad \leq (X + Y)^{\frac{1}{2}} \log\left(\frac{1}{(Ak)^H}\right) \log\left(\frac{1}{(Bk)^H}\right) \log(k) k^{-\min\{1, \frac{3}{2} - \frac{\alpha_{\epsilon}}{2}\} + \alpha_{\eta} + \alpha_{\epsilon}}$$

$$1324 \quad \leq (X + Y)^{\frac{1}{2}} (H \log(Ak) + H \log(Bk)) \log(k) k^{-\min\{1, \frac{3}{2} - \frac{\alpha_{\epsilon}}{2}\} + \alpha_{\eta} + \alpha_{\epsilon}},$$

1325 where the first inequality is due to Lemma **B.10** and the second inequality comes from Lemma
1326 **E.1**. \square

B.3 BOUNDING THE NE GAP OF $(\mu^{k,*}, \nu^{k,*})$

Lemma B.10 (Bounding divergence difference).

$$|w_k| = \mathcal{O} \left(\frac{(X+Y)^{\frac{1}{2}} (\ln((Ak)^H) + \ln((Bk)^H))^2}{k^{\min\{1, \frac{3}{2} - \frac{\alpha\varepsilon}{2}\}}} \right).$$

Proof. Again, note that $\min_{(x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}, h \in [H]} \mu_{1:h}^k(x_h, a_h) \geq \frac{1}{(Ak)^H}$ and $\min_{(y_h, b_h) \in \mathcal{Y}_h \times \mathcal{B}, h \in [H]} \nu_{1:h}^k(y_h, b_h) \geq \frac{1}{(Bk)^H}$. Therefore, it holds that

$$\begin{aligned} & |w_k| \\ & \leq |D_\psi(\mu^{k+1,*}, \mu^{k+1}) - D_\psi(\mu^{k,*}, \mu^{k+1})| + |D_\psi(\nu^{k+1,*}, \nu^{k-1}) - D_\psi(\nu^{k,*}, \nu^{k+1})| \\ & \leq (\ln((Ak)^H) + \ln((Bk)^H)) (\|p^x \mu^{k+1,*} - p^x \mu^{k,*}\|_1 + \|p^y \nu^{k+1,*} - p^y \nu^{k,*}\|_1) \\ & \leq \mathcal{O} \left(\frac{(X+Y)^{\frac{1}{2}} (\ln((Ak)^H) + \ln((Bk)^H))^2}{k^{\min\{1, \frac{3}{2} - \frac{\alpha\varepsilon}{2}\}}} \right), \end{aligned}$$

where the second inequality is due to Lemma B.11 and the last inequality comes from Lemma B.12. \square

Lemma B.11 (Bounding divergence using ℓ_1 -norm). $\forall \mu, \mu^1, \mu^2 \in \Pi_{\max}^k$, it holds that

$$|D_\psi(\mu', \mu) - D_\psi(\mu^2, \mu)| \leq \mathcal{O}(\ln((Ak)^H) \|p^x \mu^1 - p^x \mu^2\|_1).$$

Proof.

$$\begin{aligned} & D_\psi(\mu', \mu) - D_\psi(\mu^2, \mu) \\ & = \sum_{h, (x_h, a_h)} p_{1:h}^x(x_h) \left(\mu_{1:h}^1(x_h, a_h) \log \frac{\mu_{1:h}^1(x_h, a_h)}{\mu_{1:h}(x_h \cdot a_h)} - \mu_{1:h}^2(x_h, a_h) \log \frac{\mu_{1:h}^2(x_h, a_h)}{\mu_{1:h}(x_h, a_h)} \right) \\ & = \sum_{h, (x_h, a_h)} p_{1:h}^x(x_h) \left((\mu_{1:h}^1(x_h, a_h) - \mu_{1:h}^2(x_h, a_h)) \log \frac{\mu_{1:h}^1(x_h, a_h)}{\mu_{1:h}(x_h \cdot a_h)} \right) \\ & \quad + \sum_{h, (x_h, a_h)} p_{1:h}^x(x_h) \mu_{1:h}^1(x_h, a_h) \left(\log \frac{\mu_{1:h}^1(x_h, a_h)}{\mu_{1:h}(x_h \cdot a_h)} - \log \frac{\mu_{1:h}^2(x_h, a_h)}{\mu_{1:h}(x_h \cdot a_h)} \right) \\ & \leq \mathcal{O}(\ln((Ak)^H) \|p^x \mu^1 - p^x \mu^2\|_1) - D_\psi(\mu^2, \mu^1) \\ & \leq \mathcal{O}(\ln((Ak)^H) \|p^x \mu^1 - p^x \mu^2\|_1). \end{aligned}$$

\square

Lemma B.12 (Bounding ℓ_1 -norm of the difference between $\mu^{k,*}$ and $\mu^{k+1,*}$). *The ℓ_1 -norm of the difference between $\mu^{k,*}$ and $\mu^{k+1,*}$ satisfies*

$$\|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1 = \mathcal{O} \left(\frac{(X+Y)^{\frac{1}{2}} (\ln((Ak)^H) + \ln((Bk)^H))}{k^{\min\{1, \frac{3}{2} - \frac{\alpha\varepsilon}{2}\}}} \right).$$

Proof. First note that, $\forall k, \forall (\mu, \nu) \in \Pi_{\max}^k \times \Pi_{\min}^k$, we have

$$\begin{aligned} & f_k(\mu, \nu^{k,*}) - f_k(\mu^{k,*}, \nu) \\ & = f_k(\mu, \nu^{k,*}) - f_k(\mu^{k,*}, \nu^{k,*}) + f_k(\mu^{k,*}, \nu^{k,*}) - f_k(\mu^{k,*}, \nu) \\ & \geq f_k(\mu, \nu^{k,*}) - f_k(\mu^{k,*}, \nu^{k,*}) - \nabla_\mu f_k(\mu^{k,*}, \nu^{k,*})^\top (\mu - \mu^{k,*}) \\ & \quad - f_k(\mu^{k,*}, \nu) - (-f_k(\mu^{k,*}, \nu^{k,*})) - (-\nabla_\nu f_k(\mu^{k,*}, \nu^{k,*})^\top (\nu - \nu^{k,*})) \\ & \geq \varepsilon_k D_\psi(\mu, \mu^{k,*}) + \varepsilon_k D_\psi(\nu, \nu^{k,*}) \end{aligned}$$

$$\begin{aligned}
1404 & = \varepsilon_k \text{KL} (p^x \mu, p^x \mu^{k,*}) + \varepsilon_k \text{KL} (p^y \nu, p^y \nu^{k,*}) \\
1405 & \geq \frac{1}{2} \varepsilon_k \left(\|p^x \mu - p^x \mu^{k,*}\|_1^2 + \|p^y \nu - p^y \nu^{k,*}\|_1^2 \right) \\
1406 & \geq \frac{1}{4} \varepsilon_k \|p^z \xi - p^z \xi^{k,*}\|_1^2.
\end{aligned}$$

1409 Let $\mu^{k+1,'} = p_{k+1} \bar{\mu} + (1 - p_{k+1}) \mu_{k+1}^*$. Then $\forall h, (x_h, a_h)$,

$$1411 \quad \mu^{k+1,'} (a_h | x_h) \geq p_{k+1} \frac{1}{A} + (1 - p_{k+1}) \frac{1}{A(k+1)^2} \geq \frac{1}{Ak^2},$$

1412 which means that $\mu^{k+1,'} \in \Pi_{\max}^k$. Similarly, we define $\nu^{k+1,'}$, which is such that $\nu^{k+1,'} \in \Pi_{\min}^k$.
1413 By previous analysis, we have

$$1414 \quad f_k (\mu^{k+1,'}, \nu^{k,*}) - f_k (\mu^{k,*}, \nu^{k+1,'}) \geq \frac{1}{4} \varepsilon_k \|p^z \xi^{k+1,'} - p^z \xi^{k,*}\|_1^2. \quad (12)$$

1415 On the other hand, since $(\mu^{k,*}, \nu^{k,*}) \in \Pi_{\max}^{k+1} \times \Pi_{\min}^{k+1}$, we have

$$1416 \quad f_{k+1} (\mu^{k,*}, \nu^{k+1,*}) - f_{k+1} (\mu^{k+1,*}, \nu^{k,*}) \geq \frac{1}{4} \varepsilon_{k+1} \|p^z \xi^{k,*} - p^z \xi^{k+1,*}\|_1^2. \quad (13)$$

1417 Combing both sides, we have

$$\begin{aligned}
1418 & f_k (\mu^{k+1,*}, \nu^{k,*}) - f_k (\mu^{k,*}, \nu^{k+1,*}) \\
1419 & = f_k (\mu^{k+1,'}, \nu^{k,*}) - f_k (\mu^{k,*}, \nu^{k+1,'}) + f_k (\mu^{k+1,*}, \nu^{k,*}) - f_k (\mu^{k+1,'}, \nu^{k,*}) \\
1420 & \quad + f_k (\mu^{k,*}, \nu^{k+1,'}) - f_k (\mu^{k,*}, \nu^{k+1,*}) \\
1421 & \geq \frac{1}{4} \varepsilon_k \|p^z \xi^{k+1,'} - p^z \xi^{k,*}\|_1^2 + \langle \nabla_{\mu} f_k (\mu^{k+1,'}, \nu^{k,*}), \mu^{k+1,*} - \mu^{k+1,'} \rangle \\
1422 & \quad + \langle \nabla_{\nu} f_k (\mu^{k,*}, \nu^{k+1,'}), \nu^{k+1,'} - \nu^{k+1,*} \rangle \\
1423 & \geq \frac{1}{4} \varepsilon_k \|p^z \xi^{k+1,'} - p^z \xi^{k,*}\|_1^2 - \sup_{\mu \in \Pi_{\max}^{k+1}} \|\nabla_{\mu} f_k (\mu, \nu^{k,*})\|_{\infty} \|\mu^{k+1,*} - \mu^{k+1,'}\|_1 \\
1424 & \quad - \sup_{\nu \in \Pi_{\min}^{k+1}} \|\nabla_{\nu} f_k (\mu^{k,*}, \nu)\|_{\infty} \|\nu^{k+1,'} - \nu^{k+1,*}\|_1.
\end{aligned}$$

1425 Further using the fact that

$$\begin{aligned}
1426 & \|\nabla_{\mu} f_k (\mu, \nu^{k,*})\|_{\infty} \\
1427 & = \max_{h, (x_h, a_h)} |\mathbf{G} \nu^{k,*} [(x_h, a_h)] + \varepsilon_k p_{1:h}^x (x_h) \log [p^x \mu] [(x_h, a_h)]| \\
1428 & \leq \max_{h, (x_h, a_h)} |\mathbf{G} \nu^{k,*} [(x_h, a_h)]| + |\varepsilon_k p_{1:h}^x (x_h) \log [p^x \mu] [(x_h, a_h)]| \\
1429 & \leq 1 + k^{-\alpha \varepsilon} (\ln ((Ak)^H) + \ln ((Bk)^H)) = \mathcal{O}(1),
\end{aligned}$$

1430 and

$$\begin{aligned}
1431 & \|\mu^{k+1,*} - \mu^{k+1,'}\|_1 = \|p_{k+1} (\bar{\mu} - \mu_{k+1}^*)\|_1 \leq \|p_{k+1} \bar{\mu}\|_1 + \|p_{k+1} \mu_{k+1}^*\|_1 \\
1432 & \leq p_{k+1} 2X = \mathcal{O} \left(\frac{X+Y}{k^2} \right),
\end{aligned}$$

1433 one can deduce that

$$\begin{aligned}
1434 & f_k (\mu^{k+1,*}, \nu^{k,*}) - f_k (\mu^{k,*}, \nu^{k+1,*}) \\
1435 & \geq \frac{1}{8} \varepsilon_k \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1^2 - \frac{1}{4} \varepsilon_k \|p^z \xi^{k+1,*} - p^z \xi^{k+1,*}\|_1^2 - \mathcal{O} \left(\frac{X+Y}{k^3} \right) \\
1436 & \geq \frac{1}{8} \varepsilon_k \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1^2 - \frac{1}{4} \varepsilon_k \left(2 \left(\|p^x \mu^{k+1,'} - p^x \mu^{k+1,*}\|_1^2 + \|p^y \nu^{k+1,'} - p^y \nu^{k+1,*}\|_1^2 \right) \right) \\
1437 & \quad - \mathcal{O} \left(\frac{X+Y}{k^3} \right)
\end{aligned}$$

$$\begin{aligned}
&\geq \frac{1}{8} \varepsilon_k \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1^2 - \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&\quad - \frac{1}{4} \varepsilon_k \left(4 \left(\|p_{k+1} p^x \bar{\mu}\|_1^2 + \|p_{k+1} p^x \mu^{k+1,*}\|_1^2\right) + 4 \left(\|p_{k+1} p^y \bar{\nu}\|_1^2 + \|p_{k+1} p^y \nu^{k+1,*}\|_1^2\right)\right) \\
&= \frac{1}{8} \varepsilon_k \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1^2 - \mathcal{O}\left(\frac{X+Y}{k^3}\right) - \mathcal{O}\left(\frac{1}{k^6}\right) \\
&\geq \frac{1}{8} \varepsilon_{k+1} \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1^2 - \mathcal{O}\left(\frac{X+Y}{k^3}\right).
\end{aligned}$$

Combining with Eq. (13), we have

$$\begin{aligned}
&\frac{3}{8} \varepsilon_{k+1} \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1^2 \\
&\leq f_{k+1}(\mu^{k,*}, \nu^{k+1,*}) - f_k(\mu^{k,*}, \nu^{k+1,*}) - f_{k+1}(\mu^{k,*}, \nu^{k+1,*}) + f_k(\mu^{k+1,*}, \nu^{k,*}) + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&= \bar{f}_k(\mu^{k,*}, \nu^{k+1,*}) - \bar{f}_k(\mu^{k+1,*}, \nu^{k,*}) + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \quad (\bar{f}_k(\mu, \nu) := f_{k+1}(\mu, \nu) - f_k(\mu, \nu)) \\
&= \bar{f}_k(\mu^{k,*}, \nu^{k+1,*}) - \bar{f}_k(\mu^{k+1,*}, \nu^{k+1,*}) + \bar{f}_k(\mu^{k+1,*}, \nu^{k+1,*}) - \bar{f}_k(\mu^{k+1,*}, \nu^{k,*}) + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&\leq \langle \nabla_{\mu} \bar{f}_k(\mu^{k,*}, \nu^{k+1,*}), \mu^{k,*} - \mu^{k+1,*} \rangle + \langle \nabla_{\nu} \bar{f}_k(\mu^{k+1,*}, \nu^{k,*}), \nu^{k+1,*} - \nu^{k,*} \rangle + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&= \langle \nabla_{\mu} \bar{f}_k(\mu^{k,*}, \nu^{k+1,*}) / p^x, p^x (\mu^{k,*} - \mu^{k+1,*}) \rangle + \langle \nabla_{\nu} \bar{f}_k(\mu^{k+1,*}, \nu^{k,*}) / p^y, p^y (\nu^{k+1,*} - \nu^{k,*}) \rangle + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&\leq \|\nabla_{\mu} \bar{f}_k(\mu^{k,*}, \nu^{k+1,*}) / p^x\|_{\infty} \|p^x (\mu^{k,*} - \mu^{k+1,*})\|_1 + \|\nabla_{\nu} \bar{f}_k(\mu^{k+1,*}, \nu^{k,*}) / p^y\|_{\infty} \|p^y (\nu^{k+1,*} - \nu^{k,*})\|_1 \\
&\quad + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&\leq \left(\sup_{\mu \in \Pi_{\max}^k} \|\nabla_{\mu} \bar{f}_k(\mu, \nu^{k+1,*}) / p^x\|_{\infty} + \sup_{\nu \in \Pi_{\min}^k} \|\nabla_{\nu} \bar{f}_k(\mu^{k+1,*}, \nu) / p^y\|_{\infty} \right) \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1 + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&\leq \left(\sup_{\mu \in \Pi_{\max}^k} \max_{h, (x_h, a_h)} |(\varepsilon_k - \varepsilon_{k+1}) \log[p^x \mu][x_h, a_h]| + \sup_{\nu \in \Pi_{\min}^k} \max_{h, (y_h, b_h)} |(\varepsilon_k - \varepsilon_{k+1}) \log[p^y \nu][y_h, b_h]| \right) \\
&\quad \cdot \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1 + \mathcal{O}\left(\frac{X+Y}{k^3}\right) \\
&= \mathcal{O}\left((\varepsilon_k - \varepsilon_{k+1}) (\ln((Ak)^H) + \ln((Bk)^H)) \|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1 + \frac{X+Y}{k^3} \right).
\end{aligned}$$

In what follows, we slightly abuse the notations by denoting $m_k = (Ak)^H (Bk)^H$. Solving the above equation leads to

$$\begin{aligned}
&\|p^z \xi^{k+1,*} - p^z \xi^{k,*}\|_1 \\
&\leq \frac{(\varepsilon_k - \varepsilon_{k+1}) \log(m_k) + \sqrt{(\varepsilon_k - \varepsilon_{k+1})^2 \log^2(m_k) + \varepsilon_{k+1} \frac{X+Y}{k^3}}}{\varepsilon_{k+1}} \\
&\leq \frac{(\varepsilon_k - \varepsilon_{k+1})}{\varepsilon_{k+1}} \log(m_k) + \sqrt{\frac{X+Y}{\varepsilon_{k+1} k^3}} \\
&\leq \frac{\log(m_k)}{k} + \sqrt{\frac{X+Y}{\varepsilon_{k+1} k^3}} = \mathcal{O}\left(\frac{(X+Y)^{\frac{1}{2}} \log(m_k)}{k^{\min\{1, \frac{3}{2} - \frac{\alpha}{2}\}}}\right).
\end{aligned}$$

In the last inequality of the above display, we use the fact that

$$\frac{(\varepsilon_k - \varepsilon_{k+1})}{\varepsilon_{k+1}} = \frac{k^{-\alpha_{\epsilon}}}{(k+1)^{-\alpha_{\epsilon}}} = \left(1 + \frac{1}{k}\right)^{\alpha_{\epsilon}} - 1 = \mathcal{O}\left(\frac{\alpha_{\epsilon}}{k}\right),$$

by Taylor expansion. \square

Lemma B.13. Let $\{c_i\}_{i=1}^k$ be fixed positive numbers. Then with probability at least $1 - \delta$, it holds that

$$\sum_{i=1}^k c_i \langle \mu^i, \ell^{i,x} - \hat{\ell}^{i,x} \rangle \leq XA \sum_{i=1}^k c_i \gamma_{k_i} + H \sqrt{2 \sum_{i=1}^k c_i^2 \log \frac{1}{\delta}}.$$

Proof. To begin with, notice that

$$\sum_{i=1}^k c_i \langle \mu^i, \ell^{i,x} - \hat{\ell}^{i,x} \rangle = \sum_{i=1}^k c_i \langle \mu^i, \ell^{i,x} - \mathbb{E}_{i-1} [\hat{\ell}^{i,x}] \rangle + \sum_{i=1}^k c_i \langle \mu^i, \mathbb{E}_{i-1} [\hat{\ell}^{i,x}] - \hat{\ell}^{i,x} \rangle.$$

For the first part, we have

$$\begin{aligned} & \sum_{i=1}^k c_i \langle \mu^i, \ell^{i,x} - \mathbb{E}_{i-1} [\hat{\ell}^{i,x}] \rangle \\ &= \sum_{i=1}^k c_i \sum_{h, x_h, a_h} \mu_{1:h}^i(x_h, a_h) \ell_{[(x_h, a_h)]}^{i,x} \left(1 - \frac{\mu_{1:h}^i(x_h, a_h)}{\mu_{1:h}^i(x_h, a_h) + \gamma_{k_i}} \right) \\ &\leq \sum_{i=1}^k c_i \gamma_{k_i} \sum_{h, x_h, a_h} \ell_{[(x_h, a_h)]}^{i,x} \\ &\leq \sum_{i=1}^k c_i \gamma_{k_i} XA, \end{aligned}$$

where the last inequality comes from $\ell_{[(x_h, a_h)]}^{i,x} \leq 1$ for all $(x_h, a_h) \in \mathcal{X} \times \mathcal{A}$.

For the second part, taking $\delta = \exp\left(\frac{-\varepsilon^2}{2 \sum_{i=1}^k c_i^2 H^2}\right)$, $\varepsilon = \sqrt{2 \sum_{i=1}^k c_i^2 H^2 \log\left(\frac{1}{\delta}\right)}$ and using Azuma-Hoeffding inequality, it holds with probability at least $1 - \delta$ that

$$\sum_{i=1}^k c_i \langle \mu^i, \mathbb{E}_{i-1} [\hat{\ell}^{i,x}] - \hat{\ell}^{i,x} \rangle \leq \sqrt{2 \sum_{i=1}^k c_i^2 H^2 \log\left(\frac{1}{\delta}\right)}.$$

The proof is concluded by combining the upper bounds of the two parts above. \square

Lemma B.14. Let $\delta \in (0, 1)$ and $\{\gamma_{k_i}\}_{i=1}^k \in (0, +\infty)^k$. Fix $h \in [H]$. For any coefficient sequence $\{c^i\}_{i=1}^k$ s.t. $c^i \in [0, 2\gamma_{k_i}]^{XA}$ is \mathcal{F}_{i-1} -measurable, with probability $1 - \delta$, we have

$$\sum_{i=r}^k w_i \langle c_i, \hat{\ell}_i - \ell_i \rangle \leq \max_{1 \leq i \leq k} w_i \log \frac{1}{\delta}.$$

Proof. Define $w = \max_{1 \leq i \leq k} w_i$. Hence

$$\begin{aligned} & w^i \hat{\ell}^i(x_h, a_h) \\ &= \frac{w_i \mathbb{I}_{i,h}\{x_h, a_h\} (1 - r_h^i)}{\mu_{1:h}^i(x_h, a_h) + r_i} \\ &\leq \frac{w_i \mathbb{I}_{i,h}\{x_h, a_h\} (1 - r_h^i)}{\mu_{i,h}^i(x_h, a_h) + r_i \frac{w_i (1 - r_h^i) \mathbb{I}_{i,h}\{x_h, a_h\}}{w}} \\ &= \frac{w}{2\gamma_{k_i}} \frac{2\gamma_{k_i} w_i (1 - r_h^i) \mathbb{I}_{i,h}\{x_h, a_h\}}{w \mu_{i,h}^i(x_h, a_h)} \\ &= \frac{w}{2\gamma_{k_i}} \frac{\gamma_{k_i} w_i (1 - r_h^i) \mathbb{I}_{i,h}\{x_h, a_h\}}{w \mu_{i,h}^i(x_h, a_h)} \end{aligned}$$

$$\leq \frac{w}{2\gamma_{k_i}} \log \left(1 + \frac{2\gamma_{k_i} w_i (1 - r_h^i) \mathbb{I}_{i,h} \{x_h, a_h\}}{w \mu_{i:h}^i(x_h, a_h)} \right).$$

Denote by $\hat{S}_h^i = \frac{w_i}{w} \langle c^i, \hat{\ell}_h^i \rangle$, $S_h^i = \frac{w_i}{w} \langle c^i, \ell_h^i \rangle$. Then

$$\begin{aligned} & \mathbb{E}_{i-1}[\exp(\hat{S}^i)] \\ & \leq \mathbb{E}_{i-1} \left[\exp \left(\sum_{(x_h, a_h) \in \mathcal{X} \times \mathcal{A}} \frac{c^i(x_h, a_h)}{2\gamma_{k_i}} \log \left(1 + \frac{2\gamma_{k_i} w_i (1 - r_h^i) \mathbb{I}_{i,h} \{x_h, a_h\}}{w \mu_{i:h}^i(x_h, a_h)} \right) \right) \right] \\ & \leq \mathbb{E}_{i-1} \left[\prod_{(x_h, a_h) \in \mathcal{X} \times \mathcal{A}} \left(1 + \frac{c_i(x_h, a_h) w_i (1 - r_h^i) \mathbb{I}_{i,h} \{x_h, a_h\}}{w \mu_{i:h}^i(x_h, a_h)} \right) \right] \\ & = \mathbb{E}_{i-1} \left[1 + \sum_{(x_h, a_h) \in \mathcal{X} \times \mathcal{A}} \frac{c^i(x_h, a_h) w_i (1 - r_h^i) \mathbb{I}_{i,h} \{x_h, a_h\}}{w \mu_{i:h}^i(x_h, a_h)} \right] \\ & = 1 + S_h^i \leq \exp(S_h^i). \end{aligned}$$

Finally, one can see that

$$\begin{aligned} & \mathbb{E} \left[\sum_{i=1}^k (\hat{S}_h^i - S_h^i) \geq \log \frac{1}{\delta} \right] \\ & = \mathbb{E} \left[\exp \left(\sum_{i=1}^k (\hat{S}_h^i - S_h^i) \right) \geq \frac{1}{\delta} \right] \\ & \leq \delta \mathbb{E} \left[\exp \left(\sum_{i=1}^k (\hat{S}_h^i - S_h^i) \right) \right] \\ & = \delta \mathbb{E} \left[\mathbb{E}_{k-1} \left[\exp \left(\sum_{i=1}^k (\hat{S}_h^i - S_h^i) \right) \right] \right] \\ & = \delta \mathbb{E} \left[\exp \left(\sum_{i=1}^{k-1} (\hat{S}_h^i - S_h^i) \right) \left[\mathbb{E}_{k-1} \left[\exp(\hat{S}_h^k - S_h^k) \right] \right] \right] \leq \dots \leq \delta. \end{aligned}$$

□

Lemma B.15. Let $\{c_i\}_{i=1}^k$ be fixed positive numbers. Fix $h \in [H]$. Then \forall sequence $\{q_i\}_{i=1}^k \in [0, 1]^{X^A}$ s.t. q^i is \mathcal{F}_{i-1} -measurable, with probability at least $1 - \delta$,

$$\sum_{i=1}^k c_i \langle q_i, \hat{\ell}_h^i - \ell_h^i \rangle \leq \max_{1 \leq i \leq k} \frac{c_i}{\gamma_{k_i}} \log \left(\frac{1}{\delta} \right).$$

Proof. Noticing that $\{\gamma_{k_i}\}_{i=1}^k$ is decreasing and $\|q^i\|_\infty \leq 1$, applying Lemma B.14, we arrive at

$$\sum_{i=1}^k c_i \langle q^i, \hat{\ell}_h^i - \ell_h^i \rangle = \sum_{i=1}^k \frac{c_i}{2\gamma_{k_i}} \langle 2\gamma_{k_i} q^i, \hat{\ell}_h^i - \ell_h^i \rangle \leq \max_{1 \leq i \leq k} \frac{c_i}{\gamma_{k_i}} \log \left(\frac{1}{\delta} \right).$$

□

B.4 PROOF OF THEOREM 5.1

We are now ready to present the proof of our main result.

1620 *Proof of Theorem 5.1.* Putting Lemma B.3, B.4, B.5, B.6, B.7, B.8, B.9 together, we have

$$\begin{aligned}
& D_\psi(\xi^{k+1,*}, \xi^{k+1}) \\
&= \mathcal{O} \left((XA + YB) \ln(k) k^{-\alpha_{\gamma_k} + \alpha_\epsilon} + k^{-\frac{\alpha_\eta}{2} + \frac{\alpha_\epsilon}{2}} \log\left(\frac{k^2}{\delta}\right) + k^{-\alpha_\eta + \alpha_{\gamma_k}} \log\left(\frac{k^2}{\delta}\right) \right. \\
&\quad + \left(XA \left(\log X + H \log(Ak)^2 + YB \left(\log Y + H \log(Bk) \right)^2 \right) k^{-\alpha_\eta - \alpha_\epsilon} \right. \\
&\quad + k^{\alpha_{\gamma_k} - 2\alpha_\eta} (X + Y) \log\left(\frac{1}{\delta}\right) + (X^2A + Y^2B) k^{-\alpha_\eta + \alpha_\epsilon} \\
&\quad + (X + Y)^{\frac{1}{2}} (H \log(Ak) + H \log(Bk)) (\log X + H \log(k) + \log Y + H \log(Bk)) \\
&\quad \left. \cdot \log(k) k^{-\min\{1, \frac{3}{2} - \frac{\alpha_\epsilon}{2}\} + \alpha_\eta + \alpha_\epsilon} \right) \\
&= \mathcal{O} \left(\left[k^{-\frac{1}{4}} (XA + YB) + k^{-\frac{1}{4}} + k^{-\frac{1}{4}} + (XA + YB) H^2 k^{-\frac{3}{4}} + (X + Y) k^{-\frac{7}{8}} + (X^2A + Y^2B) k^{-\frac{1}{2}} \right. \right. \\
&\quad \left. \left. + (X + Y)^{\frac{1}{2}} H^2 K^{-\frac{1}{4}} \right] (\log^2(XAk/\delta) + \log^2(YBk/\delta)) \log(K) \right).
\end{aligned}$$

1638 Moreover, note that

$$\begin{aligned}
& \text{NEGap}(\xi^k) \\
&= \sup_{\mu \in \Pi_{\max}, \nu \in \Pi_{\min}} f(\mu^k, \nu) - f(\mu, \nu^k) \\
&= f(\mu^{k,*}, \nu) - f(\mu^{k,*}, \nu) + f(\mu^k, \nu) - f(\mu, \nu^k) + f(\mu, \nu^{k,*}) - f(\mu, \nu^{k,*}) \\
&\leq \text{NEGap}(\xi^{k,*}) + (\mu^k - \mu^{k,*})^\top \mathbf{G}\nu + \mu^\top \mathbf{G}(\nu^{k,*} - \nu^k) \\
&\leq \text{NEGap}(\xi^{k,*}) + \langle p^x(\mu^k - \mu^{k,*}), \mathbf{G}\nu/p^x \rangle + \langle p^y(\nu^k - \nu^{k,*}), \mathbf{G}^\top \mu/p^y \rangle \\
&\leq \text{NEGap}(\xi^{k,*}) + \|p^x(\mu^k - \mu^{k,*})\|_1 \|\mathbf{G}\nu/p^x\|_\infty + \|p^y(\nu^k - \nu^{k,*})\|_1 \|\mathbf{G}^\top \mu/p^y\|_\infty \\
&\leq \text{NEGap}(\xi^{k,*}) + X \|p^x(\mu^k - \mu^{k,*})\|_1 + Y \|p^y(\nu^k - \nu^{k,*})\|_1 \\
&\leq \varepsilon_k H (\ln(XA) + \ln(YB)) + \mathcal{O} \left(\frac{XAH}{k} + \frac{YBH}{k} \right) \\
&\quad + \mathcal{O} \left(X \sqrt{\text{KL}(p^x \mu^{k,*}, p^x \mu^k)} + Y \sqrt{\text{KL}(p^y \nu^{k,*}, p^y \nu^k)} \right) \\
&\leq \varepsilon_k H (\ln(XA) + \ln(YB)) + \mathcal{O} \left(\frac{XAH}{k} + \frac{YBH}{k} \right) \\
&\quad + \mathcal{O} \left((X + Y) \sqrt{\text{KL}(p^x \mu^{k,*}, p^x \mu^k) + \text{KL}(p^y \nu^{k,*}, p^y \nu^k)} \right) \\
&\leq \varepsilon_k H (\ln(XA) + \ln(YB)) + \mathcal{O} \left(\frac{XAH}{k} + \frac{YBH}{k} \right) + \mathcal{O} \left((X + Y) \sqrt{\text{KL}(p^z \xi^{k,*}, p^z \xi^k)} \right) \\
&\leq \varepsilon_k H (\ln(XA) + \ln(YB)) + \mathcal{O} \left(\frac{XAH}{k} + \frac{YBH}{k} \right) + \mathcal{O} \left((X + Y) \sqrt{D_\psi(\xi^{k,*}, \xi^k)} \right),
\end{aligned}$$

1664 where $\mathbf{G}\nu/p^x \in \mathbb{R}^{X^A}$ is defined as $(\mathbf{G}\nu/p^x)[(x_h, a_h)] = (\mathbf{G}\nu)[(x_h, a_h)]/p_{1:h}^x(x_h)$ and similarly
1665 for $\mathbf{G}^\top \mu/p^y$.

1666 Therefore, we can see that

$$\begin{aligned}
& \text{NEGap}(\mu^k, \nu^k) \\
&= \mathcal{O} \left((X + Y) \left[k^{-\frac{1}{8}} (XA + YB)^{\frac{1}{2}} + (XA + YB)^{\frac{1}{2}} H k^{-\frac{3}{8}} + (X^2A + Y^2B)^{\frac{1}{2}} k^{-\frac{1}{4}} + (X + Y)^{\frac{1}{4}} H k^{-\frac{1}{8}} \right] \right. \\
&\quad \left. \cdot (\log(XAk/\delta) + \log(YBk/\delta)) \log^{\frac{1}{2}}(k) + k^{-\frac{1}{8}} H (\ln(XA) + \ln(YB)) + \frac{XAB}{k} + \frac{YBH}{k} \right) \\
&= \tilde{\mathcal{O}} \left((X + Y) \left[k^{-\frac{1}{8}} (XA + YB)^{\frac{1}{2}} + (XA + YB)^{\frac{1}{2}} H k^{-\frac{3}{8}} + (X^2A + Y^2B)^{\frac{1}{2}} k^{-\frac{1}{4}} + (X + Y)^{\frac{1}{4}} H k^{-\frac{1}{8}} \right] \right)
\end{aligned}$$

$$\begin{aligned}
& + \frac{(XAH + YBH)}{k} \\
& = \tilde{\mathcal{O}} \left((X + Y)k^{-\frac{1}{8}} \left[(XA + YB)^{\frac{1}{2}} + (X + Y)^{\frac{1}{4}}H \right] \right),
\end{aligned}$$

where the last equality holds when $k \geq \max\{H^4, (X^2A + Y^2B)^4 / (XA + YB)^4, (XA + YB)^{8/7} / (X + Y)^{10/7}\}$. \square

C LAST-ITERATE CONVERGENCE RATE IN EXPECTATION

Theorem C.1. *With the same condition as in Theorem 5.1, Algorithm 1 guarantees that*

$$\mathbb{E} [\text{NEGap}(\mu^k, \nu^k)] = \tilde{\mathcal{O}} \left(\left((X + Y)^{\frac{1}{4}}H + \sqrt{(X^2A + Y^2B)} \right) k^{-\frac{1}{6}} \right).$$

Proof. With the same arguments as in the proof of Theorem 5.1, we have

$$\begin{aligned}
& D_\psi(\xi^{k+1,x}, \xi^{k+1}) \\
& \leq (1 - \eta_k \varepsilon_k) D_\psi(\xi^{k,\star}, \xi^k) + \eta_k^2 (X\bar{\tau}_k + Y\bar{\tau}_k) + \eta_k^2 (X^2A + Y^2B) + \eta_k \rho_k + \eta_k \sigma_k + \omega_k \\
& \quad + \eta_k^2 X A \varepsilon_k^2 (\log X + H \log(Ak))^2 + \eta_k^2 Y B \varepsilon_k^2 (\log Y + H \log(Bk))^2.
\end{aligned}$$

Taking conditional expectation $\mathbb{E}_{k-1}[\cdot]$ on both sides and by noticing the fact that $\mathbb{E}_{k-1}[\tau_k] < 0$, $\mathbb{E}_{k-1}[\rho_k] = 0$, and $\mathbb{E}_{k-1}[\sigma_k] = 0$, we have

$$\begin{aligned}
& \mathbb{E}_{k-1} [D_\psi(\xi^{k+1,x}, \xi^{k+1})] \\
& \leq (1 - \eta_k \varepsilon_k) D_\psi(\xi^{k,\star}, \xi^k) + \eta_k^2 (X^2A + Y^2B) + \mathbb{E}_{k-1}[\omega_k] \\
& \quad + \eta_k^2 X A \varepsilon_k^2 (\log X + H \log(Ak))^2 + \eta_k^2 Y B \varepsilon_k^2 (\log Y + H \log(Bk))^2.
\end{aligned}$$

Expanding the recursion in the above display leads to

$$\begin{aligned}
& \mathbb{E} [D_\psi(\xi^{k+1,\star}, \xi^{k+1})] \\
& \leq \mathbb{E} \left[\sum_{i=1}^k w_k^i \omega_i \right] + XA (\log X + H \log(Ak))^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2 + YB (\log Y + H \log(Bk))^2 \sum_{i=1}^k w_k^i (\eta_i \varepsilon_i)^2 \\
& \quad + \sum_{i=1}^k w_k^i \eta_i^2 (X^2A + Y^2B) \\
& \leq (X + Y)^{\frac{1}{2}} (H \log(Ak) + H \log(Bk)) (\log X + H \log(Ak) + \log Y + H \log(Bk)) \\
& \quad \cdot \log(k) k^{-\min\{1, \frac{3}{2} - \frac{\alpha_\varepsilon}{2}\} - \alpha_\eta + \alpha_\varepsilon} \\
& \quad + \left(XA (\log X + H \log(Ak))^2 + YB (\log Y + H \log(Bk))^2 \right) k^{-\alpha_\eta - \alpha_\varepsilon} + (X^2A + Y^2B) k^{-\alpha_\eta + \alpha_\varepsilon} \\
& = \tilde{\mathcal{O}} \left((X + Y)^{\frac{1}{2}} H^2 k^{-\min\{1, \frac{3}{2} - \frac{\alpha_\varepsilon}{2}\} + \alpha_\eta + \alpha_\varepsilon} + (XA + YB) H^2 k^{-\alpha_\eta - \alpha_\varepsilon} + (X^2A + Y^2B) k^{-\alpha_\eta + \alpha_\varepsilon} \right).
\end{aligned}$$

Hence,

$$\begin{aligned}
& \text{NEGap}(\mu^k, \nu^k) \\
& = \tilde{\mathcal{O}} \left(\varepsilon_k H + \frac{XAH}{k} + \frac{YBH}{k} \right. \\
& \quad + (X + Y) \left[(X + Y)^{\frac{1}{4}} H k^{(-\min\{1, \frac{3}{2} - \frac{\alpha_\varepsilon}{2}\} + \alpha_\eta + \alpha_\varepsilon)/2} + \sqrt{(XA + YB)} + H k^{\frac{-\alpha_\eta - \alpha_\varepsilon}{2}} \right. \\
& \quad \left. \left. + \sqrt{(X^2A + Y^2B)} k^{\frac{-\alpha_\eta + \alpha_\varepsilon}{2}} \right] \right) \\
& = \tilde{\mathcal{O}} \left(k^{-\frac{1}{6}} H + \frac{XAH}{k} + \frac{YBH}{k} + (X + Y) \left[(X + Y)^{\frac{1}{4}} H k^{-\frac{1}{6}} + \sqrt{(XA + YB)} H k^{-\frac{1}{3}} \right] \right)
\end{aligned}$$

$$\begin{aligned}
& + \sqrt{X^2A + Y^2B}k^{-\frac{1}{6}}) \\
& = \tilde{\mathcal{O}}\left((X + Y)\left[(X + Y)^{\frac{1}{4}}H + \sqrt{(X^2A + Y^2B)}\right]k^{-\frac{1}{6}}\right).
\end{aligned}$$

□

D PROOF OF LOWER BOUND OF LAST-ITERATE CONVERGENCE

Proof of Theorem 5.3. Let $\text{NEGap}_k := \text{NEGap}(\mu^k, \nu^k)$ with (μ^k, ν^k) as the policy profile generated by some algorithm Alg. Suppose that Alg leans the IIEFG with the last-iterate convergence rate of $\text{NEGap}_k = \Theta(f(X, A)k^{-\alpha})$ for some $\alpha \in (0, 1)$, where $f^{\text{Alg}}(X, A)$ denotes the polynomial dependence on X and A of NEGap_k .

Fix some $K \geq \max(XA, YB)$. Consider the regret defined as follows (Kozuno et al., 2021; Bai et al., 2022; Fiegel et al., 2023):

$$\text{Reg}_K(\text{Alg}) = \sup_{\mu \in \Pi_{\max}} \sum_{k=1}^K \langle \mu^k - \mu, \mathbf{G}\nu^k \rangle,$$

where $\{\nu^k\}_{k \in [K]}$ is potentially generated by an adversary. Then, one can deduce that

$$\begin{aligned}
\text{Reg}_K(\text{Alg}) &= \sup_{\mu \in \Pi_{\max}} \sum_{k=1}^K \langle \mu_k - \mu, \mathbf{G}\nu_k \rangle & (14) \\
&\leq \sum_{k=1}^K \sup_{\mu \in \Pi_{\max}} \langle \mu_k - \mu, \mathbf{G}\nu_k \rangle \\
&= \sum_{k=1}^K \sup_{\mu \in \Pi_{\max}} \mu_k^\top \mathbf{G}\nu_k - \mu^\top \mathbf{G}\nu_k \\
&\leq \sum_{k=1}^K \sup_{\mu \in \Pi_{\max}, \nu \in \Pi_{\min}} \mu_k^\top \mathbf{G}\nu - \mu^\top \mathbf{G}\nu_k \\
&= \sum_{k=1}^K \text{NEGap}_k \\
&= \Theta\left(f(X, A) \sum_{k=1}^K k^{-\alpha}\right) \\
&= \Theta(f(X, A)K^{1-\alpha}). & (15)
\end{aligned}$$

On the other hand, by Theorem 6 of Bai et al. (2022) (see also Theorem 3.1 of Fiegel et al. (2023)), we have

$$\text{Reg}_K(\text{Alg}) \geq \Omega(\sqrt{AXK}). \quad (16)$$

Combining Eq. (14) and Eq. (16), we have

$$\Omega(\sqrt{AXK}) \leq \Theta(f(X, A)K^{1-\alpha}).$$

We now further consider the following three cases:

- If $\alpha > \frac{1}{2}$, then $\sqrt{AX} \leq f(X, A)K^{\frac{1}{2}-\alpha}$. However, this does not hold for any f , when K is large enough;
- If $\alpha = \frac{1}{2}$, it must hold that $\sqrt{AX} \leq f(X, A)$;
- If $\alpha < \frac{1}{2}$, then $\sqrt{AX} \leq f(X, A)K^{\frac{1}{2}-\alpha}$. This holds for all f , including $f(X, A) = 1$ when K is large enough. In this case, the “minimal” f is $f(X, A) = 1$, implying that the minimal possible convergence rate of NEGap_k in this case is $\text{NEGap}_k = \Theta(k^{-\alpha})$.

Taking the above three cases into account, the minimal possible convergence rate is

$$\begin{aligned} & \min \left\{ \Theta \left(\sqrt{XA} k^{-\frac{1}{2}} \right), \Theta \left(k^{-\alpha} \right) \right\} \quad \left(\alpha > \frac{1}{2} \right) \\ & = \Theta \left(\sqrt{XA} k^{-\frac{1}{2}} \right). \end{aligned}$$

Analogously, we can prove that $\text{NEGap}_k \geq \Theta(\sqrt{YB} k^{-\frac{1}{2}})$. Therefore, we have

$$\text{NEGap}_k \geq \Theta \left(\left(\sqrt{XA} + \sqrt{YB} \right) k^{-\frac{1}{2}} \right).$$

The proof is concluded by noticing that the above holds for all algorithms. \square

E AUXILIARY LEMMAS

Lemma E.1 (Lemma 1 of Cai et al. (2023)). *Let $0 < h < 1, 0 \leq k \leq 2$, and let $t \geq \left(\frac{24}{1-h} \ln \frac{12}{1-h} \right)^{\frac{1}{1-h}}$. Then*

$$\sum_{i=1}^t \left(i^{-k} \prod_{j=i+1}^t (1 - j^{-h}) \right) \leq 9 \ln(t) t^{-k+h}.$$

Lemma E.2 (Lemma 2 of Cai et al. (2023)). *Let $0 < h < 1, 0 \leq k \leq 2$, and let $t \geq \left(\frac{24}{1-h} \ln \frac{12}{1-h} \right)^{\frac{1}{1-h}}$. Then*

$$\max_{1 \leq i \leq t} \left(i^{-k} \prod_{j=i+1}^t (1 - j^{-h}) \right) \leq 4t^{-k}.$$

Lemma E.3 (Lemma 20 of Bai et al. (2020)). *Let c_1, c_2, \dots, c_t be fixed positive numbers. Then with probability at least $1 - \delta$,*

$$\sum_{i=1}^t c_i \langle x_i, \ell_i - \widehat{\ell}_i \rangle = \mathcal{O} \left(A \sum_{i=1}^t \beta_i c_i + \sqrt{\ln(A/\delta) \sum_{i=1}^t c_i^2} \right).$$

F OPTIMIZATION PROBLEM IN EQ. (3)

Algorithm 3 Frank-Wolfe-type Algorithm for Solving Eq. (3) (max-player)

- 1: **Input:** Policy μ^k used in episode k , constrained policy space Π_{\max}^{k+1} , learning rate η^{k+1} , regularizer ψ , loss estimator $\widehat{\ell}^k$, number of iterations T .
 - 2: **Initialize:** $\mu^{(1)} = \mu^k$, $\phi(\mu) = \eta^{k+1} \langle \mu, \widehat{\ell}^k \rangle + D_\psi(\mu, \mu^k)$.
 - 3: **for** $t = 1, \dots, T$ **do**
 - 4: Compute $g^{(t)} = \nabla \phi(\mu^{(t)})$.
 - 5: Compute $\widehat{\mu}^{(t)} = \arg \min_{\mu \in \Pi_{\max}^{k+1}} \langle \mu, g^{(t)} \rangle$ by Algorithm 4.
 - 6: Let $\delta = \frac{2}{1+t}$.
 - 7: Update $\mu^{(t+1)} = (1 - \delta)\mu^{(t)} + \delta\widehat{\mu}^{(t)}$.
 - 8: **end for**
 - 9: **Return** $\mu^{(T)}$.
-

In this section, we provide Algorithm 3 and Algorithm 4, which compute an approximate solution to Eq. (3).

Algorithm 4 Computing Linear Minimizer in Algorithm 3 (max-player)

```

1836 1: Input:  $\Pi_{\max}^{k+1}, g^{(t)}$ .
1837
1838 2: Initialize:  $G^{(t)}(x_h, a_h) = 0, \mu(a_h|x_h) = 0, \forall (x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}, \forall h \in [H]$ .
1839
1840 3: for  $h = H, \dots, 1$  do
1841
1842 4:   for  $x_h \in \mathcal{X}_H$  do
1843
1844 5:     Compute
1845
1846     
$$G^{(t)}(x_h, a_h) = \sum_{x_{h+1} \in C(x_h, a_h), a_{h+1} \in \mathcal{A}} \mu(a_{h+1}|x_{h+1}) \left( g^{(t)}(x_{h+1}, a_{h+1}) + G^{(t)}(x_{h+1}, a_{h+1}) \right).$$

1847
1848 6:     Set  $\mu(a_h|x_h) = \frac{1}{A^{(k+1)}}, \forall a_h \in \mathcal{A}$ .
1849
1850 7:     Set  $\mu(a'_h|x_h) = 1 - \frac{A-1}{A^{(k+1)}}$ , where  $a'_h = \arg \min_{a \in \mathcal{A}} g^{(t)}(x_h, a) + G^{(t)}(x_h, a)$ .
1851
1852 8:   end for
1853
1854 9: end for
1855
1856 10: Return  $\mu$ .

```

Computation Complexity Suppose there are K episodes. Let $w = \max_{h \in [H], (x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}} |C(x_h, a_h)|$, where $C(x_h, a_h)$ is the set of immediate descendant infosets of (x_h, a_h) as defined in Section 3. Then the computation complexity of our Algorithm 2 and Algorithm 1 will be of $\mathcal{O}(wXA)$ and of $\mathcal{O}(wXA + K(XA + \text{Oracle}))$, where Oracle denotes the computation complexity of an oracle algorithm to solve our Eq. (3). If Algorithm 3 and Algorithm 4 are adopted to solve an approximate solution to Eq. (3), then Oracle will be of $\mathcal{O}(wXAT)$ where T is the number of iterations in Algorithm 3 and the total computation complexity of our Algorithm 1 will be of $\mathcal{O}(wXATK)$.

G EXPERIMENTS

In this section, we present the empirical evaluations of our Algorithm 1. Since we are not aware of any other algorithm that can also learn the (approximate) NE policy profile in IIEFGs with provable *last-iterate convergence* guarantees under bandit feedback, we compare our algorithm against previous algorithms that converge to the (approximate) NE policy profile in IIEFGs with only *average-iterate convergence* guarantees including IXOMD (Kozuno et al., 2021), BalancedOMD (Bai et al., 2022) and BalancedFTRL (Fiegel et al., 2023). Since these algorithms are only devised to obtain the average-iterate convergence for learning IIEFGs, the last-iterate convergence of these algorithms for learning IIEFGs is not theoretically guaranteed.

Environments We consider four standard IIEFG instances including Lewis Signaling, Kuhn Poker (Kuhn, 1950), Leduc Poker (Southey et al., 2012) and Liars Dice. All the implementation of these games are from the OpenSpiel library (Lanctot et al., 2019).

Implementation Details For our algorithm, to save the computation costs, instead of using our Algorithm 3 and Algorithm 4 to solve Eq. (3) in Algorithm 1, we use a lazy update of our Algorithm 1, where only the policy of the experienced trajectory of infoset action pairs $\{(x_h^k, a_h^k)\}_{h \in [H]}$ in each episode k are updated. For the remaining infoset action pairs that are not experienced by the max-player in episode k , the losses contributed by the entropy regularization (*i.e.*, the second term in our constructed entropy regularized loss estimator) of these infoset action pairs will be accumulated and will be used to update these infoset action pairs once they are experienced in some future episode, coming from the observation that the losses contributed by the entropy regularization are much smaller than the importance-weighted losses constructed using the rewards in the game (*i.e.*, the first term in our constructed entropy regularized loss estimator). In this way, the resulting computation complexity of our algorithm will only be of $\mathcal{O}(wXA + KXA)$ for running our algorithm in K episodes where $w = \max_{h \in [H], (x_h, a_h) \in \mathcal{X}_h \times \mathcal{A}} |C(x_h, a_h)|$ ($C(x_h, a_h)$ is the set of immediate descendant infosets of (x_h, a_h) as defined in Section 3). We adopt the implementation of all the baselines by Fiegel et al. (2023).² Besides, we consider a (logarithmic) grid search on the learning

²<https://github.com/anon17893/IIG-tree-adaptation>.

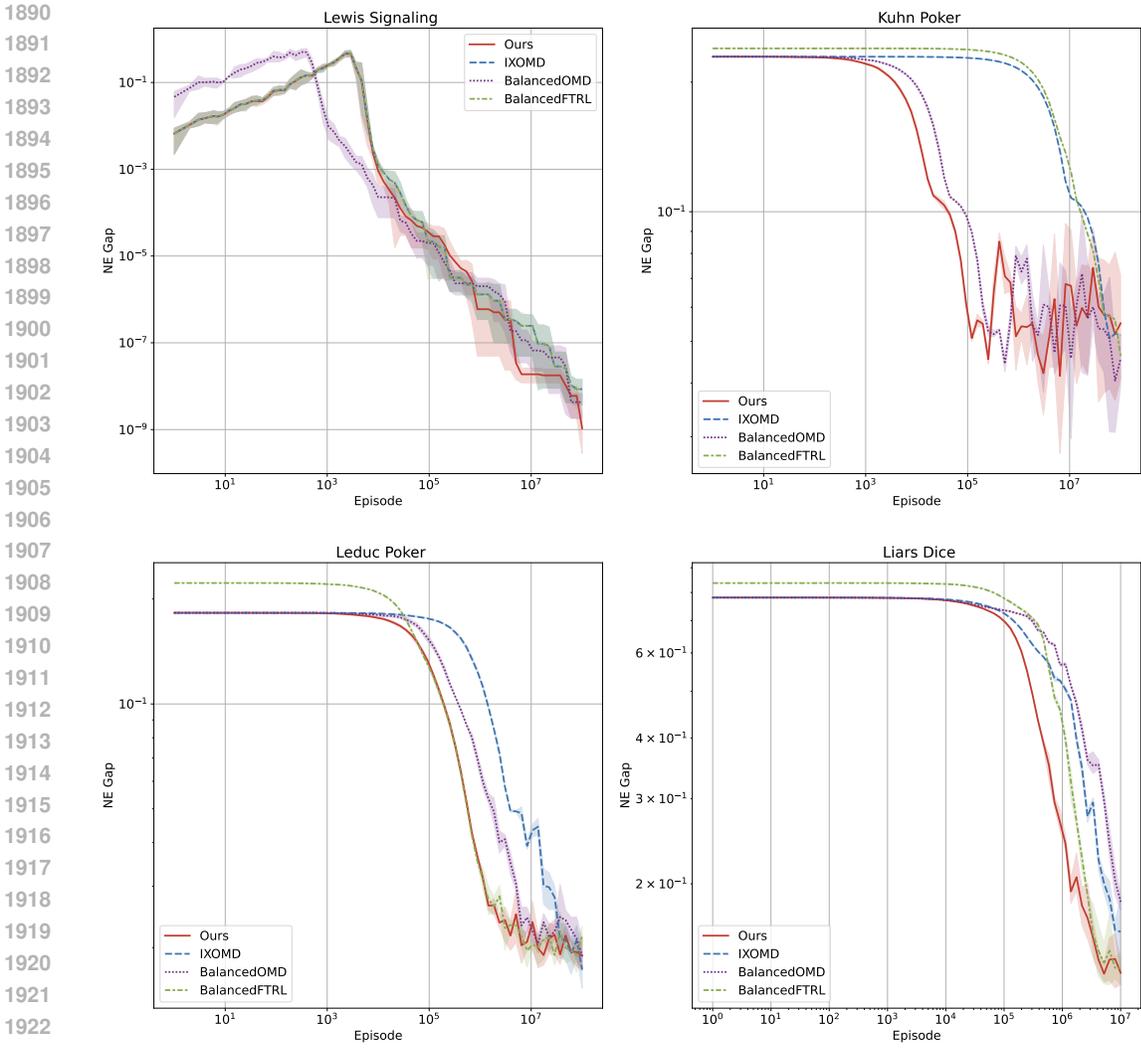


Figure 2: Experiment results of our Algorithm 1 against IXOMD (Kozuno et al., 2021), Balance-dOMD (Bai et al., 2022) and BalancedFTRL (Fiegel et al., 2023). The curves show the last-iterate convergence results of the NE gap defined in Eq. (2) against the number of episodes and are averaged over 5 different seeds.

rates for all the algorithms, following Fiegel et al. (2023). All the experiments are conducted on a server with an Intel Xeon Gold CPU and 251GiB system memory. The running of all the algorithms including our algorithm costs approximately 10 hours, 12 hours, 13 hours, and 16 hours on Lewis Signaling, Kuhn Poker, Leduc Poker, and Liars Dice, respectively.

Results The experimental results are shown in Figure 2. Our algorithm obtains the best or the competitive performance across all four IIEFG instances. In particular, our algorithm converges faster than all the baseline algorithms on Kuhn Poker and Liars Dice and also converges as fast as the empirically best baseline algorithm on Lewis Signaling and Leduc Poker. Though some baseline algorithms work relatively well on some game instances, we would like to note again that these algorithms are not theoretically guaranteed to converge to the NE policy profile with the last-iterate convergence. We speculate that this might also be the reason why some baseline algorithms perform relatively well in some instances but poorly in the remaining ones. For instance, the BalancedFTRL algorithm performs well on Leduc Poker while converging very slowly on Kuhn Poker. Analogously,

1944 BalancedOMD converges relatively well on Kuhn Poker and Leduc Poker but converges the most
1945 slowly on Liars Dice.

1946
1947 Moreover, in general, it appears that the advantage of our algorithm becomes more pronounced
1948 in IIEFG instances with larger infoset spaces \mathcal{X} (and action spaces \mathcal{A}) over previous algorithms.
1949 This observation aligns with the intuition that in such instances, the baseline algorithms, which
1950 solely have average-iterate convergence theoretical guarantees, face greater difficulty in achieving
1951 last-iterate convergence to the NE. This challenge may arise because these algorithms are more
1952 susceptible to getting stuck in suboptimal policy profiles, due to lack of the last-iterate convergence
1953 theoretical guarantees.

1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997